

CoDGraD: A Code-based Distributed Gradient Descent Scheme for Decentralized Convex Optimization

Elie Atallah, Nazanin Rahnavard, and Qiyu Sun

Department of Electrical and Computer Engineering, University of Central Florida, Orlando, FL

Department of Mathematics, University of Central Florida, Orlando, FL

Emails: elieatallah@knights.ucf.edu, nazanin@ece.ucf.edu, Qiyu.Sun@ucf.edu

Abstract—In this paper, we consider a large network containing many regions such that each region is equipped with a worker with some data processing and communication capability. For such a network, some workers may become stragglers due to the failure or heavy delay on computing or communicating. To resolve the above straggling problem, a coded scheme that introduces certain redundancy for every worker was recently proposed, and a gradient coding paradigm was developed to solve convex optimization problems when the network has a centralized fusion center. In this paper, we propose an iterative distributed algorithm, referred as Code-Based Distributed Gradient Descent algorithm (CoDGraD), to solve convex optimization problems over distributed networks. In each iteration of the proposed algorithm, an active worker shares the coded local gradient and approximated solution of the convex optimization problem with non-straggling workers at the adjacent regions only. In this paper, we also provide the consensus and convergence analysis for the CoDGraD algorithm and we demonstrate its performance via numerical simulations.

Index Terms—distributed optimization, gradient coding, consensus, distributed networks

I. INTRODUCTION

CONVEX optimization on a network of large size has played a significant role for solving various problems, such as big-data processing in machine learning, distributed parameter estimation in wireless sensor networks, distributed sampling and signal reconstruction, distributed design of filter banks, distributed spectrum sensing in cognitive radio networks, source localization in cellular networks [1, 2, 3, 4, 5, 6, 7]. The objective functions f in such optimization problems,

$$f(\mathbf{x}) = \sum_{l=1}^m f_l(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^N, \quad (\text{I.1})$$

are often the summation of some local objective functions $f_l, 1 \leq l \leq m$, related to a partition of the network. In this paper, we consider the scenario that each region of the partition is equipped with a worker that has some data processing and communication ability, while the network has a fusion center with limited capacity or it does not have a fusion center at all.

For a network with centralized data processing facility, the following optimization problem

$$\bar{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \sum_{l=1}^m f_l(\mathbf{x}), \quad (\text{I.2})$$

associated with the objective function $f = \sum_{l=1}^m f_l$ in (I.1) has been well studied, see [1, 4, 8, 9, 10, 11, 12, 13, 14, 16, 17, 18] and references therein for various algorithms implementable in a strong fusion center or in local workers distributed over the network. Denote the gradient of g on \mathbb{R}^N by ∇g . For a network equipped with one data processing unit only, a conventional approach to the optimization problem (I.2) is the *gradient descent* algorithm,

$$\mathbf{x}(k+1) = \mathbf{x}(k) - \frac{\alpha_k}{m} \sum_{l=1}^m \nabla f_l(\mathbf{x}(k)), \quad k \geq 0, \quad (\text{I.3})$$

where $\{\alpha_k\}_{k=0}^{\infty}$ is a positive sequence chosen appropriately [19, 20, 21, 22, 23, 24, 25, 26, 27]. Our illustrative examples of step sizes $\alpha_k, k \geq 0$, are

$$\alpha_k = (k+a)^{-\theta}, \quad k \geq 0, \quad (\text{I.4})$$

for some $\theta \in (1/2, 1)$ and $a \geq 1$.

Several distributed versions of the gradient descent algorithm (I.3) have been proposed, including the Adapt-Then-Combine algorithm (ATC) and the Combine-Then-Adapt algorithm (CTA) [13, 16], where implementation of the worker at agent l is given by

$$\begin{cases} \mathbf{y}_l(k) = \mathbf{x}_l(k) - \alpha_k \nabla f_l(\mathbf{x}_l(k)) \\ \mathbf{x}_l(k+1) = \sum_{j \in \mathcal{N}_l} w_{lj}(k) \mathbf{y}_j(k) \end{cases} \quad (\text{I.5})$$

and

$$\begin{cases} \mathbf{y}_l(k) = \sum_{j \in \mathcal{N}_l} w_{lj}(k) \mathbf{x}_j(k) \\ \mathbf{x}_l(k+1) = \mathbf{y}_l(k) - \alpha_k \nabla f_l(\mathbf{y}_l(k)) \end{cases} \quad (\text{I.6})$$

respectively, where \mathcal{N}_l contains all adjacent agents of the agent l for data sharing, and $(w_{lj})_{1 \leq l, j \leq m}$ is a consensus matrix. The above ATC and CTA algorithms may reach consensus over all nodes to the optimal solution $\bar{\mathbf{x}}$ in (I.2). They are essentially the same as the gradient descent algorithm (I.3) if the graph to describe topological structure of the network is complete and the consensus weights $w_{lj} = 1/m, 1 \leq l, j \leq m$. However, comparing with the gradient descent algorithm (I.3), in the implementation of the ATC and CTA algorithms,

we circumvent the expansive evaluation of gradient ∇f of the global objective function by evaluating gradients $\nabla f_l, 1 \leq l \leq m$, of local objective functions at each worker node and then communicating local gradients to the neighbors with nonzero consensus weights.

In applications such as distributed learning and optimization over the cloud, some workers in the network may become inactive due to the failure or heavy delay on computing or communicating [28, 29, 31]. To resolve the above problem, uncoded and coded local gradients have been proposed in [32, 33, 34] to recover the full gradient from local gradients on active nodes $\Delta \subset \{1, \dots, m\}$. Without loss of generality, we assume that $\Delta \subset \{1, \dots, n\}$, where $n \leq m$ is the number of active nodes. In this paper, we follow the coded scheme in [32] and consider the paradigm that for every active node i , the global objective function can be recovered from coded objective functions g_j on non-stragglers relative to node i , i.e.,

$$f(\mathbf{x}) = \sum_{j=1}^n a(i, j) g_j(\mathbf{x}) \quad (\text{I.7})$$

for some decoding matrix $\mathbf{A} = (a(i, j))_{1 \leq i, j \leq n}$. The above requirement is met if the coded scheme for non-stragglers is given by

$$g_i(\mathbf{x}) = \sum_{l=1}^m b(i, l) f_l(\mathbf{x}), \quad 1 \leq i \leq n, \quad (\text{I.8})$$

and the coding matrix $\mathbf{B} = (b(i, l))_{1 \leq i \leq n, 1 \leq l \leq m}$ satisfies

$$\mathbf{AB} = \mathbf{1}_{n \times m}, \quad (\text{I.9})$$

where $\mathbf{1}_{n \times m}$ is the $n \times m$ matrix with all entries taking value 1.

The coded scheme (I.7) and (I.8) has been used in [32] to solve the global optimization problem (I.2), where the worker at each region evaluates the coded local gradients and sends them to the worker at the master node; then the worker at the master node aggregates a weighted sum of coded local gradients to form the gradient of the global objective function, and applies a centralized gradient descent approach similar to (I.3) to update the approximation to the optimal solution $\bar{\mathbf{x}}$ in (I.2); and finally the worker at the master node sends the updated approximation to all local workers on the network for the next iteration. In this paper, based on the coded scheme (I.7) and (I.8), we propose a *distributed* algorithm to solve the convex optimization problem (I.2).

A. Objectives and Contributions

Since the focus of distributed optimization is overwhelmingly occupied with first-order methods, it is worth trying to enhance such methods rather than employing methods of another type. Our initial aim of this paper is to improve the performance of distributed gradient descent algorithm through utilizing coding. In particular, we focus on the leveraging of coding on the performance of the distributed optimization algorithm and how the convergence rate of these algorithms can be better enhanced through using an appropriate coding scheme based on the network topology.

While distributed optimization is implementable through distributing the local functions among the nodes to sum up to the global optimization function as in (I.1), we follow here an alternative route. We decompose the global function among the nodes in a coded manner and try to implement an algorithm which allows optimization under such scheme (I.7). To this end, we adapt the stochastic form of the decoding matrix \mathbf{A} in our proposed algorithm, carry the negativity in some of its coefficients to the gradient operation, and then utilize gradient descent and ascent steps, see (II.6) and (II.3).

In this paper, we do not impose any algebraic structure of the coding matrix \mathbf{B} and the decoding matrix \mathbf{A} , and thus our approach applies for a broader class of coding/decoding matrices as long as the allowed number of stragglers is fulfilled, except that the normalized decoding matrix $|\tilde{\mathbf{A}}|$ in (II.8) is assumed to have simple eigenvalue one, see Proposition II.2. Our work has the assumptions of strongly convex global function f , Lipschitz continuous (coded) gradients ∇g_i , and uniformly bounded (coded) gradients ∇f_i . We believe that if the coding/decoding matrices have additional structure properties [36], our algorithm could converge under weak requirements on the objective functions and coded gradients.

In this paper, we utilize exact gradient coding for fusion centralized networks and investigate the implications of such coding used in a multi-agent setting. This serves well to our initial aim of showing how enhancement to the convergence rate can be accomplished, see the consensus and convergence conclusions in Sections IV and V of our proposed algorithm. The paradigm of approximate gradient coding is discussed in [35]. It could be an interesting problem to extend our convergence conclusions in the above setting.

In this paper, we work on distributed optimization problem for multi-agent systems on fixed static network topology. Thus, the active nodes are the non-straggler nodes which are fixed throughout the proposed distributed algorithm and where coded local functions are used according to the coding scheme (I.7) instead of the brute decomposition as in (I.1). We postpone the consideration of time-varying straggler networks to a later endeavor.

The implementation of distributed algorithms on static/time-varying networks with stragglers is of importance. We wish that this work may serve as a starting point for a full-fledged investigation to distributed algorithms on static/time-varying networks with stragglers, with applications to the engineering field, such as federated decentralized learning and distributed machine learning.

B. Organization and notations

In Section II, we formulate our code-based distributed gradient descent algorithm (CoDGraD) to solve the optimization problem (I.2). In Section III we describe the coordinated distributed formation of the network and its data coding. Then in Sections IV and V, we consider the consensus and convergence properties of the proposed CoDGraD algorithm. Afterwards, we demonstrate the performance of CoDGraD through simulation in Section VI. We conclude the paper in Section VII.

In what follows, we use bold capital letters, bold lower case letters and lower case letters for matrices, vectors and scalar variables, respectively. Denote the positive part of a real number t by $t_+ = \max(t, 0)$, Denote the matrix of dimension $n \times m$ with all entries taking value one by $\mathbf{1}_{n \times m}$ (and $\mathbf{1}_n$ when $m = 1$) and the unit matrix of size $n \times n$ by \mathbf{I}_n , and the zero matrix of size $n \times m$ by $\mathbf{0}_{n \times m}$ respectively. Denote the transpose of an matrix \mathbf{A} by \mathbf{A}^T , and the transpose and standard ℓ^p -norm of a vector \mathbf{x} by \mathbf{x}^T and $\|\mathbf{x}\|_p, 1 \leq p \leq \infty$, respectively.

II. A CODE-BASED DISTRIBUTED GRADIENT DESCENT ALGORITHM

Let the coded objective functions $g_j, 1 \leq j \leq n$, and the decoding matrix $\mathbf{A} = (a(i, j))_{1 \leq i, j \leq n}$ be as in the coded scheme (I.7) and (I.8). Set decoding weights

$$w_i = \left(\sum_{j \in \Gamma_i} |a(i, j)| \right)^{-1}, \quad 1 \leq i \leq n, \quad (\text{II.1})$$

where

$$\Gamma_i = \{j, a(i, j) \neq 0\}. \quad (\text{II.2})$$

In [30, 32], all workers not in Γ_i are considered as ‘‘stragglers’’ or ‘‘non-neighboring workers’’ for an active node i , since coded information at those workers are not used to evaluate the gradient ∇f of the global objective function f at the node i . We remark that in this paper active nodes are those workers over the network that don’t witness delay or failure on computing or communicating. Being in a fixed static network implementation then all nodes in the graph of the network topology are considered as active nodes.

