

Fully asynchronous distributed optimization with linear convergence over directed networks*

SHA Xingyu¹, ZHANG Jiaqi², YOU Keyou¹✉

1. Department of Automation, Beijing National Research Center for Information Science
and Technology, Tsinghua University, Beijing 100084, China

2. Beijing Sankuai Online Technology Co., Ltd., Beijing 100102, China

Abstract: We study distributed optimization problems over a directed network, where nodes aim to minimize the sum of local objective functions via directed communications with neighbors. Many algorithms are designed to solve it for synchronized or randomly activated implementation, which may create deadlocks in practice. In sharp contrast, we propose a fully asynchronous push-pull gradient (APPG) algorithm, where each node updates without waiting for any other node by using possibly delayed information from neighbors. Then, we construct two novel augmented networks to analyze asynchrony and delays, and quantify its convergence rate from the worst-case point of view. Particularly, all nodes of APPG converge to the same optimal solution at a linear rate of $\mathcal{O}(\lambda^k)$ if local functions have Lipschitz-continuous gradients and their sum satisfies the Polyak-Łojasiewicz condition (convexity is not required), where $\lambda \in (0,1)$ is explicitly given and the virtual counter k increases by one when any node updates. Finally, the advantage of APPG over the synchronous counterpart and its linear speedup efficiency are numerically validated via a logistic regression problem.

Key words: fully asynchronous; distributed optimization; linear convergence; Polyak-Łojasiewicz condition

CLC number: O221.2 **Document code:** A **Article ID:** 2097 – 0137 (2023) 05 – 0001 – 23

1 Introduction

As data get larger and more spatially distributed, distributed optimization over a network of computing nodes (aka. agents or workers) has found numerous applications in multi-agent problems (Chang et al., 2015), machine learning (Assran et al., 2019), and resources allocation (Zhang et al., 2020b). It aims to minimize the sum of local objective functions, i. e.,

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} f(\mathbf{x}) := \sum_{i=1}^n f_i(\mathbf{x}), \quad (1)$$

where n is the number of nodes and the local objective function f_i is only known by node i . Nodes are expected to solve (1) by only communicating with neighbors that are defined by the network. For example, f_i often takes the form $f_i(\mathbf{x}) := \sum_{\xi \in \mathcal{D}_i} F_i(\mathbf{x}; \xi)$ in empirical risk minimization problems (Assran et al., 2019), where \mathcal{D}_i is a local data-

* Received: 2023 – 03 – 31

Accepted: 2023 – 04 – 11

Published online: 2023 – 08 – 30

Supported by National Natural Science Foundation of China(62033006, 62203254)

✉ Corresponding author: YOU Keyou(youky@tsinghua.edu.cn)

SHA Xingyu(shaxy18@mails.tsinghua.edu.cn); ZHANG Jiaqi(jqzhang61@gmail.com)

set of node i , \mathbf{x} is the model parameter to be optimized, and $F_i(\mathbf{x}; \xi)$ is the loss of a single sample ξ . See Fig. 1 for an illustration.

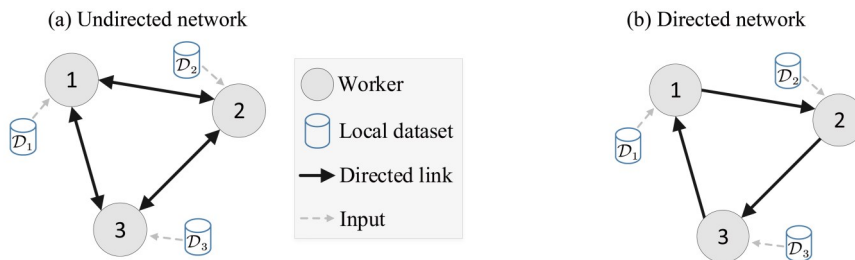


Fig. 1 Undirected and directed peer-to-peer networks

For large-scale optimization problems, it is crucial to design an easily implementable algorithm that is robust to heterogeneous nodes and communication delays. Many existing works focus on synchronous algorithms where all nodes essentially start to compute each iteration simultaneously (c.f. Fig. 2(a)) by using a global synchronization scheme that is often not amenable to the distributed setting.

Asynchronous updates are demonstrated to perform better than the synchronous counterparts (Nedić, 2010; Lian et al., 2018; Wu et al., 2018; Xu et al., 2018; Assran et al., 2020; Zhang et al., 2020a). A popular one is gossip-based (Nedić, 2010; Lian et al., 2018; Xu et al., 2018) where a pair of neighbors is randomly selected to concurrently update via information exchange, see Fig. 2(b). However, this may create deadlocks in practice (Lian et al., 2018), especially for networks with many cycles. It is also sensitive to communication delays and cannot work on directed networks.

To simultaneously address these issues, this work considers the *fully asynchronous* setting (c.f. Fig. 2(c)) (Tian et al., 2020; Zhang et al., 2020a; Assran et al., 2021) over *directed* networks and proposes an asynchronous push-pull gradient (APPG) algorithm to distributedly solve (1). In APPG, a node starts to update without waiting for others with locally accessed (possibly stale) information. It does not need any network synchronization, and can tolerate uneven update frequencies and communication delays among nodes.

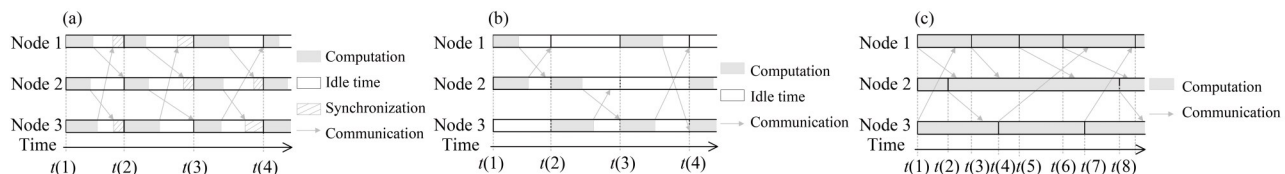


Fig. 2 (a) A synchronous algorithm over directed network; (b) A gossip-based algorithm; (c) A fully asynchronous algorithm

To theoretically evaluate its performance, we develop an augmented network approach and use the machinery of linear matrix inequalities (LMIs) to capture some key quantities via a novel λ -sequence. If all local functions f_i have Lipschitz continuous gradients and the global objective function f satisfies the Polyak-Łojasiewicz (PL) condition (no convexity requirement), we prove from the worst-case point of view that APPG converges linearly to an optimal solution at a rate $\mathcal{O}(\lambda^k)$, where $\lambda \in (0,1)$ is explicitly provided in terms of the asynchrony level and delay bounds, and the virtual counter k increases by one no matter which node updates. Note that the convergence guarantee for gossip-based algorithms is generally given in the stochastic sense.