To solve the optimization problem (I.2), we propose an iterative distributed algorithm with the implementation of the worker at the region i given by

$$\begin{cases} \mathbf{v}_i(k) = \nabla g_i(\mathbf{x}_i(k)), \\ \mathbf{y}_i^+(k) = \mathbf{x}_i(k) - \alpha_k \mathbf{v}_i(k), \\ \mathbf{y}_i^-(k) = \mathbf{x}_i(k) + \alpha_k \mathbf{v}_i(k), \\ \mathbf{x}_i(k+1) = \sum_{j \in \Gamma_i} w_i \{ (a(i, j))_+ \mathbf{y}_j^+(k) \\ + (-a(i, j))_+ \mathbf{y}_j^-(k) \}, \quad k \geq 0, \end{cases} \quad (\text{II.3})$$

where initials $\mathbf{x}_i(0)$ are chosen randomly or set zero initially, and step sizes $\alpha_k, k \geq 0$, [14, 20, 23] are so chosen that

$$\sum_{k=0}^{\infty} \alpha_k = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty. \quad (\text{II.4})$$

We call this algorithm as a code-based distributed gradient descent algorithm and use CoDGraD for abbreviation. The implementation of the CoDGraD is described in Algorithm 1.

For our purpose, we consider the coded scheme (I.7) and (I.8) that matches our network topology. Here the network topology is described by an undirected graph $\mathcal{G} = (V, E)$, where the worker in each region is represented by a vertex in V and each edge $(l, l') \in E$ means that workers in the regions l and $l' \in V$ have direct communication for data sharing. Therefore the topological matching property of the coded scheme (I.7) and (I.8) is satisfied if the decoding weight $a(i, j)$

Algorithm 1 The CoDGraD Algorithm

Input: Set tolerance value ϵ for halting the algorithm and the number of iteration $k = 0$. Initialize an estimate $\mathbf{x}_i(0)$ of the optimal solution $\bar{\mathbf{x}}$ and approximation error $e_i(0) = \epsilon$ at the worker i .

- 1: **while** $e_i(k) \geq \epsilon$ **do** {Halting is done at each node independently with no coordination}
- 2: At the worker i , use (II.3) to update the estimate $\mathbf{x}_i(k)$.
- 3: Find error $e_i(k+1) = \|\mathbf{x}_i(k+1) - \mathbf{x}_i(k)\|$ for all $1 \leq i \leq n$
- 4: $k = k + 1$
- 5: **end while**

Output: $\mathbf{x}_i(k)$ and k .

takes zero value whenever there is no direct communication between active workers in the region i and j , i.e.,

$$a(i, j) = 0 \quad \text{if} \quad (i, j) \notin E.$$

This means that any non-straggling node of any active node i which is contributing in the decoding process must be a neighbor of the active node i in the whole network, i.e., $\Gamma_i \subset \mathcal{N}_i$, where Γ_i is given in (II.2). However, the converse is not necessarily true as a non-straggler node may be used in the coding procedure and does not contribute to the active node i .

For the above scenario of the coded scheme (I.7) and (I.8), the active worker at each region first evaluates the coded local gradient, then it updates its estimate through gradient descent/ascent step, next it shares the updated estimate approximations with neighboring active workers at the adjacent regions, henceforth updating its estimate towards the optimal solution $\bar{\mathbf{x}}$, Hence the CoDGraD algorithm is implementable in the network that does not have a fusion center at all.

For the decoding matrix \mathbf{A} in (I.7), we define its normalized decoding matrix $\mathbf{A}_{\text{sde}} = (\tilde{a}(i, j))_{1 \leq i, j \leq n}$ of size $n \times n$ by

$$\tilde{a}(i, j) = w_i a(i, j), \quad 1 \leq i, j \leq n, \quad (\text{II.5})$$

and its row stochastic decoding matrix \mathbf{A}_{sde} of size $2n \times 2n$ by

$$\mathbf{A}_{\text{sde}} = \begin{pmatrix} \tilde{\mathbf{A}}_+ & \tilde{\mathbf{A}}_- \\ \tilde{\mathbf{A}}_+ & \tilde{\mathbf{A}}_- \end{pmatrix}, \quad (\text{II.6})$$

where $w_i, 1 \leq i \leq n$, are decoding weights in (II.1), and

$$\tilde{\mathbf{A}}_+ = ((\tilde{a}(i, j))_+)_{1 \leq i, j \leq n} \quad \text{and} \quad \tilde{\mathbf{A}}_- = ((-\tilde{a}(i, j))_+)_{1 \leq i, j \leq n}$$

are positive/negative parts of the normalized decoding matrix $\tilde{\mathbf{A}} = (\tilde{a}(i, j))_{1 \leq i, j \leq n}$ respectively.

In this paper, we consider the **consensus** and **convergence** properties of $\mathbf{x}_i(k), k \geq 0$, in the proposed CoDGraD algorithm (II.3) under the following assumptions: (i) The row stochastic matrix \mathbf{A}_{sde} in (II.6) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin; (ii) The global objective function f is a differentiable strongly convex; (iii) The (coded) local objective functions $g_i, 1 \leq i \leq n$, have bounded gradients; and (iv) the (coded) local objective functions $g_j, 1 \leq j \leq n$ are differentiable and have continuous gradients.

A. Conditions on the Network Topology and Coding Scheme

The code-decode scheme in (I.9) presented in this paper is essentially the same as the scheme in [32], where n, m, s denote the number of workers, samples and stragglers respectively. Let $\Gamma_i, 1 \leq i \leq n$, be given in (II.2) and denote their complements by $\Gamma_i^c, 1 \leq i \leq n$. Therefore we require that stragglers to a node i are contained in $\Gamma_i^c, 1 \leq i \leq n$.

In the following, we consider two scenarios of network topology in which the matrix \mathbf{A}_{sde} used in the CoDGraD algorithm (II.3) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin. The first scenario is of full cyclic assignment of the active workers after certain permutation, i.e., the decoding matrix $\mathbf{A} = (a(i, j))_{1 \leq i, j \leq n}$ with the first row having its first $n - s$ entries assigned as non-zero, and as we move down the rows, the positions of the $n - s$ non-zero entries shifting one step to the right, and cycle around until the last row. Mathematically, the decoding matrix \mathbf{A} satisfies the following conditions:

$$a(i, j) \neq 0 \quad (\text{II.7})$$

for all (i, j) satisfying either $i \leq j \leq \min(i + n - s, n)$ or $j + s \leq i$, cf., [32, eq.10]. For the above scenario, we have

Proposition II.1. *Let $1 \leq s \leq n - 1$ and the decoding matrix \mathbf{A} satisfy (II.7). Then the matrix \mathbf{A}_{sde} used in the CoDGraD algorithm (II.3) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin.*

Denote the normalized decoding matrix of the decoding matrix \mathbf{A} by

$$|\tilde{\mathbf{A}}| = (|\tilde{a}(i, j)|)_{1 \leq i, j \leq n}, \quad (\text{II.8})$$

where $\tilde{a}(i, j), 1 \leq i, j \leq n$, are given in (II.5). By (II.5), the normalized decoding matrix $|\tilde{\mathbf{A}}|$ has row stochastic property. To prove Proposition II.1, we need an equivalence between the eigenvalue properties for the row stochastic matrices \mathbf{A}_{sde} and $|\tilde{\mathbf{A}}|$.

Proposition II.2. *The algebraic multiplicities of nonzero eigenvalues of stochastic decoding matrices \mathbf{A}_{sde} and $|\tilde{\mathbf{A}}|$ are the same.*

The proof of the above proposition will be given in Appendix A. We assume that Proposition eigenvalueone.pr holds and we give the proof of Proposition II.1 below.

Proof of Proposition II.1. Set $\mathbf{B} = |\tilde{\mathbf{A}}|$, and write $B^k = (b_k(i, j))_{1 \leq i, j \leq n}$, $k \geq 1$, and also $\mathbf{B} = (b(i, j))_{1 \leq i, j \leq n}$ for $k = 1$. Then \mathbf{B} is a row stochastic matrix with nonzero diagonal entries. By Perron-Frobenius theorem and Proposition II.2, it suffices to prove that \mathbf{B} is irreducible, i.e., for any $1 \leq i, j \leq n$, there exists $k \geq 1$ such that

$$b_k(i, j) \neq 0. \quad (\text{II.9})$$

By the assumption on the decode matrix \mathbf{A} , we have

$$b(i, j) \neq 0 \text{ if } i \leq j \leq i + n - s \text{ and if } j + s \leq i. \quad (\text{II.10})$$

Observe that for any $1 \leq i, j \leq n$, we have

$$\begin{aligned} b_{k+1}(i, j) &= \sum_{l=1}^n b(i, l)b_k(l, j) \geq \sum_{l=i}^{\min(i+n-s, n)} b(i, l)b_k(l, j) \\ &\quad + \sum_{l=1}^{i-s} b(i, l)b_k(l, j), \quad k \geq 1. \end{aligned} \quad (\text{II.11})$$

By (II.10) and (II.11), we can prove by induction on $k \geq 1$ that $b_k(i, j) \neq 0$ for all (i, j) satisfying either $i \leq j \leq \min(i + k(n - s), n)$ or $j + ks \leq i$. This proves (II.9) for all $k \geq n/s$ and completes the proof. \square

Let $\mathcal{G}_{\mathbf{A}} = (V, E_{\mathbf{A}})$ be the graph to describe the network topology in which there is an edge $(l, l') \in E_{\mathbf{A}}$ if and only if $a(i, j) \neq 0$. Define the minimum out/in degree of the network graph $\mathcal{G}_{\mathbf{A}}$ by $\delta_{\text{out}}(\mathcal{G}_{\mathbf{A}}) = \min_{1 \leq i \leq n} \#\{j, a(i, j) \neq 0\}$ and $\delta_{\text{in}}(\mathcal{G}_{\mathbf{A}}) = \min_{1 \leq i \leq n} \#\{j, a(j, i) \neq 0\}$ respectively. Next we consider the scenario of the network topology that

$$\delta(\mathcal{G}_{\mathbf{A}}) := \min(\delta_{\text{out}}(\mathcal{G}_{\mathbf{A}}), \delta_{\text{in}}(\mathcal{G}_{\mathbf{A}})) > n/2. \quad (\text{II.12})$$

In the above scenario, for each active node there are at least $\delta(\mathcal{G})$ non-straggler to receive information from and to send information to.

Proposition II.3. *If (II.12) holds, then the matrix \mathbf{A}_{sde} used in the CoDGraD algorithm (II.3) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin.*

Proof. Following the argument used in the proof of Proposition II.1, it suffices to prove

$$b_2(i, j) \neq 0 \text{ for all } 1 \leq i, j \leq n. \quad (\text{II.13})$$

Set $C_1 = \{l, b(i, l) \neq 0\}$ and $C_2 = \{l, b(l, j) \neq 0\}$. By the assumption on the decoding matrix \mathbf{A} , they contain at least $\delta(\mathcal{G}) > n/2$ elements contained in C_1 and C_2 respectively, and hence there exists $l_0 \in C_1 \cap C_2$. Hence

$$b_2(i, j) = \sum_{l=1}^n b(i, l)b_k(l, j) \geq b(i, l_0)b(l_0, j) > 0.$$

This proves (II.13) and completes the proof. \square

III. COORDINATED DISTRIBUTED CODING

A. Network Formation

1) *Network Detection:* A node gets activated and decides to form a network coded CoDGraD implementation. We signify this node as the *coordinator node*. The coordinator node sends a message containing a label identifying the CoDGraD implementation, the transmitting node (i.e., which is itself), the accumulated path of the message (i.e., currently the node itself), and its public key. When a neighboring node receives that message it transmits a message containing the label of the received message, the current transmitting node (i.e., which is the current node), the accumulated path of message (i.e., appending previous message path with the current transmitting node) and an encrypted message block. This neighboring node also sends in its encrypted message block part; its symmetric key with its identifier both encrypted with the coordinator

node public key. This process of transmitting a message and receiving a message then retransmitting continues henceforth until one of the two halting criteria is met.

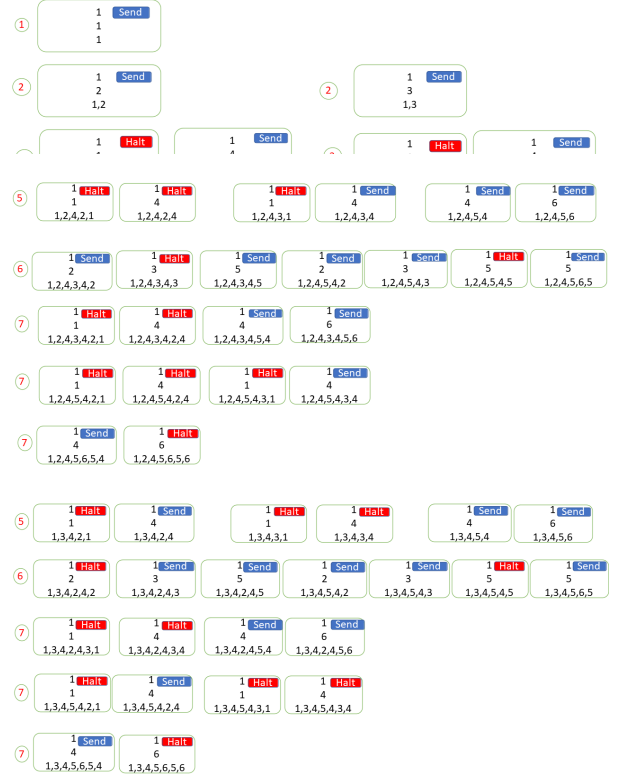
- 1) If the to-be-transmitted message at a node contains in its accumulated path an edge which is traversed twice in the same direction, then this message will not be retransmitted.
- 2) A desired time that takes into the consideration the size of the network to be designed and its desired performance is set. After that time, the coordinator node would have received messages related to its CoDGrAD implementation, that contain the label, the transmitting nodes and the accumulated paths.

Every time criterion 2 is met, the coordinator node forms from the accumulated paths the adjacency matrix of the interacting nodes. And according to the designed coding scheme and the allowed number of stragglers, the coordinator node decides on the nodes that need to join the network thus satisfying the stragglers and data partitions' redundancy thresholds. Thus it forms the network adjacency matrix and its related coding schemes encoding and decoding matrices \mathbf{A}_{sde} and \mathbf{B} , i.e., the decoding matrix used for the weighing matrix and the gradient coding matrix used for information privacy. The coordinator node will also decrypts all encrypted data in the encrypted message block of the message, thus retrieving the symmetric keys of each node of the network that were encrypted with its public key.

2) *Network Forming and Coding:* Following the protocol described in the preceding subsection, the coordinator node sends a new message containing the label of the network, an encrypted information containing the weights relative to the node that will join the network (i.e., \mathbf{A}_{sde} row), and also an encrypted information containing the coefficients of the coded gradients of the node to join the network (i.e., \mathbf{B} row for computing ∇g_i). Subsequently, each node will recover this encrypted information through privacy symmetric keys between the coordinator node and itself.

More precisely, the coordinator node sends in the encrypted message block of the message the row of \mathbf{A}_{sde} identified with node i , for all nodes $i \in \mathcal{G}$, each encrypted with the respective symmetric key of node i that was decrypted in the previous step. It also sends the weight of ∇f_i (i.e., \mathbf{B}_{ji}) that will be used in coding node i primary gradient in its neighbor j and the symmetric key of node j , both encrypted with the symmetric key of node i which was decrypted in the previous step. It performs this operation for all nodes i and their respective neighbors $j \in \mathcal{N}_i$, accordingly. In this message there are also the transmitting node and accumulated path information as before.

This message might also contain the adjacency matrix of the designed network and the adjacency matrix of the larger network that also contains nodes that are not allowed to join the network at this time due to not meeting the straggler threshold. All this adjacency information is also encrypted in such a way that only nodes of the designed network can access this information and only the information related to their neighborhoods.

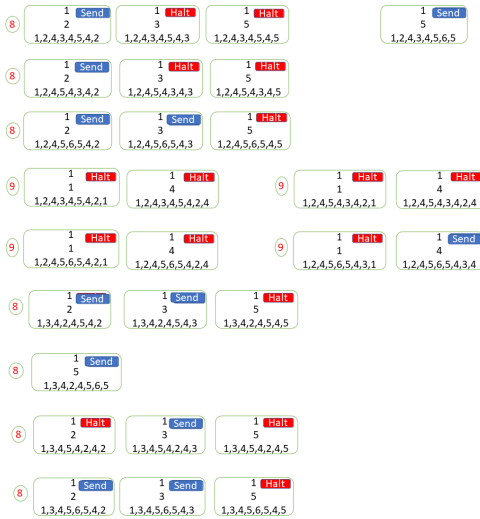


Meanwhile, the coordinator node and all other nodes of the network will separate their raw data into two parts. The primary information containing the initial raw data at the node and the secondary information containing the coded information formed from neighboring nodes due to CoDGrAD coding. When the coordinator node sends the previously described message it also sends its coded primary information to its neighboring vicinity. Thus, it sends its partition of raw data weighted with the weights \mathbf{B}_{j1} for all $j \in \mathcal{N}_1$, each encrypted with the symmetric key of node j decrypted previously by the coordinator node. And obviously, the coordinator node (i.e., node 1) will have access of the row of \mathbf{A}_{sde} found in the previous step. As for the other nodes that are allowed by the coordinator node to join the network, when they receive this message they send the same message to their neighboring nodes with the new transmitting node information and the accumulated paths. They will also decrypt using their symmetric keys all the information related to them in the encrypted message block. Thus, each node i will be able to decrypt the row of \mathbf{A}_{sde} identified with it, and the weight \mathbf{B}_{ji} that will be used in coding the primary data partitions (i.e., the coded gradients ∇g_j) at node j , together with the symmetric key of its neighbor j , for all of its neighbors $j \in \mathcal{N}_i$. At the first reception of this message the receiving node i also transmits in conjunction to the above described message (i.e., probably in a different channel) its primary raw information weighted by weights \mathbf{B}_{ji} for all $j \in \mathcal{N}_i$, each encrypted with the symmetric key of node $j \in \mathcal{N}_i$ decrypted previously. In the same way when a node j receives the encrypted primary raw data containing the data partitions used for evaluating the coded gradients ∇g_j from each neighboring node $i \in \mathcal{N}_j$, it

decrypts with its symmetric key the coded data part of each of its neighbors i and stores it in its secondary information.

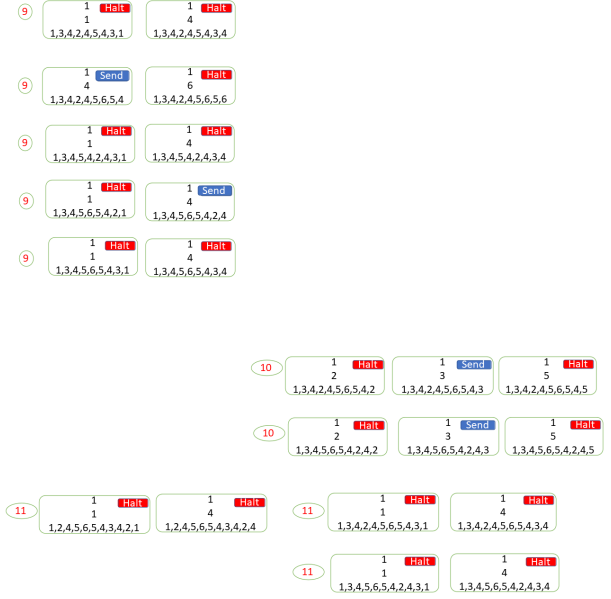
Thus, as it was mentioned before, by now each node will be able, according to the received encrypted coded primary information from neighboring nodes and the encrypted coefficients of the coding scheme, to form its weights matrix row part (i.e., \mathbf{A}_{sde} row) and the coded secondary information needed in gradient coding (i.e., computes ∇g_i).

The message at a node is not transmitted again, if as before, it contains in its accumulated path an edge crossed in the same direction twice. Meanwhile, the primary information is transmitted only once or according to a designed protocol, one example would be, when all its neighbors send a message signifying that they received that information. This process continues until the coordinator node receives all messages containing the sent encrypted information with the accumulated paths and conceives that all nodes in the desired network have coded their data and formed their weighting coefficients. Then it decides to implement the CoDGraD algorithm.



3) *Adding nodes to already formed network:* While the CoDGraD algorithm is in process, when a node detects a message from a new active node it sends an encrypted message containing the updated neighborhood of this node with the new out of network node. This new node detection can be perceived directly through the new node sensing a source different from its usual neighbors or can also be recognized if the node detects a new neighbor not in the neighborhood adjacency matrix row which was communicated to it by the coordinator node in the previous step. Note that the latter policy is usually used if the update of the new adjacency matrix is accomplished only at the coordinator node and the first if any node in the network can perform such adjacency matrix upgrade. This follows for all nodes that detect new nodes. And when these messages are received by other nodes in the network they send this encrypted information to neighboring nodes only when they receive these messages for the first time. When the coordinator node receives these updated neighborhoods in an encrypted manner it deciphers the adjacency matrix of the resulting new network and if the new nodes meet the straggler threshold they are added to the network and a new

updated network is formed. Henceforth, the coordinator node, as in stage 2, sends the encrypted messages of the adjacency matrices and coding schemes and allow the formation of a new network with new coding containing the allowed new nodes and thus implementing the CoDGraD on this new network.



It is worth mentioning that in our analysis we focused on static topology with fixed straggler nodes and thus fixed weighing and gradient coding matrices. Although we encrypted the information which allows upgrading to dynamic networks with privacy due to encrypting, however, we restrict our analysis to one weighing matrix and one coding scheme. Dynamic networks with the same nodes require protocol that uses the encoding matrix \mathbf{A} of all possible $s+1$ -combinations, where s is the maximum number of allowed stragglers, and the process of adding new nodes needs the use of new weighing matrices (i.e., a new coding scheme), so we will leave their analysis to a future work. Although we can also use one coding scheme and one weighing matrix for a dynamic network, the one corresponding to the $n-s$ non-straggler combinations but its performance will be considerably degraded.

Meanwhile, it is also worth noting that we could have allowed any node to upgrade the coding scheme whence it receives information about a neighborhood of a new node with the allowed straggler threshold. But we have restricted that to one node and specifically the coordinator node in order to preserve encrypted privacy keys designed due to that node. And thus not allowing the first approach since then we need to disclose the whole adjacency matrix to all nodes. By performing that we would be unable to preserve privacy through allowing only neighborhood information to be disclosed to each node of the network while keeping the whole information to the coordinator node, the CoDGraD implementer. However, we can allow other nodes to send this information if they keep it encrypted through encrypting keys between them and the rest of the nodes.

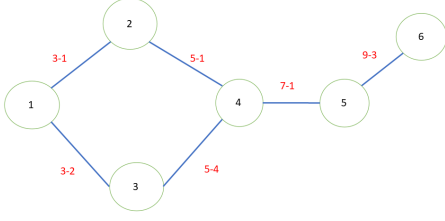


Fig. 1. Example of a 6-node network where coordinated distributed coded network forming is performed. The coefficients on the edges shows the corresponding message received by the coordinator node 1 that allows it to infer that edge connection.

Remark III.1. We have provided in this section an example of the messaging protocol used in network forming on the 6-node network described in Fig 1.

IV. CONSENSUS PROPERTY OF THE CODE-BASED DISTRIBUTED GRADIENT DESCENT ALGORITHM

In this section, we consider the consensus property of $\mathbf{x}_i(k)$, $k \geq 1$, in the CoDGrAD algorithm (II.3).