Then, we implement APPG on a multi-core server with the Message Passing Interface (MPI) to solve a multi-class logistic regression problem over the *Covtype* dataset. The result confirms that its empirical convergence rate in running time is faster than its synchronous counterpart, and achieves linear speedup efficiency with respect to the node number. Its robustness to slow computing nodes is also validated, which is essential in heterogeneous

large-scale networks.

The rest of the paper is organized as follows. Section 2 briefly reviews related works. Section 3 formulates the problem and proposes APPG. We provide the theoretical guarantees for APPG in Section 4. To prove them, we first develop a time-varying augmented network approach in Section 5, and then establish the LMIs in Section 6. In Section 7, we conduct numerical experiments on APPG. Some concluding remarks are drawn in Section 8.

2 Related work

The past decade has witnessed an increasing attention on distributed optimization, especially for the design of synchronous algorithms over undirected networks, e. g. , DGD (Nedić et al. , 2009), DIGing (Nedić et al. , 2017), NIDS (Jakovetić, 2019; Li et al. , 2019; Sun et al. , 2019; Yuan et al. , 2019; Li et al. , 2020; Xin et al. , 2020; Yi et al. , 2020; Xu et al. , 2021). Extensions to directed networks while maintaining linear convergence are challenging due to the network unbalancedness. To resolve it, epigraph reformulation (Xie et al. , 2018) and the push-sum-based method (Nedić et al. , 2015; Assran et al. , 2019; Scutari et al. , 2019) have been proposed. For example, Push-DIGing (Nedić et al. , 2017) combines the push-sum method and gradient tracking (Nedić et al. , 2017; Xu et al. , 2018) to achieve linear convergence for strongly convex and Lipschitz smooth functions over directed networks. However, it involves nonlinear operators per update and may cause numerical issues in implementation, which is resolved in the push-pull/ \mathcal{AB} algorithm (Xin et al. , 2018; Saadatniaki et al. , 2020; Pu et al. , 2021) by simultaneously adopting a row-stochastic matrix and a column-stochastic matrix for information mixing.

In the meanwhile, asynchronous algorithms are proposed to alleviate the synchronization overhead. The gossip-based AsynDGM (Xu et al. , 2018) assumes that each pair of neighbors is activated to update with a fixed probability, and has been extended in Lian et al. (2018) where the activated nodes admit a doubly-stochastic mixing matrix. Nedić (2010) considers that each node is activated by a Poisson process. More general models are studied in Zhao et al. (2015); Wu et al. (2018) by assuming nodes or edges are independently randomly activated with fixed probabilities. It is unclear whether these methods can be applied to directed networks especially when there exist communication delays. Moreover, the coordination between neighbors is indispensable.

Inspired by the seminal works (Tsitsiklis et al. , 1986; Li et al. , 1987), the fully asynchronous setting has emerged as a more scalable and easier implementable alternative. An ADMM-based asynchronous algorithm is studied in Chang et al. (2016) for the master-slave networks. For peer-to-peer networks, an initial attempt is made in Assran et al. (2021) to extend the synchronous gradient-push algorithm (Nedić et al. , 2015) to the fully asynchronous setting. However, they cannot achieve exact convergence to an optimal solution if nodes have different update frequencies. To address it, Zhang et al. (2020a) proposes a novel adaptive mechanism to dynamically adjust stepsize, which is also adopted in Spiridonoff et al. (2020) to analyze stochastic algorithms. Nonetheless, the algorithms in Zhang et al. (2020a); Spiridonoff et al. (2020) only have sublinear convergence rates even if the objective function is strongly convex and Lipschitz smooth. Asy-SONATA (Tian et al. , 2020) exploits the perturbed push-sum method with gradient tracking (Nedić et al. , 2017; Xu et al. , 2018) in the fully asynchronous setting, which converges linearly for strongly convex problems and sublinearly for non-convex problems. The differences between Tian et al. (2020) and our work include: (i) APPG further converges linearly even under the PL condition, which holds for some important non-convex problems, e. g. , the policy optimization for LQR (Fazel et al. , 2018). (ii) APPG uses uncoordinated constant stepsizes while Asy-SONATA can only use uncoordinated diminishing stepsizes, which however leads to a sublinear convergence rate. (iii) We provide explicit convergence rate for APPG in terms of key parameters of the problem, while an explicit rate seems unclear for Asy-SONATA.

Notation: We use the following notations:

(i) $[A]_{ij}$ denotes the element in row i and column j of A . $|\mathcal{A}|$ denotes the cardinality of set \mathcal{A} . $\lfloor x \rfloor$ denotes the largest integer less than or equal to x . \mathbb{R}^n and \mathbb{N} denote the set of n -dimensional real numbers and natural numbers, respectively.

(ii) $\|\cdot\|_2$ denotes the l_2 -norm of a vector or matrix. $\|\cdot\|_F$ denotes the matrix Frobenius norm. $\mathcal{O}(\cdot)$ denotes the big-O notation.

(iii) $\mathbf{1}_n$ and $\mathbf{0}_n$ denote respectively the n -dimensional vector with all ones and all zeros. The subscript may be omitted if the dimension is clear from the context.

(iv) $\nabla f(x)$ denotes the gradient of a function f at x .

(v) \mathbf{a} is called a stochastic vector if it is nonnegative and $\mathbf{a}^T \mathbf{1} = 1$. A is called a row-stochastic matrix if A is nonnegative and $A \mathbf{1} = \mathbf{1}$. A is column-stochastic if A^T is row-stochastic.

(vi) $[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$ and $[\mathbf{a}_1; \mathbf{a}_2; \dots; \mathbf{a}_n]$ denote the horizontal stack and vertical stack of $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$, respectively.

(vii) $\Pi_{\mathcal{X}}(\mathbf{x})$ denotes the projection of \mathbf{x} onto the set \mathcal{X} .