Theorem IV.1. Let $\mathbf{x}_i(k)$, $k \geq 1$, be in the CoDGrAD algorithm (II.3). If the row stochastic matrix \mathbf{A}_{sde} in (II.6) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin, the sequence $\{\alpha_k\}_{k=0}^{\infty}$ satisfies (II.4), and the local objective functions g_i , $1 \leq i \leq n$, have bounded gradients, i.e., there exists a positive constant M such that

$$\|\nabla g_i(\mathbf{x})\|_2 \leq M, \quad \mathbf{x} \in \mathbb{R}^N, \quad (\text{IV.1})$$

then

$$\lim_{k \rightarrow \infty} (\mathbf{x}_i(k) - \mathbf{x}_j(k)) = 0 \quad (\text{IV.2})$$

and

$$\lim_{k \rightarrow \infty} (f(\mathbf{x}_i(k)) - f(\mathbf{x}_j(k))) = 0, \quad 1 \leq i, j \leq n. \quad (\text{IV.3})$$

Let the row stochastic matrix \mathbf{A}_{sde} in (II.6) have simple eigenvalue one and $\lambda_m(\mathbf{A}_{\text{sde}})$, $1 \leq m \leq 2n$, be its eigenvalues listed in the order that

$$1 = \lambda_1(\mathbf{A}_{\text{sde}}) > |\lambda_2(\mathbf{A}_{\text{sde}})| \geq \dots \geq |\lambda_{2n}(\mathbf{A}_{\text{sde}})|. \quad (\text{IV.4})$$

Write $\mathbf{A}_{\text{sde}} = (q(i, j))_{1 \leq i, j \leq 2n}$ and for $1 \leq i \leq n$, set

$$\mathbf{z}_i(k) = \mathbf{z}_{i+n}(k) = \mathbf{x}_i(k), \quad k \geq 0. \quad (\text{IV.5})$$

Then the CoDGrAD algorithm (II.3) can be rewritten as

$$\mathbf{z}_i(k+1) = \sum_{j=1}^{2n} q(i, j) (\mathbf{z}_j(k) - \alpha_k \mathbf{h}_j(k)), \quad 1 \leq i \leq 2n, \quad (\text{IV.6})$$

where

$$\mathbf{h}_i(k) = \begin{cases} \nabla g_i(\mathbf{z}_i(k)) & \text{if } 1 \leq i \leq n \\ -\nabla g_{i-n}(\mathbf{z}_i(k)) & \text{if } n+1 \leq i \leq 2n. \end{cases} \quad (\text{IV.7})$$

Set

$$\mathbf{z}(k) := (\mathbf{z}_i(k))_{1 \leq i \leq 2n} \quad \text{and} \quad \mathbf{h}(k) := (\mathbf{h}_i(k))_{1 \leq i \leq 2n} \quad (\text{IV.8})$$

with vectors $\mathbf{z}_i(k)$ and $\mathbf{h}_i(k) \in \mathbb{R}^N$, $1 \leq i \leq 2n$, as their i -th entries respectively. Then the iterative algorithm (IV.6) can be reformulated in a matrix form:

$$\mathbf{z}(k+1) = \mathbf{A}_{\text{sde}} \mathbf{z}(k) - \alpha_k \mathbf{A}_{\text{sde}} \mathbf{h}(k), \quad k \geq 0. \quad (\text{IV.9})$$

Define

$$\tilde{\mathbf{z}}(k) = \mathbf{z}(k) - \mathbf{P} \mathbf{z}(k), \quad k \geq 0, \quad (\text{IV.10})$$

where

$$\mathbf{P} = \mathbf{1}_{2n} (\mathbf{a}_{\text{sde}})^T \quad (\text{IV.11})$$

and \mathbf{a}_{sde} is the stationary probability vector invariant under the row stochastic matrix \mathbf{A}_{sde} , i.e., the left eigenvector of \mathbf{A}_{sde} associated with eigenvalue one that satisfies

$$\mathbf{a}_{\text{sde}}^T \mathbf{A}_{\text{sde}} = \mathbf{a}_{\text{sde}}^T \quad \text{and} \quad \mathbf{a}_{\text{sde}}^T \mathbf{1}_{2n} = 1. \quad (\text{IV.12})$$

Then the consensus property (IV.2) reduces to establishing

$$\lim_{k \rightarrow \infty} \|\tilde{\mathbf{z}}(k)\|_{2, \infty} = 0, \quad (\text{IV.13})$$

where $\|\mathbf{z}\|_{2, \infty} = \sup_{1 \leq i \leq 2n} \|\mathbf{z}_i\|_2$ for a vector $\mathbf{z} = (\mathbf{z}_i)_{1 \leq i \leq 2n}$ with entries $\mathbf{z}_i \in \mathbb{R}^N$, $1 \leq i \leq 2n$.

By (IV.11) and (IV.12), we have

$$\mathbf{P} \mathbf{A}_{\text{sde}} = \mathbf{A}_{\text{sde}} \mathbf{P} = \mathbf{P} \quad \text{and} \quad \mathbf{P}^2 = \mathbf{P}. \quad (\text{IV.14})$$

This together with (IV.9) implies that

$$\tilde{\mathbf{z}}(k+1) = (\mathbf{A}_{\text{sde}} - \mathbf{P}) \tilde{\mathbf{z}}(k) - \alpha_k (\mathbf{A}_{\text{sde}} - \mathbf{P}) \mathbf{h}(k). \quad (\text{IV.15})$$

Applying (IV.15) repeatedly yields

$$\tilde{\mathbf{z}}(k) = (\mathbf{A}_{\text{sde}} - \mathbf{P})^k \tilde{\mathbf{z}}(0) - \sum_{l=0}^{k-1} \alpha_l (\mathbf{A}_{\text{sde}} - \mathbf{P})^{k-l} \mathbf{h}(l), \quad k \geq 1. \quad (\text{IV.16})$$

Therefore, we have the following estimate for $\|\tilde{\mathbf{z}}(k)\|_{2, \infty}$, $k \geq 1$ in Proposition 1, see Appendix B for a detailed proof.

Proposition IV.2. Let \mathbf{A}_{sde} , $\lambda_2(\mathbf{A}_{\text{sde}})$ and \mathbf{P} be as in (II.6), (IV.4) and (IV.11) respectively. Assume that the row stochastic matrix \mathbf{A}_{sde} has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin, and that the local objective functions g_i , $1 \leq i \leq n$, satisfy (IV.1). Then there exists a positive constant C_1 such that

$$\|\tilde{\mathbf{z}}(k)\|_{2, \infty} \leq C_1 M \sum_{l=0}^{k-1} \left(\frac{1 + |\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l + C_1 \left(\frac{1 + |\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^k \|\tilde{\mathbf{z}}(0)\|_{2, \infty} \quad (\text{IV.17})$$

hold for all $k \geq 1$.

By (IV.11), we can write

$$\mathbf{P} \mathbf{z}(k) = \bar{\mathbf{x}}(k) \mathbf{1}_{2n} \quad \text{for some } \bar{\mathbf{x}}(k) \in \mathbb{R}^N. \quad (\text{IV.18})$$

Observe that $\|\tilde{\mathbf{z}}(k)\|_{2, \infty} = \max_{1 \leq i \leq 2n} \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2$, $k \geq 1$. Then by Proposition IV.2 we obtain the following estimate about the consensus property of $\mathbf{x}_i(k)$ for different $1 \leq i \leq n$:

$$\begin{aligned} & \max_{1 \leq i \leq n} \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2 \\ & \leq C_1 \left(\frac{1 + |\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^k \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2 \\ & \quad + C_1 M \sum_{l=0}^{k-1} \left(\frac{1 + |\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l, \quad k \geq 1. \end{aligned} \quad (\text{IV.19})$$

Set $\gamma = \frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \in (1/2, 1)$. Applying (IV.19) to our illustrative example (I.4) of step sizes $\alpha_k = (k+1)^{-\theta}$, $k \geq 0$ for some $1/2 < \theta \leq 1$, we can find a positive constant C such that

$$\begin{aligned} & \max_{1 \leq i \leq n} \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2 \\ & \leq C_1 \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2 \gamma^k \\ & \quad + C_1 M \sum_{l=0}^{k-1} \gamma^{k-l} (k-l+a)^\theta (k+a)^{-\theta} \\ & \leq C(k+a)^{-\theta}, \quad k \geq 0, \end{aligned} \quad (\text{IV.20})$$

where the first inequality follows from (IV.19) and the observation that $a \geq 1$ and $(k+a)^\theta \leq (l+a)^\theta (k-l+a)^\theta$, $0 \leq l \leq k$, and the second estimate holds by the boundedness of the sequence $\gamma^k (k+a)^\theta$, $k \geq 0$, and the convergence of the series $\sum_{m=0}^{\infty} \gamma^m (m+a)^\theta$.

We finish this section with the proof of Theorem IV.1.

Proof of Theorem IV.1. By (IV.4) and (IV.17), we have

$$\begin{aligned} & \left(\sum_{k=0}^{\infty} \max_{1 \leq i \leq n} \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2^2 \right)^{1/2} \\ & \leq C_1 M \left(\sum_{k=0}^{\infty} \left(\sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l \right)^2 \right)^{1/2} \\ & \quad + C_1 \left(\sum_{k=0}^{\infty} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{2k} \right)^{1/2} \\ & \quad \times \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2 \\ & \leq \frac{2C_1 M}{1-|\lambda_2(\mathbf{A}_{\text{sde}})|} \left(\sum_{k=0}^{\infty} \alpha_k^2 \right)^{1/2} \\ & \quad + \frac{2C_1}{\sqrt{1-|\lambda_2(\mathbf{A}_{\text{sde}})|}} \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2, \end{aligned} \quad (\text{IV.21})$$

where the first inequality is obtained from (IV.17) and the triangle inequality for the space of square-summable sequence,

$$\left\{ \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^k \right\}_{k=0}^{\infty} \quad \text{and} \quad \left\{ \sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l \right\}_{k=0}^{\infty},$$

and the second estimate follows from

$$\begin{aligned} & \sum_{k=0}^{\infty} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{2k} = \left(1 - \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^2 \right)^{-1} \\ & \leq \left(1 - \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right) \right)^{-1} \left(1 + \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^2 \right)^{-1} \\ & \leq \frac{4}{1-|\lambda_2(\mathbf{A}_{\text{sde}})|}, \end{aligned}$$

and

$$\begin{aligned} & \sum_{k=0}^{\infty} \left(\sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l \right)^2 \\ & \leq \sum_{k=0}^{\infty} \left(\sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l^2 \right) \\ & \quad \times \left(\sum_{l'=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l'} \right) \\ & \leq \frac{2}{1-|\lambda_2(\mathbf{A}_{\text{sde}})|} \sum_{k=0}^{\infty} \left(\sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l^2 \right) \\ & \leq \frac{2}{1-|\lambda_2(\mathbf{A}_{\text{sde}})|} \sum_{l=0}^{\infty} \sum_{k=l+1}^{\infty} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l^2 \\ & \leq \left(\frac{2}{1-|\lambda_2(\mathbf{A}_{\text{sde}})|} \right)^2 \sum_{l=0}^{\infty} \alpha_l^2. \end{aligned}$$

Therefore the limit in (IV.2) follows as the sequence $\{\alpha_k\}_{k=0}^{\infty}$ is square summable by (II.4).

From (I.7) and (II.1) it follows that

$$\|\nabla f(\mathbf{x})\|_2 \leq W^{-1} \sup_{1 \leq j \leq n} \|\nabla g_j(\mathbf{x})\|_2 \leq MW^{-1},$$

where $W = \max_{1 \leq i \leq n} w_i$. This together with Proposition IV.2 implies that

$$\begin{aligned} & |f(\mathbf{x}_i(k)) - f(\bar{\mathbf{x}}(k))| \\ & \leq \sup_{0 \leq t \leq 1} \|\nabla f(t\mathbf{x}_i(k) + (1-t)\bar{\mathbf{x}}(k))\|_2 \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2 \\ & \leq C_1 M^2 W^{-1} \sum_{l=0}^{k-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^{k-l} \alpha_l \\ & \quad + C_1 M W^{-1} \left(\frac{1+|\lambda_2(\mathbf{A}_{\text{sde}})|}{2} \right)^k \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2 \\ & \rightarrow 0 \quad \text{as } k \rightarrow \infty \end{aligned} \quad (\text{IV.22})$$

for all $1 \leq i \leq n$. Then the convergence (IV.3) of the difference $f(\mathbf{x}_i(k))$, $k \geq 1$, between different $1 \leq i \leq n$ follows. \square

V. CONVERGENCE PROPERTY OF THE CODE-BASED DISTRIBUTED GRADIENT DESCENT ALGORITHM

In this section, we consider the convergence of $\mathbf{x}_i(k)$, $k \geq 1$, in the CoDGraD algorithm (II.3),

$$\lim_{k \rightarrow \infty} \mathbf{x}_i(k) = \bar{\mathbf{x}}, \quad 1 \leq i \leq n, \quad (\text{V.1})$$

where $\bar{\mathbf{x}}$ is the solution of the optimization problem (I.2). By (IV.2), (IV.18) and (IV.21), it suffices to show that $\bar{\mathbf{x}}(k)$, $k \geq 1$, converges to $\bar{\mathbf{x}}$.