3 Problem formulation and the APPG

3.1 Problem formulation

We aim to solve (1) over a directed network (digraph), which is denoted by $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, 2, \dots, n\}$ is the set of nodes and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges. The directed edge $(i, j) \in \mathcal{E}$ if node i can directly send information to node j . Let $\mathcal{N}_{\text{in}}^i = \{j | (j, i) \in \mathcal{E}\} \cup \{i\}$ denote the set of in-neighbors of node i and $\mathcal{N}_{\text{out}}^i = \{j | (i, j) \in \mathcal{E}\} \cup \{i\}$ denote the set of out-neighbors of i . A path from node i to node j is a sequence of consecutively directed edges from node i to node j . Then, \mathcal{G} is *strongly connected* if there exists a directed path between any pair of nodes.

For a distributed algorithm over \mathcal{G} , each node i has a local state vector \mathbf{x}_i and iteratively update it via directed communications with neighbors, the objective of which is to ensure that all local states $\mathbf{x}_i, i \in \mathcal{V}$ converge to an optimal solution of (1). In this work, we make the following assumptions.

Assumption 1

(a) The digraph \mathcal{G} is strongly connected.

(b) All local functions f_i are β -Lipschitz smooth, i.e., there exists a $\beta > 0$ such that

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall i \in \mathcal{V}, \mathbf{x}, \mathbf{y} \in \mathbb{R}^m.$$

(c) The global objective function f has at least one minimizer and satisfies the Polyak-Łojasiewicz condition (Karimi et al., 2016) with parameter $\alpha > 0$, i.e., $\mathcal{X}^* = \underset{\mathbf{x} \in \mathbb{R}^m}{\operatorname{argmin}} f(\mathbf{x}) \neq \emptyset$, and

$$2\alpha(f(\mathbf{x}) - f^*) \leq \|\nabla f(\mathbf{x})\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^m,$$

where $f^* = \min_{\mathbf{x} \in \mathbb{R}^m} f(\mathbf{x})$.

Assumptions 1(a)-(b) are standard in the distributed smooth optimization over directed networks (Nedić et al., 2017; Pu et al., 2021). The PL condition in Assumption 1(c) is satisfied in some important *non-convex* problems such as the policy optimization for LQR (Fazel et al., 2018). It is strictly weaker than the strongly convex condition that is commonly used to derive the linear convergence of gradient-based methods (Nedić et al., 2017; Tian et al., 2020). Particularly, the strong convexity implies the uniqueness of the minimizer, which is clearly not the case for the PL condition.

3.2 The APPG

We propose APPG to solve (1), which is given in Algorithm 1 from the view of a single node. Note that we

do not introduce any iteration index to emphasize the fact of fully asynchronous implementation.

Algorithm 1 The APPG — from the view of node i

(i) **Initialization:** Each node i selects local stepsize γ_i , initializes \mathbf{x}_i as an arbitrary vector in \mathbb{R}^m , computes $\mathbf{g}_i = \mathbf{y}_i = \nabla f_i(\mathbf{x}_i)$, and creates local buffers \mathcal{X}_i and \mathcal{Y}_i . Then it broadcasts $\tilde{\mathbf{x}}_i := \mathbf{x}_i$ and $\tilde{\mathbf{y}}_i := \mathbf{y}_i / |\mathcal{N}_{\text{out}}^i|$ to its out-neighbors.

(ii) **Repeat**

a) Keep receiving $\tilde{\mathbf{x}}_j$ and $\tilde{\mathbf{y}}_j$ from in-neighbors of node i and copy to \mathcal{X}_i and \mathcal{Y}_i respectively^①, until node i is activated to update.

b) Update \mathbf{x}_i and \mathbf{y}_i as

$$\begin{cases} \mathbf{x}_i \leftarrow \text{avg}(\mathcal{X}_i), \\ \mathbf{g}_i^- \leftarrow \mathbf{g}_i, \quad \mathbf{g}_i \leftarrow \nabla f_i(\mathbf{x}_i), \\ \mathbf{y}_i \leftarrow \text{sum}(\mathcal{Y}_i) + \mathbf{g}_i - \mathbf{g}_i^-, \\ \tilde{\mathbf{x}}_i \leftarrow \mathbf{x}_i - \gamma_i \mathbf{y}_i, \end{cases} \quad (2)$$

where $\text{avg}(\mathcal{X}_i)$ returns the average^② of vectors in \mathcal{X}_i , $\text{sum}(\mathcal{Y}_i)$ takes the sum of vectors in \mathcal{Y}_i .

c) Broadcast $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{y}}_i := \mathbf{y}_i / |\mathcal{N}_{\text{out}}^i|$ to all out-neighbors of i , after which *empty* both \mathcal{X}_i and \mathcal{Y}_i .

(iii) **Until** a stopping criterion is satisfied. e. g. , node i stops if $\|\mathbf{y}_i\|_2 < \varepsilon$ for some predefined $\varepsilon > 0$.

(iv) **Return** \mathbf{x}_i .

Our novel idea lies in the use of local buffers to achieve asynchronous updates. Particularly, a node just keeps receiving messages from its in-neighbors and storing them to its local buffers until it is activated to compute a new update via (2), which involves only summation and (weighted) average of data in the buffers, and a local gradient computation. Clearly, each buffer may contain zero, one or multiple receptions from the same in-neighbor, which is unavoidable in the fully asynchronous setting. Another feature of APPG is that the node computes a new update by using all messages in the buffers instead of only using the latest reception (Tian et al. , 2020). This enables APPG to be robust to bounded communication delays, out of sequence issues and there is no need to use time-stamps (Tian et al. , 2020). Note that APPG can also use only the latest information by carefully assigning weights to receptions. In addition, the local buffers are only conceptually needed and can be waived in implementation since the average and summation operators in (2) can be recursively computed. After (2), the node broadcasts the updated vectors to its out-neighbors and discards all the used messages. Such a process is repeated until a local stopping criterion is satisfied.

Different from synchronous or gossip-based algorithms, APPG does not require any global clock or coordination among neighbors, and each node does not wait for others for new updates. For example, a node can simply start to compute a new update once it completes the current one or receives a message. Thus, there is no deadlock problem. Moreover, nodes do not need to use the same stepsize.

3.3 The idea of APPG

This subsection aims to intuitively explain the linear convergence of the APPG. To this end, we first show how APPG works if nodes are forced to update synchronously. In this case, we can use a global iteration index k to record the update progress.

^① In-neighbors include i itself, i. e. , $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{y}}_i$ are copied to \mathcal{X}_i and \mathcal{Y}_i ;

^② The average can be weighted, e. g. , one may assign higher weights to more recent messages to potentially improve the convergence rate in practice.

Let $X(k)$, $Y(k)$ and $\nabla f(X(k))$ be the stacked local states and gradients at the k -th iteration, i. e. ,

$$\begin{aligned} X(k) &= [\mathbf{x}_1(k), \mathbf{x}_2(k), \dots, \mathbf{x}_n(k)]^\top \in \mathbb{R}^{n \times m}, \\ Y(k) &= [\mathbf{y}_1(k), \mathbf{y}_2(k), \dots, \mathbf{y}_n(k)]^\top \in \mathbb{R}^{n \times m}, \\ \nabla f(X(k)) &= [\nabla f_1(\mathbf{x}_1(k)), \dots, \nabla f_n(\mathbf{x}_n(k))]^\top \in \mathbb{R}^{n \times m}, \end{aligned}$$

and $\mathbf{F} = \text{diag}(\gamma_1, \dots, \gamma_n)$. Then, we obtain that

$$X(k+1) = A(X(k) - \mathbf{F}Y(k)), \quad (3)$$

$$Y(k+1) = BY(k) + \nabla f(X(k+1)) - \nabla f(X(k)) \quad (4)$$

where the row-stochastic matrix A and column-stochastic matrix B result from $\text{avg}(\cdot)$ and $\text{sum}(\cdot)$ in (2), respectively. Clearly, (3)-(4) reduces to the algorithm(Xin et al. , 2018; Saadatniaki et al. ,2020;Pu et al. ,2021), which has been proved to converge linearly for strongly convex and Lipschitz smooth functions.

The key to the linear convergence of (3)-(4) is the introduction of \mathbf{y}_i to distributedly track the gradient of f . To see it, we left multiply (4) with $\mathbf{1}^\top$, use the column-stochasticity of B and notice $Y(0) = \nabla f(X(0))$. This implies that

$$\mathbf{1}^\top Y(k) = \mathbf{1}^\top \nabla f(X(k)). \quad (5)$$

Now suppose that $X(k)$ and $Y(k)$ have already converged to X^∞ and Y^∞ , respectively. It follows from (4) that $Y^\infty = BY^\infty$. Jointly with (5), we obtain that

$$\mathbf{y}_i^\infty = \pi_i^B (\mathbf{1}^\top \nabla f(X^\infty))^\top$$

where π^B is the Perron vector of B , i. e. , $B\pi^B = \pi^B$. Moreover, the row-stochasticity of A implies that $X^\infty = \mathbf{1}(\mathbf{x}^\infty)^\top$ and

$$\left(\mathbf{1}^\top \nabla f(X^\infty) \right)^\top = \left(\mathbf{1}^\top \nabla f(\mathbf{1}(\mathbf{x}^\infty)^\top) \right)^\top = \nabla f(\mathbf{x}^\infty).$$

That is, $\mathbf{y}_i^\infty = \pi_i^B \nabla f(\mathbf{x}^\infty)$. Substituting X^∞ and Y^∞ into (3) implies

$$X(k+1) = A\left(X^\infty - \mathbf{F}\pi^B \nabla f(\mathbf{x}^\infty)^\top\right). \quad (6)$$

Let π^A be Perron vector of A^\top . We left multiply (6) with $(\pi^A)^\top$ and notice that $X^\infty = \mathbf{1}(\mathbf{x}^\infty)^\top$. Then,

$$\bar{\mathbf{x}}(k+1) = \mathbf{x}^\infty - (\pi^A)^\top \mathbf{F}\pi^B \nabla f(\mathbf{x}^\infty) = \mathbf{x}^\infty - \rho \nabla f(\mathbf{x}^\infty), \quad (7)$$

where $\rho = (\pi^A)^\top \mathbf{F}\pi^B$ and $\bar{\mathbf{x}}(k+1) = (\pi^A)^\top X(k+1)$. Clearly, (7) is a gradient descent update, which converges linearly under Assumption 1. It also shows that the limiting point \mathbf{x}^∞ must be an optimal point \mathbf{x}^* and \mathbf{y}_i converges to $\mathbf{y}_i^\infty = \pi_i^B \nabla f(\mathbf{x}^\infty) = \pi_i^B \nabla f(\mathbf{x}^*) = 0$. Moreover, the smaller the γ_i , the closer \mathbf{x}_i to an optimal solution. Therefore, \mathbf{y}_i can serve as a stopping criterion in Algorithm 1.

In the fully asynchronous setting, \mathbf{y}_i plays a similar role. However, a theoretical understanding of APPG is much more complicated. The information delays make the key relation (5) invalid and the uncoordinated updates degrade the tracking performance of \mathbf{y}_i . In essence, APPG is a multi-timescale decision-making problem. To resolve it, we develop an augmented network approach to prove its linear convergence by associating each node with some virtual nodes, under which the asynchronous updates and communications of nodes are transformed into synchronous operations over the augmented network. Moreover, the key technique(Xin et al. , 2018; Saadatniaki et al. , 2020;Pu et al. , 2021) cannot be applied since the transformed system is time-varying and lacks important properties(e. g. the irreducibility of the weighting matrix). To this end, we further introduce an absolute probability sequence and a λ -sequence for the proof.

4 Linear convergence of APPG

Two assumptions on the asynchrony and communication delays are needed for the convergence of APPG.

Assumption 2(Bounded activation time interval) Let t_i and t_i^+ be any two consecutive activation time of node i . There exist two positive constants $\underline{\tau}$ and $\bar{\tau}$ such that $0 < \underline{\tau} \leq |t_i^+ - t_i| \leq \bar{\tau} < +\infty$ for all $i \in \mathcal{V}$.

Assumption 2 is easily satisfied and desirable in practice. In fact, both the lower and upper bounds exist naturally since computing update consumes time and can be finished in finite time. If violated, e. g. , some node is broken, then the information from this node can no longer be accessed, and hence it is impossible to find an optimal solution of (1).

Assumption 3(Bounded transmission delays) For any $(i, j) \in \mathcal{E}$, the transmission delay from node i to node j is bounded by a constant $\tau > 0$.

Note that transmission delays can be time-varying, and the parameters $\underline{\tau}$, $\bar{\tau}$ and τ are not needed for implementing APPG.