Theorem V.1. *Assume that the row stochastic matrix \mathbf{A}_{sde} in (II.6) has simple eigenvalue one and all other eigenvalues contained in the open unit complex disk centered at the origin, the sequence $\{\alpha_k\}_{k=0}^{\infty}$ satisfies (II.4), the local objective functions g_j , $1 \leq j \leq n$, satisfy (IV.1) and*

$$\|\nabla g_i(x) - \nabla g_i(y)\|_2 \leq L \|x - y\|_2, \quad x, y \in \mathbb{R}^N, \quad 1 \leq i \leq n, \quad (\text{V.2})$$

for some positive constant L , and the global objective function f is strongly convex in the sense that

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle_N \geq A \|\mathbf{x} - \mathbf{y}\|_2^2 \quad (\text{V.3})$$

for all $\mathbf{x}, \mathbf{y} \in B(\bar{\mathbf{x}}, C_2)$, where $\langle \cdot, \cdot \rangle_N$ is the inner product on \mathbb{R}^N , $B(\bar{\mathbf{x}}, C_2)$ is the ball with center $\bar{\mathbf{x}}$ and radius C_2 , and

$$\begin{aligned} C_2 := & \exp\left(\sum_{j=0}^{\infty} \alpha_j^2\right) \left\{ \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 \right. \\ & + \frac{8C_1^2 L^2}{(1 - |\lambda_2(\mathbf{A}_{\text{sde}})|)^2} \max_{1 \leq i \leq n} \|\mathbf{x}_i(0) - \bar{\mathbf{x}}(0)\|_2^2 \\ & \left. + \frac{M^2(1 + 8C_1^2)}{(1 - |\lambda_2(\mathbf{A}_{\text{sde}})|)^2} \left(\sum_{j=0}^{\infty} \alpha_j^2\right)\right\}. \end{aligned}$$

Then $\bar{\mathbf{x}}(k), k \geq 1$, in (IV.18) converges to the solution $\bar{\mathbf{x}}$ of the optimization problem (I.2),

$$\lim_{k \rightarrow \infty} \bar{\mathbf{x}}(k) = \bar{\mathbf{x}}. \quad (\text{V.4})$$

Combining Theorems IV.1 and V.1, we have the following result on the convergence of the proposed CoDGrad algorithm.

Theorem V.2. *Let the decoding matrix \mathbf{A}_{sde} and the objective function f be as in Theorem V.1, and step sizes $\{\alpha_k\}_{k=0}^{\infty}$ in the CoDGrad algorithm satisfy (II.4). Then the sequences $\{\mathbf{x}_i(k)\}_{k=1}^{\infty}, 1 \leq i \leq n$ in the CoDGrad algorithm converge to the optimal point $\bar{\mathbf{x}}$, i.e., $\lim_{k \rightarrow \infty} \mathbf{x}_i(k) = \bar{\mathbf{x}}, 1 \leq i \leq n$.*

To prove Theorem V.1, we need a technical lemma, which follows the probability property for the vector \mathbf{a}_{sde} in (IV.12). For the completeness of this paper, we include a detailed proof in Appendix C.

Lemma V.3. *Let \mathbf{a}_{sde} and $w_i, 1 \leq i \leq n$, be as in (IV.12) and (II.1) respectively. Set*

$$\mathbf{w} = (w_1, \dots, w_n, w_1, \dots, w_n)^T \quad \text{and} \quad \tilde{w} = \mathbf{a}_{\text{sde}}^T \mathbf{w}. \quad (\text{V.5})$$

Then

$$0 < \min_{1 \leq i \leq n} w_i \leq \tilde{w} \leq \max_{1 \leq i \leq n} w_i. \quad (\text{V.6})$$

By (IV.11) and (IV.18), we have

$$\bar{\mathbf{x}}(k) = \mathbf{a}_{\text{sde}}^T \mathbf{z}(k). \quad (\text{V.7})$$

This together with (IV.9) leads to the following iterative algorithm for $\bar{\mathbf{x}}(k), k \geq 0$:

$$\bar{\mathbf{x}}(k+1) = \bar{\mathbf{x}}(k) - \alpha_k \mathbf{a}_{\text{sde}}^T \mathbf{h}(k), \quad k \geq 0. \quad (\text{V.8})$$

For local objective functions $g_i, 1 \leq i \leq n$, with Lipschitz gradients (V.2), we observe that $\mathbf{a}_{\text{sde}}^T \mathbf{h}(k)$ is an inexact estimate of the scaled global gradient $\tilde{w} \nabla f(\bar{\mathbf{x}}(k))$, see Appendix D for a detailed proof.

Proposition V.4. *Let $\mathbf{A}_{\text{sde}}, \lambda_2(\mathbf{A}_{\text{sde}}), \mathbf{P}$ and \tilde{w} be as in (II.6), (IV.4), (IV.11) and (V.5) respectively. Assume that the row stochastic matrix \mathbf{A}_{sde} has simple eigenvalue one and all other eigenvalues contained in the open unit disk centered at the origin, and the local objective functions $g_j, 1 \leq j \leq n$, satisfy (V.2). Then*

$$\|\mathbf{a}_{\text{sde}}^T \mathbf{h}(k) - \tilde{w} \nabla f(\bar{\mathbf{x}}(k))\|_2 \leq L \|\mathbf{a}_{\text{sde}}\|_1 \|\tilde{\mathbf{z}}(k)\|_{2, \infty}. \quad (\text{V.9})$$

By Proposition V.4, the iterative algorithm (V.8) can be considered as the gradient descent algorithm (I.3) with inexact gradient update. Then following a standard argument, we have the boundedness of $\bar{\mathbf{x}}(k), k \geq 1$, when the objective function f is convex, see Appendix E for a detailed proof.

Proposition V.5. *Let the matrix \mathbf{A}_{sde} , the sequence $\{\alpha_k\}_{k=0}^{\infty}$ and the local objective functions $g_j, 1 \leq j \leq n$, be as in Theorem V.1. If the global objective function f is convex, then*

$$\|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 \leq C_2 \quad \text{for all } k \geq 0, \quad (\text{V.10})$$

where C_2 is the constant in Theorem V.1.

The estimate in (V.10) can be improved if the objective function f has the strongly convex property (V.3), see Appendix E for a detailed proof.

Proposition V.6. *Let the matrix \mathbf{A}_{sde} , the sequence $\{\alpha_k\}_{k=0}^{\infty}$, the local objective functions $g_j, 1 \leq j \leq n$, and the global objective function f be as in Theorem V.1. Then there exists a positive constant C such that*

$$\begin{aligned} \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 \leq & C \exp\left(-\tilde{w}A \sum_{j=0}^{k-1} \alpha_j\right) \left\{ \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 \right. \\ & + \sum_{j=0}^{k-1} \exp\left(\tilde{w}A \sum_{l=0}^j \alpha_l\right) \left(M^2 \alpha_j^2 \right. \\ & \left. \left. + L^2 \max_{1 \leq i \leq n} \|\mathbf{x}_i(j) - \bar{\mathbf{x}}(j)\|_2^2\right) \right\}, \quad (\text{V.11}) \end{aligned}$$

hold for all $k \geq 2$.

Applying (V.11) for the illustrative examples (I.4) of step sizes, and using (IV.20), we can find a positive constant C such that

$$\begin{aligned} \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 \leq & \exp\left(-\tilde{w}A \sum_{l=0}^{k-1} (l+a)^{-\theta}\right) \left\{ \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 \right. \\ & + (M^2 + L^2 C^2) \sum_{j=0}^{k-1} \exp\left(\tilde{w}A \sum_{l=0}^j (l+a)^{-\theta}\right) (j+a)^{-2\theta} \left. \right\} \\ \leq & \exp\left(-\frac{\tilde{w}A}{1-\theta} ((k+a)^{1-\theta} - a^{1-\theta})\right) \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 \\ & + (M^2 + L^2 C^2) \exp\left(-\frac{\tilde{w}A}{1-\theta} (k+a)^{1-\theta}\right) \\ & \times \sum_{j=0}^{k-1} \exp\left(\frac{\tilde{w}A}{1-\theta} (j+1+a)^{1-\theta}\right) (j+a)^{-2\theta} \\ \leq & \tilde{C}^2 (k+a)^{-\theta}, \quad k \geq 1, \quad (\text{V.12}) \end{aligned}$$

where the first estimate holds by (V.11), the second one follows from the observation

$$\sum_{l=m}^{k-1} (l+a)^{-\theta} \leq \frac{(k+a)^{1-\theta} - (m+a)^{1-\theta}}{1-\theta}$$

for $0 \leq m \leq k-1$, and the third one is obtained from the boundedness of the sequences $(k+a)^{\theta} \exp(-\frac{\tilde{w}A}{1-\theta} (k+a)^{1-\theta})$

and $(k+a+1)^{1-\theta} - (k+a)^{1-\theta}, k \geq 0$, and the integrability of $\int_0^\infty \exp(-\frac{\tilde{w}A}{1-\theta}t)(t+1)^{\theta/(1-\theta)}dt$, and the following estimate

$$\begin{aligned} & \sum_{j=0}^{k-1} \exp\left(\frac{\tilde{w}A}{1-\theta}(j+1+a)^{1-\theta}\right)(j+a)^{-2\theta} \\ & \leq 3^{2\theta} \int_0^k \exp\left(\frac{\tilde{w}A}{1-\theta}(x+a+1)^{1-\theta}\right)(x+a+1)^{-2\theta} dx \\ & \leq 3^{2\theta} (1-\theta)^{-1} \int_0^{a_k} \exp\left(\frac{\tilde{w}A}{1-\theta}((k+a+1)^{1-\theta}-t)\right) \\ & \quad \times ((k+a+1)^{1-\theta}-t)^{-\theta/(1-\theta)} dt \\ & \leq 3^{2\theta} (1-\theta)^{-1} \exp\left(\frac{\tilde{w}A}{1-\theta}(k+a+1)^{1-\theta}\right)(k+a+1)^{-\theta} \\ & \quad \times \int_0^{a_k} \exp\left(-\frac{\tilde{w}A}{1-\theta}t\right)(t+1)^{\theta/(1-\theta)} dt \\ & \leq 3^{2\theta} (1-\theta)^{-1} \exp\left(\frac{\tilde{w}A}{1-\theta}(k+a+1)^{1-\theta}\right)(k+a)^{-\theta} \\ & \quad \times \int_0^\infty \exp\left(-\frac{\tilde{w}A}{1-\theta}t\right)(t+1)^{\theta/(1-\theta)} dt, \end{aligned}$$

where $a_k = (k+1+a)^{1-\theta} - (a+1)^{1-\theta} \leq (k+1+a)^{1-\theta} - 1$. Therefore

$$\|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\| \leq \tilde{C}(k+a)^{-\theta/2}, \quad k \geq 1. \quad (\text{V.13})$$

This implies that $\bar{\mathbf{x}}(k), k \geq 0$, has faster convergence to the limit $\bar{\mathbf{x}}$ when we select larger $\theta \in (1/2, 1)$, however comparing the convergence rate $O(\log k/\sqrt{k})$ for the uncoded incremental gradient methods [1, 10], the convergence rate in (V.13) is always slow, even it is very close when θ is close to one. After careful verification in the above estimation, we observe that the constant \tilde{C} could be smaller if the second eigenvalue $\lambda_2(\mathbf{A}_{\text{sde}})$ associated with the decoding matrix \mathbf{A} is smaller.

Set $W = \max_{1 \leq i \leq n} w_i$. Observe that

$$f(\bar{\mathbf{x}}(k)) - f(\bar{\mathbf{x}}) = \langle \nabla f(t\bar{\mathbf{x}}(k)) + (1-t)\bar{\mathbf{x}}, \bar{\mathbf{x}}(k) - \bar{\mathbf{x}} \rangle$$

for some $0 \leq t \leq 1$, and that

$$\|\nabla f(\mathbf{x})\|_2 \leq \|\nabla f(\mathbf{x}) - \nabla f(\bar{\mathbf{x}})\|_2 \leq LW^{-1}\|\mathbf{x} - \bar{\mathbf{x}}\|_2,$$

where the second inequality follows from (V.2) and the row stochastic property for the matrix \mathbf{A}_{sde} . This together with Proposition V.6 proves that

$$\begin{aligned} f(\bar{\mathbf{x}}) & \leq f(\bar{\mathbf{x}}(k)) \leq f(\bar{\mathbf{x}}) + LW^{-1}\|\mathbf{x} - \bar{\mathbf{x}}\|_2 \\ & \leq f(\bar{\mathbf{x}}) + LW^{-1} \exp\left(-\tilde{w}A \sum_{j=0}^{k-1} \alpha_j\right) \\ & \quad \times \left\{ \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \sum_{j=0}^{k-1} \exp\left(\tilde{w}A \sum_{l=0}^j \alpha_l\right) \right. \\ & \quad \left. \times \left(M^2 \alpha_j^2 + L^2 \max_{1 \leq i \leq n} \|\mathbf{x}_i(j) - \bar{\mathbf{x}}(j)\|_2^2\right) \right\} \quad (\text{V.14}) \end{aligned}$$

hold for all $k \geq 2$. For the case that step size $\alpha_k, k \geq 0$ chosen as in (I.4), we can use the similar argument used to prove (V.13) and show that

$$|f(\bar{\mathbf{x}}(k)) - f(\bar{\mathbf{x}})| \leq C(k+a)^{-\theta}, \quad k \geq 1, \quad (\text{V.15})$$

for some positive constant C .

We finish this section with the proof of Theorem V.1 under the assumption that Proposition V.6 holds.

Proof of Theorem V.1. By (II.4) and (IV.21), we have

$$\sum_{j=0}^{\infty} M^2 \alpha_j^2 + L^2 \max_{1 \leq i \leq n} \|\mathbf{x}_i(j) - \bar{\mathbf{x}}(j)\|_2^2 < \infty, \quad (\text{V.16})$$

and

$$\lim_{k \rightarrow \infty} \exp\left(-\tilde{w}A \sum_{j=l}^k \alpha_j\right) = 0, \quad l \geq 0. \quad (\text{V.17})$$

Combining (V.11), (V.16) and (V.17) proves the desired limit in (V.4) by the dominated convergence theorem. \square

VI. NUMERICAL SIMULATIONS

In this section, we consider the following unconstrained convex optimization problem on a network

$$\arg \min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{G}\mathbf{x} - \mathbf{y}\|_2^2, \quad (\text{VI.1})$$

where the network contains n regions with each region of the partition equipped with a worker, \mathbf{G} is a random matrix of size $Q \times N$ whose entries are independent and identically distributed standard normal random variables, and

$$\mathbf{y} = \mathbf{G}\mathbf{x}_o \in \mathbb{R}^Q \quad (\text{VI.2})$$

has entries of \mathbf{x}_o being identically independent random variables sampled from the uniform bounded random distribution between -1 and 1 . The solution $\bar{\mathbf{x}}$ of the above optimization problem is the least squares solution of the overdetermined system $\mathbf{y} = \mathbf{G}\mathbf{x}_o, \mathbf{x}_o \in \mathbb{R}^N$. In this section, we demonstrate the performance of the CoDGrad algorithm (II.3) to solve the convex optimization problem (VI.1) and also compare it with the performance of the conventional distributed gradient descent algorithm (DGD) under the CTA prototype (I.6).

Assume that the network has n active nodes. Then we can repartition the network into n regions around those n nodes, and accordingly, the random measurement matrix \mathbf{G} , the measurement data \mathbf{y} , and the objective function $f(\mathbf{x}) := \|\mathbf{G}\mathbf{x} - \mathbf{y}\|_2^2$ in (VI.1) as follows,

$$f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}) := \sum_{i=1}^n \|\mathbf{G}_i \mathbf{x} - \mathbf{y}_i\|_2^2.$$

In our simulations, we assume that the repartitioned regions have the same size, i.e., the number of rows in \mathbf{G}_i and lengths of vectors \mathbf{y}_i , for all i where $1 \leq i \leq n$ are all equal. Shown in Figure 2 are two undirected graphs to describe data exchanging structure for active nodes of a 3-node and 5-node network respectively.

In our simulations, we take $(Q, N) = (225, 75)$ for Figure 3 and $(M, N) = (250, 50)$ for Figure 4. We use absolute error

$$\text{AE} := \max_{1 \leq i \leq n} \|\mathbf{x}_i(k) - \mathbf{x}_o\|_2 / \|\mathbf{x}_o\|_2$$

and consensus error

$$\text{CE} := \max_{1 \leq i \leq n} \|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|_2 / \|\mathbf{x}_o\|_2$$

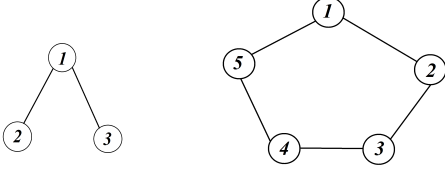


Fig. 2. Data exchanging structures with three/five-node network

to measure the performance of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6), where n is the number of active nodes in the network.

In the first simulation where there are 3 nodes with its topology described on the left of Figure 2, we take the coding matrix \mathbf{B} and the decoding matrix \mathbf{A} as follows:

$$\mathbf{B} = \begin{pmatrix} 1 & -\frac{5}{4} & 0 \\ 0 & 1 & \frac{4}{9} \\ \frac{9}{5} & 0 & 1 \end{pmatrix} \text{ and } \mathbf{A} = \begin{pmatrix} 0 & 1 & \frac{5}{9} \\ 1 & \frac{9}{4} & 0 \\ -\frac{4}{5} & 0 & 1 \end{pmatrix}. \quad (\text{VI.3})$$

The above coding/decoding matrix pair satisfies (I.9) and the corresponding row stochastic matrix in (II.6) is

$$\mathbf{A}_{\text{sde}} = \begin{pmatrix} 0 & 9/14 & 2/7 & 0 & 0 & 0 \\ 4/13 & 9/13 & 0 & 0 & 0 & 0 \\ 0 & 0 & 5/9 & 4/9 & 0 & 0 \\ 0 & 9/14 & 5/14 & 0 & 0 & 0 \\ 4/13 & 9/13 & 0 & 0 & 0 & 0 \\ 0 & 0 & 5/9 & 4/9 & 0 & 0 \end{pmatrix}.$$

In Figure 3, we present the performance of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6) with absolute and consensus metric being the average of the corresponding metrics over 100 trials, where random measurement matrix \mathbf{G} has independent and identically distributed standard normal random variables as its entries, the original vector \mathbf{x}_o is identically independent random variables uniform distributed in $[-1, 1]$, and step sizes are $\alpha_k = (k + 300)^{-0.75}$ and $(k + 500)^{-0.85}$, $k \geq 0$, respectively.

In the second simulation, the network has 5 active nodes with data exchanging structure described on the right of Figure 2. In that simulation, the coding/decoding matrices are given by

$$\mathbf{B} = \begin{pmatrix} 1 & 2 & 1/2 & 0 & 0 \\ 0 & -1 & 3 & 4 & 0 \\ 0 & 0 & -5/2 & -3 & 1 \\ 1 & 0 & 0 & 1/5 & 13/5 \\ 2 & 1 & 0 & 0 & 4 \end{pmatrix} \quad (\text{VI.4})$$

and

$$\mathbf{A} = \begin{pmatrix} 1/2 & 1/4 & 0 & 0 & 1/4 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & -1 & -8/5 & 1 & 0 \\ 0 & 0 & -2/5 & -1 & 1 \\ 2 & 0 & 0 & 5 & -3 \end{pmatrix} \quad (\text{VI.5})$$

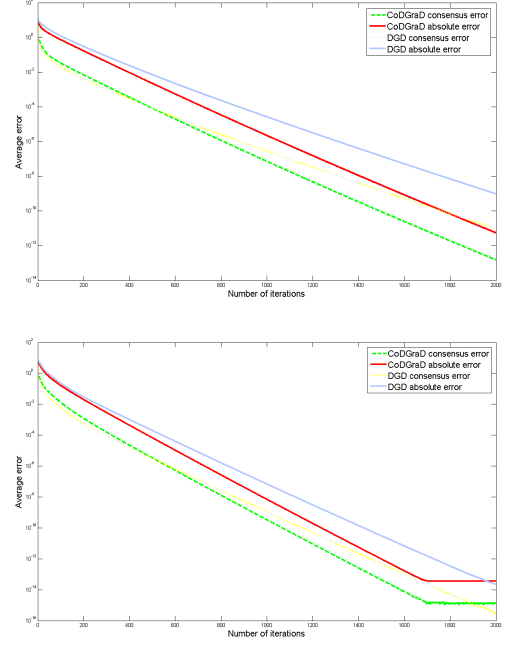


Fig. 3. Performance comparison of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6) over a three active nodes network with $(Q, N) = (225, 225)$ and step sizes $\alpha_k = (k + 300)^{-0.75}$ (top) and $(k + 500)^{-0.85}$, $k \geq 0$ (bottom).

respectively. The above coding/decoding matrix pair satisfies (I.9) and the corresponding row stochastic matrix in (II.6) is

$$\mathbf{A}_{\text{sde}} = \begin{pmatrix} \frac{1}{2} & \frac{1}{4} & 0 & 0 & \frac{1}{4} & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{5}{18} & 0 & 0 & \frac{5}{18} & \frac{8}{18} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{5}{12} & 0 & 0 & \frac{1}{6} & \frac{5}{12} & 0 \\ \frac{1}{5} & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 & \frac{3}{10} \\ \frac{1}{2} & \frac{1}{4} & 0 & 0 & \frac{1}{4} & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{5}{18} & 0 & 0 & \frac{5}{18} & \frac{8}{18} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{5}{12} & 0 & 0 & \frac{1}{6} & \frac{5}{12} & 0 \\ \frac{1}{5} & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 & \frac{3}{10} \end{pmatrix}.$$

Shown in Figure 4 is the performance of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6), where the absolute metric and consensus metric are the average of the corresponding metrics over 100 trials with random measurement matrix \mathbf{G} and the original vector \mathbf{x}_o being selected as in the first simulation, and step sizes being $\alpha_k = (k + 800)^{-0.90}$ and $(k + 500)^{-0.95}$, $k \geq 0$, respectively.

From the above simulations, we observe that the CoDGrAD algorithm (II.3) has much better performance than the CTA algorithm (I.6) in reaching consensus. Even though $\bar{\mathbf{x}}(k)$ satisfies a gradient descent algorithm (V.8) with an inexact global gradient, see Proposition V.4, our simulations indicate that the CoDGrAD algorithm (II.3) still has comparable performance in the absolute error with the DGD algorithm under the CTA prototype (I.6). Also we can conceive from the simulations that the CoDGrAD algorithm (II.3) has faster convergence for a smaller exponent $\theta \in (1/2, 1]$, which confirms its convergence rate estimate in (IV.20) and (V.13). On the other hand, our

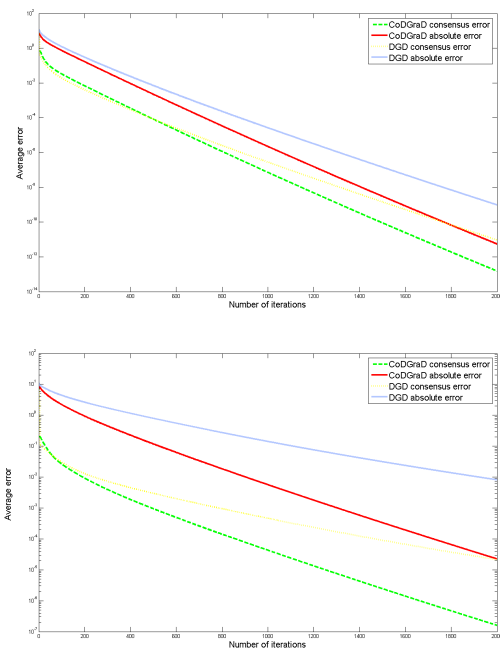


Fig. 4. Performance comparison of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6) over a five node network with $(Q, N) = (225, 225)$ and step sizes $\alpha_k = (k + 800)^{-0.9}$ (top) and $(k + 500)^{-0.95}, k \geq 0$ (bottom).

simulations also indicate that decreasing the exponent θ moves the CoDGrAD algorithm (II.3) into the instability phase, which could directly be related to the sparsity of the network (i.e., the graph degree of the corresponding network and the degree distribution of vertices). It is worth mentioning that we can adequately calibrate this instability by increasing the value of a in our illustrative examples of step sizes (I.4) for a fixed exponent θ . Thus we can anticipate in Figure 3 that the increase in the value to $\theta = 0.85$ degraded the convergence so that the CoDGrAD algorithm became closer in performance to the DGD algorithm. While in Figure 4 we realize that the lower value of $\theta = 0.75$ is impermissible since the CoDGrAD algorithm will considerably enter the instability region while a higher value of $\theta = 0.9$ favors a better convergence rate and the highest value of $\theta = 0.95$ degraded the convergence again.