Let $\mathcal{T} = \{t(k)\}_{k \geq 1}$ be an increasing sequence of updating time of all nodes, i. e. , $t \in \mathcal{T}$ if some node starts to update at time t . Denote the state of node i just before time $t(k)$ by $x_i(k)$ and $y_i(k)$. Then, the following lemma introduces an important quantity b to characterize the information delay in terms of \mathcal{T} .

Lemma 1 The following statements hold.

(i) Under Assumption 2, let $b_1 = (n - 1)\lceil \bar{\tau}/\underline{\tau} \rceil + 1$. Each node is activated at least once within the time interval $(t(k), t(k + b_1)]$.

(ii) Under Assumptions 2-3, let $b_2 = n\lceil \tau/\underline{\tau} \rceil$ and $b = b_1 + b_2$. The information sent from node i at time $t(k)$ can be received by node j before time $t(k + b_2)$ and used to compute an update before time $t(k + b)$ for any k and $(i, j) \in \mathcal{E}$.

Proof (i) Suppose that node i is not activated during the time interval $(t(p), t(q)]$, $p, q \in \mathbb{N}$ but is activated at $t(q + 1)$. It follows from Assumption 2 that $t(q) - t(p) \leq \bar{\tau}$. Moreover, any other node can be activated at most $\lceil (t(q) - t(p))/\underline{\tau} \rceil \leq \lceil \bar{\tau}/\underline{\tau} \rceil$ times during the time interval $(t(p), t(q)]$, which implies $q - p \leq (n - 1)\lceil \bar{\tau}/\underline{\tau} \rceil$. Hence, the first part of the result follows.

(ii) Suppose that node i sends information at time $t(p)$, $p \in \mathbb{N}$ and node j receives it in the time interval $(t(q), t(q + 1)]$, $q \in \mathbb{N}$. It follows from Assumption 3 that $t(q) - t(p) \leq \tau$. Moreover, Assumption 2 implies that any node can be activated at most $\lceil \tau/\underline{\tau} \rceil$ times during the time interval $[t(p), t(q)]$, i. e. , $q - p + 1 \leq n\lceil \tau/\underline{\tau} \rceil$, and hence $q + 1 \leq p + n\lceil \tau/\underline{\tau} \rceil$. The result follows by letting $p = k$. Jointly with Lemma 1(i), the rest of proof follows immediately.

The quantity b measures the largest possible information staleness caused by asynchrony and delays. For synchronous algorithms without communication delays, we have $b = 1$. For cyclic updating modes where node i updates at $t(nl + i)$, $\forall l \in \mathbb{N}$, we can set $b = n$. Note that a larger b does not imply that the algorithm updates slower in wall-clock time. For example, b is small in synchronous algorithms but $t(k + 1) - t(k)$ may be large due to idle time, while in fully asynchronous algorithms b is often large but $t(k + b) - t(k)$ can be small. Our main theoretical result is given below.

Theorem 1 Suppose that Assumptions 1-3 hold. Let $\bar{\gamma} = \max_i \gamma_i$ and $\underline{\gamma} = \min_i \gamma_i$. If $\underline{\gamma} > 0$ and $\bar{\gamma} <$

$\sqrt{\underline{\gamma}} \left(\frac{\alpha \log_2(1 - \theta)^{-1}}{16n^{b+2}(b+2)^2 \beta} \right)^4$, then there exists an optimal solution $\mathbf{x}^* \in \mathcal{X}^*$ such that

$$\|\mathbf{x}_i(k) - \mathbf{x}^*\|_2 = \mathcal{O}(\lambda^k), \quad \|\mathbf{y}_i(k)\|_2 = \mathcal{O}(\lambda^k), \quad \forall i \in \mathcal{V} \quad (8)$$

where

$$\max \left\{ \left(1 - \theta\right)^{\frac{1}{n(b+3)^2 \log_2 n}}, 1 - \frac{1}{8} \underline{\gamma} \alpha n^{-b} \right\} \leq \lambda < 1, \quad (9)$$

α, β are given in Assumption 1, b is defined in Lemma 1, and $\theta = (nb)^{-nb}$.

Theorem 1 shows that the local decision vector of each node in APPG converges to the same optimal solution at a linear rate, which is characterized in terms of key parameters n, b, α, β and γ_i . The quantity θ in Theorem 1 reflects the speed of achieving consensus and depends on the network topology and asynchrony. Let $v(k) := \max_i x_i(k) - \min_i x_i(k)$ be the maximum difference among nodes at time $t(k)$, then we can show that $v(k + nb) \leq (1 - \theta)v(k)$ when applying APPG on the consensus problem with $f_i = 0, \forall i$. Clearly, the convergence rate of APPG cannot exceed the consensus speed, which is reflected in the first term of the left-hand-side of (9).

To obtain a deterministic convergence rate, Theorem 1 adopts the *worst-case* point of view under the setting that: (a) The underlying network is heavily unbalanced, i. e., the differences between the numbers of in-neighbors and out-neighbors of a node is large. (b) Some nodes compute much faster than the others (match the lower and upper bounds of Assumption 2, respectively). (c) Communication delays are based on the upper bound in Assumption 3. The values of b and θ in Lemma 1 and Theorem 1 are given in this case. Thus, the theoretical rate is expected to be very conservative, but the practical performance is empirically much better. For the synchronous delay-free case, we have $b = 1$ and $\theta = n^{-n}$, and the convergence rate in Theorem 1 reduces to

$\mathcal{O} \left(\max \left\{ \left(1 - n^{-n}\right), 1 - \frac{\alpha \gamma}{8n} \right\}^k \right)$, where the first term in $\max(\cdot)$ matches the consensus result in Olshevsky et al.

(2011, Theorem 8. 1), and the second term is close to the push-pull method (Pu et al., 2021, Remark 6).

Note that the virtual counter k in (8) increases by one no matter which node updates, and hence more nodes generally lead to faster increase of k . To some extent, this suggests a linear speedup efficiency of APPG, which is confirmed via experiments in Section 7.

5 The time-varying augmented system

This section develops our time-varying augmented network approach for the convergence analysis.

5.1 Construction of the augmented digraph

Let $\mathcal{T}_i \subseteq \mathcal{T}$ be the sequence of activation time of node i , i. e., $t \in \mathcal{T}_i$ if node i computes an update at time t . Then, it is clear that

$$\left[\mathbf{x}_i(k+1), \mathbf{y}_i(k+1), \mathbf{g}_i(k+1), \mathbf{g}_i^-(k+1) \right] = \left[\mathbf{x}_i(k), \mathbf{y}_i(k), \mathbf{g}_i(k), \mathbf{g}_i^-(k) \right], \quad \forall t(k) \notin \mathcal{T}_i.$$

To handle bounded time-varying transmission delays and asynchrony, we associate each node i with two types of virtual nodes, and each type has b virtual nodes, where b is given in Lemma 1(ii). We denote the above two types of virtual nodes by $\{v_{x,i}^{(1)}, \dots, v_{x,i}^{(b)}\}$ and $\{v_{y,i}^{(1)}, \dots, v_{y,i}^{(b)}\}$, respectively. We call the first type virtual nodes x -type nodes, which is to deal with the staleness of the state $\tilde{\mathbf{x}}_i, i \in \mathcal{V}$. The second type with subscript y is called y -type nodes, which is to handle the staleness of $\tilde{\mathbf{y}}_i, i \in \mathcal{V}$. Then, we construct an augmented time-varying digraph $\tilde{\mathcal{G}}(k) = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}}(k))$ to represent the communication topology of all these nodes at time $t(k)$, where $\tilde{\mathcal{V}}$ contains $n(2b + 1)$ nodes, including n nodes of \mathcal{G} and $2nb$ virtual nodes.

The edge set $\tilde{\mathcal{E}}(k)$ is described as follows. We first note that there is no edge between any x -type node and any y -type node. For the x -type virtual nodes, the edges $(i, v_{x,i}^{(1)}), (v_{x,i}^{(1)}, v_{x,i}^{(2)}), \dots, (v_{x,i}^{(b-2)}, v_{x,i}^{(b-1)})$ and $(v_{x,i}^{(b-1)}, v_{x,i}^{(b)})$ always exist for all $k \in \mathbb{N}$ and $i \in \mathcal{V}$, see Fig. 3. If $(i, j) \in \mathcal{E}$ in \mathcal{G} and node j receives $\tilde{\mathbf{x}}_i$ at time $t(k)$, then some of the edges $(v_{x,i}^{(1)}, j), (v_{x,i}^{(2)}, j), \dots, (v_{x,i}^{(b)}, j)$ and (i, j) are included in $\mathcal{E}(k)$ (c.f. Fig. 4(a)), depending on the transmission delay of the received message. For example, if node j receives $\tilde{\mathbf{x}}_i(t-u)$ and $\tilde{\mathbf{x}}_i(t-v)$ at time $t(k)$ for some $u, v > 1$, then $(v_{x,i}^{(u-1)}, j), (v_{x,i}^{(v-1)}, j) \in \tilde{\mathcal{E}}(k)$. If $u = 1$, which means that there is no communication delay, then $(i, j) \in \tilde{\mathcal{E}}(k)$.

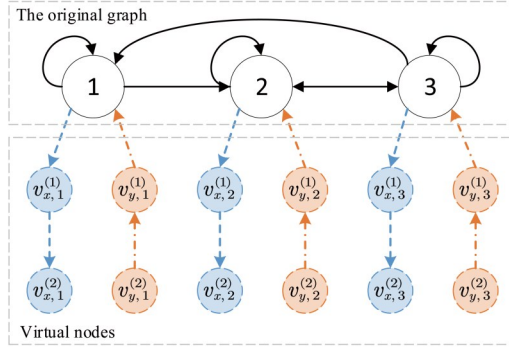


Fig. 3 An augmented graph with virtual nodes to address delays of the original graph

Fig. 4(a) illustrates such an augmented graph^① at some time $t(k)$, where node 1 uses the 2-steps delayed information $\tilde{x}_3(k-3)$ and the latest information $\tilde{x}_1(k-1)$ to compute $x_1(k)$, and hence $(v_{x,3}^{(2)}, 1) \in \tilde{\mathcal{E}}(k)$. Node 2 uses $\tilde{x}_2(k-1)$ and the 1-step delayed information $\tilde{x}_1(k-2)$ and $\tilde{x}_3(k-2)$ to compute $x_2(k)$. Node 3 uses the latest information $\tilde{x}_2(k-1)$ and $\tilde{x}_3(k-1)$ to compute $x_3(k)$. The corresponding row-stochastic matrix $\tilde{A}(k)$ is as demonstrated in Fig.4(b).

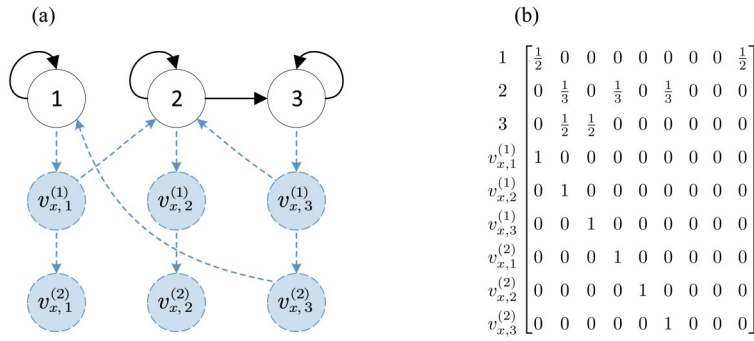


Fig. 4 (a) The topology of the x -type virtual nodes in the augmented graph at some time $t(k)$;
(b) The corresponding row-stochastic matrix

The topology of the y -type virtual nodes is similarly developed with reversed edge directions (c.f. Fig. 5(a)), which is the main motivation of using two types of virtual nodes. Firstly, edges $(v_{y,j}^{(1)}, j)$, $(v_{y,j}^{(2)}, v_{y,j}^{(1)})$, \dots , and $(v_{y,j}^{(b)}, v_{y,j}^{(b-1)})$ are always included in $\tilde{\mathcal{E}}(k)$ with the reversed directions of x -type nodes, see Fig. 3. Secondly, if $(i, j) \in \mathcal{E}$ in \mathcal{G} and $k \in \mathcal{T}_i$, then only one edge in edges $(i, v_{y,j}^{(1)})$, $(i, v_{y,j}^{(2)})$, \dots , $(i, v_{y,j}^{(nb)})$ and (i, j) is included in $\mathcal{E}(k)$, which also depends on the transmission delay of \tilde{y}_i sent from node i to node j . At time $t(k)$, suppose that node i sends $y_i(k)$ to node j , which is received at $t(k+u)$ for $u > 1$, then $(i, v_{y,j}^{(u-1)}) \in \tilde{\mathcal{E}}(k)$ and the delay is $u-1$. Similarly, if there is no communication delay, i. e., $u=1$, then $(i, j) \in \tilde{\mathcal{E}}(k)$. Fig. 5(a) illustrates such an augmented graph at some time $t(k)$, which represents that node 1 sends $\tilde{y}_1(k)$ to node 2 and node 2 use it at $t(k+3)$ to compute $y_2(k+3)$. Other edges can be interpreted similarly. The corresponding column-stochastic matrix $\tilde{B}(k)$ is as