In these simulations, we have compared the performance of the CoDGrAD algorithm (II.3) and the DGD algorithm (I.6) over the described 3-node and 5-node networks with subsystems on nodes being overdetermined. Therefore, a least squares solutions on each node will correspond to the unique solution of a strongly convex function and will consequently correspond to the least squared solution of the whole network, that is, the unique solution of the strongly convex global function f . It is observed that the errors decrease significantly in 2000 iterations where CoDGrAD outperforms DGD in reaching the unique minimizer (i.e., absolute error) and in reaching consensus. While for both algorithms the consensus error decreases at a higher rate than the absolute error meaning that the workers become closer in their estimates while they all drift towards the unique solution. Our further simulations indicate that convergence behaviors of the CoDGrAD algo-

rithm (II.3) and the DGD algorithm (I.6) depends directly on maximal condition number of matrices $\mathbf{G}_i, 1 \leq i \leq n$, cf. (IV.22) and (V.14) where A is closely related to the maximal condition number in the current setting.

VII. CONCLUSIONS

In this paper, we proposed the Code-Based Distributed Gradient Descent algorithm (II.3) to solve a convex optimization problem over a large network with some workers being stragglers due to the failure or heavy delay on computing or communicating. The proposed algorithm is a distributed version of gradient descent algorithm with inexact gradient updating, and it has better performance in reaching consensus as we apply the row stochastic matrix associated with the coding/decoding scheme. The convergence rate of the proposed CoDGrAD algorithm depends on the topological structure of the network, the second largest eigenvalue of row stochastic matrix in magnitude, and the updating step sizes in the algorithm. Moreover, our coding scheme does not necessarily comply with the conventional paradigm of decomposing the global convex function onto a summand of local convex functions and hence our coding/decoding scheme may shed new light on distributed inexact (stochastic) gradient descent algorithms. We wish that this work on CoDGrAD will serve as a starting point for a full-fledged investigation on static and time-varying networks especially in the field of federated decentralized learning that we would like to continue resolving in the coming future.

Acknowledgement: The authors would like to thank all reviewers for their constructive comments for the improvement of the manuscript. This work is partially supported by the National Science Foundation (DMS-1816313).

APPENDIX

A. Proof of Proposition II.2

Observe that $\tilde{\mathbf{A}}_+ + \tilde{\mathbf{A}}_+ = |\tilde{\mathbf{A}}|$. Then, in order to establish the equivalence in the proposition, it suffices to prove that for any positive integer $k \geq 1$ and nonzero complex number λ , the null space of $\left(\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix} - \lambda \mathbf{I}_{2n} \right)^k$ and the one of $(\mathbf{A} + \mathbf{B} - \lambda \mathbf{I}_n)^k$, to be denoted by $\mathcal{N}_k^1(\lambda)$ and $\mathcal{N}_k^2(\lambda)$ respectively, have the same dimension, where \mathbf{A} and \mathbf{B} are two square matrices of size $n \times n$.

Set $\mathbf{C} = \mathbf{A} + \mathbf{B}$. By induction on $j \geq 1$, we can show that

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix}^j = \begin{pmatrix} \mathbf{C}^{j-1} & \mathbf{0}_{n \times n} \\ \mathbf{0}_{n \times n} & \mathbf{C}^{j-1} \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix}. \quad (\text{A.1})$$

Therefore any $\mathbf{u} \in \mathcal{N}_k^2(\lambda)$, i.e., $(\mathbf{C} - \lambda \mathbf{I}_n)^k \mathbf{u} = \mathbf{0}_{n \times 1}$, we have

$$\begin{aligned} & \left(\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix} - \lambda \mathbf{I}_{2n} \right)^k \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix} \\ &= \sum_{j=1}^k \binom{k}{j} (-\lambda)^{k-j} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix}^j \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix} + (-\lambda)^k \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix} \\ &= \begin{pmatrix} (\mathbf{C} - \lambda \mathbf{I}_n)^k \mathbf{u} \\ (\mathbf{C} - \lambda \mathbf{I}_n)^k \mathbf{u} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_{n \times 1} \\ \mathbf{0}_{n \times 1} \end{pmatrix}, \end{aligned}$$

where the second equality follows from (A.1). This proves that

$$\mathcal{N}_k^1(\lambda) \supset \left\{ \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix}, \mathbf{u} \in \mathcal{N}_k^2(\lambda) \right\}. \quad (\text{A.2})$$

On the other hand, for any $\mathbf{u}, \mathbf{w} \in \mathbb{C}^n$ satisfying

$$\left(\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix} - \lambda \mathbf{I}_{2n} \right)^k \begin{pmatrix} \mathbf{u} \\ \mathbf{w} \end{pmatrix} = \begin{pmatrix} \mathbf{0}_{n \times 1} \\ \mathbf{0}_{n \times 1} \end{pmatrix},$$

we obtain from (A.1) that

$$\begin{aligned} (-\lambda)^k \begin{pmatrix} \mathbf{u} \\ \mathbf{w} \end{pmatrix} &= - \sum_{j=1}^k \binom{k}{j} (-\lambda)^{k-j} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{pmatrix}^j \begin{pmatrix} \mathbf{u} \\ \mathbf{w} \end{pmatrix} \\ &= - \sum_{j=1}^k \binom{k}{j} (-\lambda)^{k-j} \begin{pmatrix} \mathbf{C}^{j-1}(\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{w}) \\ \mathbf{C}^{j-1}(\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{w}) \end{pmatrix}. \end{aligned}$$

This implies that $\mathbf{w} = \mathbf{u}$. Substituting the above equality back, we get

$$(-\lambda)^k \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix} = - \sum_{j=1}^k \binom{k}{j} (-\lambda)^{k-j} \begin{pmatrix} \mathbf{C}^j \mathbf{u} \\ \mathbf{C}^j \mathbf{u} \end{pmatrix}$$

which implies that $(\mathbf{C} - \lambda \mathbf{I}_n)^k \mathbf{u} = \mathbf{0}_{n \times 1}$. Hence

$$\mathcal{N}_k^1(\lambda) \subset \left\{ \begin{pmatrix} \mathbf{u} \\ \mathbf{u} \end{pmatrix}, \mathbf{u} \in \mathcal{N}_k^2(\lambda) \right\}. \quad (\text{A.3})$$

Combining (A.2) and (A.3) completes the proof.

B. Proof of Proposition IV.2

Denote the spectrum of a square matrix \mathbf{A} by $\sigma(\mathbf{A})$. By the assumption on the matrix \mathbf{A}_{sde} , its spectrum $\sigma(\mathbf{A}_{\text{sde}})$ satisfies

$$\sigma(\mathbf{A}_{\text{sde}}) \subset \{1\} \cup \{z, |z| < 1\}, \quad (\text{A.4})$$

and the eigenspace associated with eigenvalue one is given by

$$N(\mathbf{A}_{\text{sde}} - \mathbf{I}) = \text{span}\{\mathbf{1}_{2n}\}. \quad (\text{A.5})$$

Combining (IV.11), (IV.14), (A.4) and (A.5), we obtain that the spectrum of $\mathbf{A}_{\text{sde}} - \mathbf{P}$ is contained in the open unit disk,

$$\sigma(\mathbf{A}_{\text{sde}} - \mathbf{P}) = (\sigma(\mathbf{A}_{\text{sde}}) \setminus \{1\}) \cup \{0\} \subset \{z, |z| \leq |\lambda_2(\mathbf{A}_{\text{sde}})|\}. \quad (\text{A.6})$$

Therefore there exists a positive constant C_1 such that

$$\|(\mathbf{A}_{\text{sde}} - \mathbf{P})^k\|_{\mathcal{B}^\infty} \leq C_1 \left(\frac{1 + 2|\lambda_2(\mathbf{A}_{\text{sde}})|}{3} \right)^k, \quad k \geq 1, \quad (\text{A.7})$$

where $\|\mathbf{A}\|_{\mathcal{B}^\infty} = \sup_{\|\mathbf{x}\|_\infty=1} \|\mathbf{A}\mathbf{x}\|_\infty$.

By (I.8), (IV.1), (IV.7) and (IV.8), we have

$$\|\mathbf{h}(k)\|_{2,\infty} \leq \sup_{1 \leq i \leq n} \sup_{\mathbf{x} \in \mathbb{R}^N} \|\nabla g_i(\mathbf{x})\|_2 \leq M. \quad (\text{A.8})$$

Then combining (IV.16), (A.7) and (A.8) completes the proof.

C. Proof of Lemma V.3

By (IV.14), we have $(\mathbf{A}_{\text{sde}} - \mathbf{P})^k = \mathbf{A}_{\text{sde}}^k - \mathbf{P}$, $k \geq 1$. Therefore $\lim_{k \rightarrow \infty} \mathbf{A}_{\text{sde}}^k = \mathbf{P}$ by (A.7). This, together with the observation that all entries of $\mathbf{A}_{\text{sde}}^k$, $k \geq 1$ are nonnegative. Hence the required estimate implies that all entries of \mathbf{a}_{sde} are nonnegative and the desired estimate on weights follows.

D. Proof of Proposition V.4

By (IV.7), (IV.18) and (V.2), we have

$$\begin{aligned} \|\mathbf{h}_i(k) - \nabla g_i(\bar{\mathbf{x}}(k))\|_2 &= \|\nabla g_i(\mathbf{z}_i(k)) - \nabla g_i(\bar{\mathbf{x}}(k))\|_2 \\ &\leq L\|\mathbf{z}_i(k) - \bar{\mathbf{x}}(k)\|_2 \leq L\|\tilde{\mathbf{z}}(k)\|_{2,\infty} \end{aligned}$$

for $1 \leq i \leq n$, and

$$\|\mathbf{h}_i(k) + \nabla g_{i-n}(\bar{\mathbf{x}}(k))\|_2 \leq L\|\tilde{\mathbf{z}}(k)\|_{2,\infty}$$

for $n+1 \leq i \leq 2n$. Therefore

$$\|\mathbf{h}(k) - \tilde{\mathbf{h}}(k)\|_{2,\infty} \leq L\|\tilde{\mathbf{z}}(k)\|_{2,\infty}, \quad (\text{A.9})$$

where

$$\tilde{\mathbf{h}}(k) = \begin{pmatrix} \nabla G(\bar{\mathbf{x}}(k)) \\ -\nabla G(\bar{\mathbf{x}}(k)) \end{pmatrix} \quad \text{and} \quad \nabla G(\mathbf{x}) = \begin{pmatrix} \nabla g_1(\mathbf{x}) \\ \vdots \\ \nabla g_n(\mathbf{x}) \end{pmatrix}.$$

Therefore

$$\begin{aligned} &\|\mathbf{a}_{\text{sde}}^T \mathbf{h}(k) - \mathbf{a}_{\text{sde}}^T \tilde{\mathbf{h}}(k)\|_2 \\ &\leq \|\mathbf{a}_{\text{sde}}\|_1 \|\mathbf{h}(k) - \tilde{\mathbf{h}}(k)\|_{2,\infty} \leq L\|\mathbf{a}_{\text{sde}}\|_1 \|\tilde{\mathbf{z}}(k)\|_{2,\infty}, \end{aligned} \quad (\text{A.10})$$

where the second estimate follows from (A.9).

Observe from (I.9) that $\mathbf{A}_{\text{sde}} \tilde{\mathbf{h}}(k) = \nabla f(\bar{\mathbf{x}}(k)) \mathbf{w}$. This together with (IV.12) implies that

$$\mathbf{a}_{\text{sde}}^T \tilde{\mathbf{h}}(k) = \mathbf{a}_{\text{sde}}^T \mathbf{A}_{\text{sde}} \tilde{\mathbf{h}}(k) = \tilde{w} \nabla f(\bar{\mathbf{x}}(k)). \quad (\text{A.11})$$

Combining (A.10) and (A.11) proves the desired estimate (V.9).

E. Proof of Proposition V.5

Set

$$\beta_k = M^2 \alpha_k^2 + L^2 \|\tilde{\mathbf{z}}(k)\|_{2,\infty}^2, \quad k \geq 0. \quad (\text{A.12})$$

By (IV.7), (IV.1), (V.8), (V.9) and (A.9), we obtain

$$\begin{aligned} \|\bar{\mathbf{x}}(k+1) - \bar{\mathbf{x}}\|_2^2 &= \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 + \alpha_k^2 \|\mathbf{a}_{\text{sde}}^T \mathbf{h}(k)\|_2^2 \\ &\quad - 2\alpha_k \langle \mathbf{a}_{\text{sde}}^T \mathbf{h}(k), \bar{\mathbf{x}}(k) - \bar{\mathbf{x}} \rangle_N \\ &= \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 + \alpha_k^2 \|\mathbf{a}_{\text{sde}}^T \tilde{\mathbf{h}}(k)\|_2^2 \\ &\quad - 2\alpha_k \langle \mathbf{a}_{\text{sde}}^T \tilde{\mathbf{h}}(k), \bar{\mathbf{x}}(k) - \bar{\mathbf{x}} \rangle_N \\ &\quad - 2\alpha_k \langle \mathbf{a}_{\text{sde}}^T \tilde{\mathbf{h}}(k), \bar{\mathbf{x}}(k) - \bar{\mathbf{x}} \rangle_N \\ &\leq \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 + M^2 \|\mathbf{a}_{\text{sde}}\|_1^2 \alpha_k^2 \\ &\quad + 2L \|\mathbf{a}_{\text{sde}}\|_1 \alpha_k \|\tilde{\mathbf{z}}(k)\|_{2,\infty} \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2 \\ &\leq (1 + \alpha_k^2) \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 + \beta_k, \end{aligned} \quad (\text{A.13})$$

where we also use the positivity of \tilde{w} in Lemma V.3, $\|\mathbf{a}_{\text{sde}}\|_1 = 1$ and the convexity of the objective function f ,

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \bar{\mathbf{x}} \rangle_N \geq 0, \quad \mathbf{x} \in \mathbb{R}^N. \quad (\text{A.14})$$

Applying (A.13) repeatedly, we get

$$\begin{aligned} &\|\bar{\mathbf{x}}(k+1) - \bar{\mathbf{x}}\|_2^2 \\ &\leq \prod_{j=0}^k (1 + \alpha_j^2) \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \beta_k + \sum_{j=0}^{k-1} \beta_j \prod_{j'=j+1}^k (1 + \alpha_{j'}^2) \\ &\leq \exp\left(\sum_{j=0}^k \alpha_j^2\right) \left(\|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \sum_{j=0}^k \beta_j \right) \\ &\leq \exp\left(\sum_{j=0}^{\infty} \alpha_j^2\right) \left(\|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \sum_{j=0}^{\infty} \beta_j \right), \quad k \geq 1. \end{aligned} \quad (\text{A.15})$$

This together with (II.4) and (IV.21) proves the desired bound (V.10) for the sequence $\bar{\mathbf{x}}(k)$, $k \geq 0$.

F. Proof of Proposition V.6

By (II.4), without a loss of generality, we assume that $0 \leq \alpha_k \leq \tilde{w}A$ for all $k \geq 0$. Then for $k \geq 1$, following the argument in (A.13) with the convexity (A.14) replaced by the strong convexity (V.3), we obtain that

$$\|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 \leq (1 - \tilde{w}A\alpha_{k-1})\|\bar{\mathbf{x}}(k-1) - \bar{\mathbf{x}}\|_2^2 + \beta_{k-1},$$

cf. (A.13). Applying the above estimate repeatedly leads to

$$\begin{aligned} \|\bar{\mathbf{x}}(k) - \bar{\mathbf{x}}\|_2^2 &\leq \prod_{j=0}^{k-1} (1 - \tilde{w}A\alpha_j) \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \beta_{k-1} \\ &\quad + \sum_{j=0}^{k-2} \beta_j \prod_{j'=j+1}^{k-1} (1 - \tilde{w}A\alpha_{j'}) \\ &\leq \exp\left(-\tilde{w}A \sum_{j=0}^{k-1} \alpha_j\right) \|\bar{\mathbf{x}}(0) - \bar{\mathbf{x}}\|_2^2 + \beta_{k-1} \\ &\quad + \sum_{j=0}^{k-2} \exp\left(-\tilde{w}A \sum_{j=j+1}^{k-1} \alpha_j\right) \beta_j, \end{aligned}$$

where $\beta_j, j \geq 0$, is given in (A.12). This proves the desired estimate (V.11).

REFERENCES

- [1] A. Nedic and D. Bertsekas (2001). Convergence rate of incremental subgradient algorithms. *Stochastic Optimization: Algorithms and Applications*, pp. 223–264. Springer.
- [2] C. R. Da Silva, B. Choi, and K. Kim (2007). Distributed spectrum sensing for cognitive radio systems. In *2007 Information Theory and Applications Workshop*, pp. 120–123. IEEE.
- [3] B. Johansson (2008). On distributed optimization in networked systems. PhD. thesis, KTH.
- [4] A. Nedic, A. Ozdaglar, and P. A. Parrilo (2010). Constrained consensus and optimization in multi-agent networks. *IEEE Transactions on Automatic Control*, **55**(4), 922–938.
- [5] D. Fu, L. Han, L. Liu, Q. Gao, and Z. Feng (2015). An efficient centralized algorithm for connected dominating set on wireless networks. *Procedia Computer Science*, **56**, 162–167.
- [6] C. Cheng, Y. Jiang, and Q. Sun (2019). Spatially distributed sampling and reconstruction. *Applied and Computational Harmonic Analysis*, **47**, 109–148.
- [7] J. Jiang, C. Cheng, and Q. Sun (2019). Nonsub-sampled graph filter banks: theory and distributed algorithms. *IEEE Transactions on Signal Processing*, **67**, 3938–3953.
- [8] A. Nedic (2015). Convergence rate of distributed averaging dynamics and optimization in networks. *Foundations and Trends in Systems and Control*, **2**, 1–100.
- [9] A. S. Bedi and K. Rajawat (2018). Asynchronous incremental stochastic dual descent algorithm for network resource allocation. *IEEE Transactions on Signal Processing*, **66**(9), 2229–2244.
- [10] M. Gürbüzbalaban, A. Ozdaglar, and P. Parrilo (2019). Convergence rate of incremental gradient and incremental newton methods. *SIAM J. Optim.*, **29**(4), 2542–2565.
- [11] N. Emirov, G. Song, and Q. Sun (2021). A divide-and-conquer algorithm for distributed optimization on networks. *arXiv:2112.02197*
- [12] N. Takahashi, I. Yamada, and A. H. Sayed (2010). Diffusion least-mean squares with adaptive combiners: formulation and performance analysis. *IEEE Transactions on Signal Processing*, **58**(9), 4795–4810.
- [13] F. S. Cattivelli and A. H. Sayed (2010). Diffusion LMS strategies for distributed estimation. *IEEE Transactions on Signal Processing*, **58**(3), 1035–1048.
- [14] D. P. Bertsekas (2011). Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. In *Optimization for Machine Learning*, S. Sra, S. Nowozin, and S. Wright, Eds., pp. 1–38. MIT Press.
- [15] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, **3**, 1–122.
- [16] A. H. Sayed, S. Barbarossa, S. Theodoridis, and I. Yamada (2013). Adaptation and learning over complex networks. *IEEE Signal Processing Magazine*, **30**(3), 14–15.
- [17] D. Needell, R. Ward, and N. Srebro (2014). Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. In *Advances in Neural Information Processing Systems*, pp. 1017–1025.
- [18] A. H. Sayed (2014). Adaptation, learning, and optimization over networks. *Foundations and Trends in Machine Learning*, **7**, 311–801.
- [19] V. J. Mathews, and Z. Xie (1993). A stochastic gradient adaptive filter with gradient adaptive step size. *IEEE Transactions on Signal Processing*, **41**(6), 2075–2087.
- [20] H. Robbins, and S. Monro (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, **22**(3), 400–407.
- [21] N. N. Schraudolph (1999). Local gain adaptation in stochastic gradient descent. In *1999 Ninth International Conference on Artificial Neural Networks ICANN 99*, pp. 569–574. IEEE.
- [22] J. H. Friedman (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*, **38**(4), 367–378.
- [23] L. Bottou (2012). Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, G. Montavon, G. B. Orr, K.-R. Müller eds, pp 421–436.
- [24] M. D. Zeiler (2012). Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*.
- [25] D. P. Kingma and J. Ba (2014). Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, 2014.

- [26] M. Hardt, B. Recht, and Y. Singer (2016). Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 1225–1234, 2016.
- [27] C. Tan, S. Ma, Y.-H. Dai, and Y. Qian (2016). Barzilai-borwein step size for stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 685–693.
- [28] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le, M. Z. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, and A. Y. Ng (2012). Large scale distributed deep networks. In *Proceeding NIPS'12 Proceedings of the 25th International Conference on Neural Information Processing Systems*, pp. 1223–1231.
- [29] Q. Ho, J. Cipar, H. Cui, S. Lee, J. K. Kim, P. B. Gibbons, G. A. Gibson, G. Ganger, and E. P. Xing (2013). More effective distributed ml via a stale synchronous parallel parameter server. In *Proceeding NIPS'13 Proceedings of the 26th International Conference on Neural Information Processing Systems*, pp. 1223–1231.
- [30] E. Atallah, N. Rahnavard (2018). A Code-Based Distributed Gradient Descent Method. In *2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton) IEEE*. pp 951–958.
- [31] M. Li, D. G. Andersen, A. J. Smola, and K. Yu (2014). Communication efficient distributed machine learning with the parameter server. In *Proceeding NIPS'14 Proceedings of the 25th International Conference on Neural Information Processing Systems*, pp. 19–27.
- [32] R. Tandon, Q. Lei, A. G. Dimakis, and N. Karampatzidakis (2017). Gradient coding: Avoiding stragglers in distributed learning. In *Proceedings of the 34th International Conference on Machine Learning, (PMLR)*, **70**, pp. 3368–3376.
- [33] W. Halbawi, N. Azizan-Ruhi, F. Salehi, and B. Hassibi (2017). Improving distributed gradient descent using reed-solomon codes. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pp. 2027–2031, IEEE.
- [34] N. Raviv, I. Tamo, R. Tandon, and A. G. Dimakis (2018). Gradient coding from cyclic mds codes and expander graphs. In *Proceedings of the 35th International Conference on Machine Learning, Stockholm, Sweden*.
- [35] M. Glasgow, and M. Wootters (2020). Approximate gradient coding with optimal decoding. *IEEE Journal on Selected Areas in Information Theory*, **2**(3), 2021, pp. 855–866.
- [36] K. Yuan, Q. Ling, W. Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*. 2016, **26**(3), pp. 1835–1854.



Elie Atallah received his Ph.D degree in Electrical Engineering from the University of Central Florida in 2019, an M.S degree in Electrical Engineering and an M.S degree in Mathematics both from University of California, Riverside in 2003 and 2006, respectively. At UCF his primary focus was in developing algorithms for distributed optimization, Compressive Sensing and Tensor Decomposition. Previously, he has been an adjunct faculty at several universities such as California Baptist University, Seminole State College, Valencia College and UCF.



Nazanin Rahnavard (S'97-M'10, SM'19) received her Ph.D. in the School of Electrical and Computer Engineering at the Georgia Institute of Technology, Atlanta, in 2007. She is currently an Associate Professor in the Department of Electrical and Computer Engineering at the University of Central Florida, Orlando, Florida. Dr. Rahnavard is the recipient of NSF CAREER award in 2011. She has interest and expertise in a variety of research topics in the communications, networking, and signal processing areas. She serves on the editorial board of the

Elsevier Journal on Computer Networks (COMNET) and on the Technical Program Committee of several prestigious international conferences.



Qiyu Sun received the Ph.D. degree in Mathematics from Hangzhou University, Hangzhou, China, in 1990. He is currently a Professor of Mathematics at the University of Central Florida, Orlando, FL, USA. His research interests include applied and computational harmonic analysis, sampling theory, phase retrieval and graph signal processing. He has published more than 120 papers. He received the 2019 Best SICON Paper Prize, presented by the Society for Industrial and Applied Mathematics (SIAM) Activity Group on Control and Systems

Theory (SIAG/CST).