^① The idea of adding virtual nodes to address asynchrony or delays was firstly adopted in Nedić et al.(2010) to study consensus problems and in Zhang et al.(2020a); Tian et al.(2020); Assran et al.(2021) for distributed optimization. Nevertheless, Nedić et al.(2010); Zhang et al.(2020a); Assran et al.(2021) use only x -type virtual nodes and involve division operators. We further introduce y -type nodes to accommodate linear update rules. Tian et al.(2020) associates virtual nodes to original edges rather than original nodes. Since there are generally more edges than nodes in a strongly connected graph and the theoretical convergence rate becomes slower as the number of virtual nodes increases, our analysis may potentially be sharper than Tian et al.(2020).

depicted in Fig. 5(b).

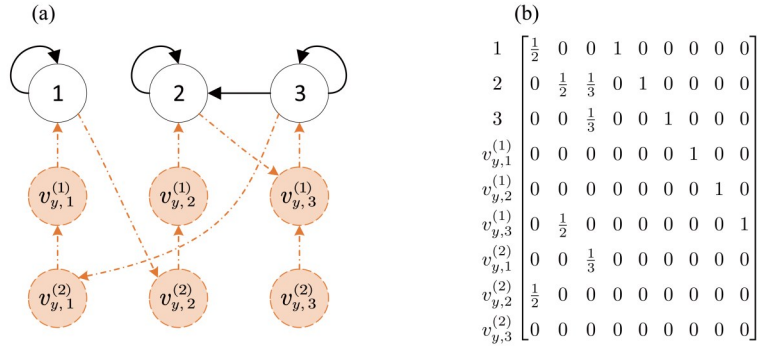


Fig. 5 (a) The topology of the y -type virtual nodes in the augmented graph at some time $t(k)$;
(b) The corresponding column-stochastic matrix

We provide a simple example to visualize the augmented graph approach. Consider that node i sends $\tilde{x}_i(k)$ and $\tilde{y}_i(k)$ to node j at time $t(k)$, and node j receives it at time $t(k+2)$, i.e., the delay is 1. In the augmented graph, this can be viewed as node i directly sends $\tilde{x}_i(k)$ to the virtual node $v_{x,i}^{(1)}$, and sends $\tilde{y}_i(k)$ to the virtual node $v_{y,j}^{(1)}$ at time $t(k)$. Nodes $v_{x,i}^{(1)}$ and $v_{y,j}^{(1)}$ respectively receive $\tilde{x}_i(k)$ and $\tilde{y}_i(k)$ at time $t(k+1)$, and immediately send them to node j at time $t(k+1)$. Finally, node j receives $\tilde{x}_i(k)$ and $\tilde{y}_i(k)$ at time $t(k+2)$. Clearly, all non-virtual nodes in $\tilde{\mathcal{G}}$ receive the same information as that in \mathcal{G} and hence their updates appear to be synchronous and delay-free.

Under the time-varying augmented digraph, we are able to rewrite the APPG in a compact form.

5.2 A compact form of the APPG over the augmented digraph

Let $\mathbf{x}_i^{(u)}(k)$ and $\mathbf{y}_i^{(u)}(k)$ denote the states of virtual node $v_{x,i}^{(u)}$ and $v_{y,i}^{(u)}$ just after time $t(k)$, and $\tilde{n} = bn$. Then, APPG can be rewritten in a compact form over $\tilde{\mathcal{G}}(k)$,

$$\begin{cases} \tilde{\mathbf{X}}(k+1) = \tilde{\mathbf{A}}(k)(\tilde{\mathbf{X}}(k) - \tilde{\mathbf{I}}\mathbf{I}_k^a\tilde{\mathbf{Y}}(k)), \\ \tilde{\mathbf{Y}}(k+1) = \tilde{\mathbf{B}}(k)\tilde{\mathbf{Y}}(k) + \mathbf{I}_k^a(\nabla(k+1) - \nabla(k)) = \tilde{\mathbf{B}}(k)\tilde{\mathbf{Y}}(k) + \nabla(k+1) - \nabla(k), \end{cases} \quad (10)$$

where $\tilde{\mathbf{I}} = \mathbf{I} \otimes \mathbf{I} \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$,

$$\begin{cases} \tilde{\mathbf{X}}(k) = [\mathbf{X}(k); \mathbf{X}^{(1)}(k); \cdots; \mathbf{X}^{(b-1)}(k)] \in \mathbb{R}^{\tilde{n} \times m}, \\ \mathbf{X}^{(u)}(k) = [\mathbf{x}_1^{(u)}(k), \cdots, \mathbf{x}_n^{(u)}(k)]^T, \\ \tilde{\mathbf{Y}}(k) = [\mathbf{Y}(k); \mathbf{Y}^{(1)}(k); \cdots; \mathbf{Y}^{(b-1)}(k)] \in \mathbb{R}^{\tilde{n} \times m}, \\ \mathbf{Y}^{(u)}(k) = [\mathbf{y}_1^{(u)}(k), \cdots, \mathbf{y}_n^{(u)}(k)]^T, \\ \nabla(k) = [\nabla \mathbf{f}(\mathbf{X}(k)); \mathbf{0}_{(\tilde{n}-n) \times m}], \end{cases} \quad (11)$$

the initial condition is $\tilde{\mathbf{X}}(0) = [\mathbf{X}(0); \mathbf{0}_{(\tilde{n}-n) \times m}]$, $\tilde{\mathbf{Y}}(0) = \nabla(0)$, and $\tilde{\mathbf{A}}(k), \tilde{\mathbf{B}}(k), \mathbf{I}_k^a \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$ are

$$[\tilde{\mathbf{A}}(k)]_{ij} = \begin{cases} \frac{1}{|\mathcal{X}_i(k)|}, & \text{if } i, v \in \mathcal{V}, j = nu + v, t(k+1) \in \mathcal{T}_i, \\ & \text{and node } i \text{ receives } \mathbf{x}_v(k-u) \text{ at } t(k+1), \\ 1, & \text{if } i \in \mathcal{V}, t(k+1) \notin \mathcal{T}_i \text{ and } j = i, \\ 1, & \text{if } i \notin \mathcal{V} \text{ and } j = i - n, \\ 0, & \text{otherwise,} \end{cases} \quad (12)$$

$$\begin{aligned}
 [\tilde{\mathbf{B}}(k)]_{ji} &= \begin{cases} \frac{1}{|\mathcal{N}_{\text{out}}^i|}, & \text{if } i, v \in \mathcal{V}, j = nu + v, t(k+1) \in \mathcal{T}_i, \\ & \text{and node } v \text{ receives } \mathbf{x}_i(k) \text{ at } t(k+u), \\ 1, & \text{if } i \in \mathcal{V}, t(k+1) \notin \mathcal{T}_i \text{ and } j = i, \\ 1, & \text{if } i \notin \mathcal{V} \text{ and } j = i - n, \\ 0, & \text{otherwise,} \end{cases} \\
 [\mathbf{I}_k^a]_{ij} &= \begin{cases} 1, & \text{if } i = j, i \in \mathcal{V}, \text{ and } t(k+1) \in \mathcal{T}_i, \\ 0, & \text{otherwise,} \end{cases}
 \end{aligned} \tag{13}$$

where $|\mathcal{X}_i(k)|$ is the number of elements in the buffer \mathcal{X}_i at time $t(k+1)$.

An example of $\tilde{\mathbf{A}}(k)$ and $\tilde{\mathbf{B}}(k)$ is illustrated in Fig. 4(b) and Fig. 5(b), respectively. \mathbf{I}_k^a is a diagonal matrix with the i -th diagonal element be 1 if node i is activated at time $t(k+1)$. The second equality in (10) follows from $\nabla(k+1) = \nabla(k)$ for any $i \in \{i | [\mathbf{I}_k^a]_{ii} = 0\}$. An important fact is that $\tilde{\mathbf{A}}(k)$ is a row-stochastic matrix and $\tilde{\mathbf{B}}(k)$ is a column-stochastic matrix by the use of two types virtual nodes, but they may not be irreducible. Moreover,

$$\mathbf{1}_n^T \tilde{\mathbf{Y}}(k) = \mathbf{1}_n^T \nabla(k) = \mathbf{1}_n^T \nabla \mathbf{f}(\mathbf{X}(k)), \tag{14}$$

which is obtained by left multiplying the second equality of (10) with $\mathbf{1}_n^T$.

Note that (10) generates the same sequence of the states \mathbf{x}_i and \mathbf{y}_i as that of APPG. Hence, it is sufficient to study the convergence of $\tilde{\mathbf{X}}(k)$ and $\tilde{\mathbf{Y}}(k)$ in (10). To this end, we define

$$\begin{cases} \Phi_t^A(k) = \tilde{\mathbf{A}}(k+t-1)\tilde{\mathbf{A}}(k+t-2)\cdots\tilde{\mathbf{A}}(k+1)\tilde{\mathbf{A}}(k), \\ \Phi_t^B(k) = \tilde{\mathbf{B}}(k+t-1)\tilde{\mathbf{B}}(k+t-2)\cdots\tilde{\mathbf{B}}(k+1)\tilde{\mathbf{B}}(k), \end{cases} \tag{15}$$

where $k, t \in \mathbb{N}$, and we adopt the convention that $\Phi_0^A(k) = \Phi_0^B(k) = \mathbf{I}$ and $\Phi_t^A(k) = \Phi_t^B(k) = \mathbf{0}$ for any $k \in \mathbb{N}$ and $t < 0$.

The following lemma states that $\Phi_t^A(k)$ and $\Phi_t^B(k)$ linearly converge to rank-one matrices.

Lemma 2 Under Assumptions 1-3, the following statements are in force.

(i) Let θ be the minimum nonzero element of $\Phi_{nb}^A(k)$ and $\Phi_{nb}^B(k)$. It holds that $\tilde{n}^{-nb} \leq \theta < 1$, where b is defined in Lemma 1(ii).

(ii) There exist stochastic vectors $\phi_t^A(k)$ and $\phi_t^B(k)$ such that

$$\left\| \Phi_t^A(k) - \mathbf{1}\phi_t^A(k)^T \right\|_2 \leq 2\rho', \quad \left\| \Phi_t^B(k) - \phi_t^B(k)\mathbf{1}^T \right\|_2 \leq 2\rho',$$

for all $k, t \in \mathbb{N}$, where $\rho = (1 - \theta)^{1/\tilde{n}} < 1$.

(iii) $\sum_{j=1}^n [\Phi_t^B(0)]_{ij} \geq n^{-b}, \forall i \in \mathcal{V}, t \in \mathbb{N}$.

Proof Part (i): In view of (12) and (13), the minimum nonzero element of $\tilde{\mathbf{A}}(k)$ and $\tilde{\mathbf{B}}(k)$ is greater than $1/\tilde{n}$ for all k . Thus, the minimum nonzero element of $\Phi_t^A(k)$ and $\Phi_t^B(k)$ is greater than \tilde{n}^{-t} for any $t > 0$.

Part (ii): In view of Lemma 1, both $\Phi_b^A(k)$ and $\Phi_b^B(k)$ are indecomposable for all k , and thus $\Phi_{nb}^A(k)$ and $\Phi_{nb}^B(k)$ have positive columns. Then, the proof is similar with that Lemma 5 (Nedić et al., 2010) or Lemma 3 (Nedić et al., 2015), and we omit it for saving space.

Part (iii): We study two cases separately. If $t < b$, then

$$[\Phi_t^B(0)]_{ii} \geq [\tilde{\mathbf{B}}(t-1)]_{ii} [\tilde{\mathbf{B}}(t-2)]_{ii} \cdots [\tilde{\mathbf{B}}(0)]_{ii} \geq (1/n)^b.$$

If $t \geq b$, it follows from a similar argument with the Lemma 2(ii) in Nedić et al. (2010) that $[\Phi_b(k)]_{ij} \geq n^{-b}$ for all $i \in \mathcal{V}$ and $j \in \tilde{\mathcal{V}}$. Then,

$$[\Phi_{b+1}(k-1)]_{ij} = \sum_{u=1}^{\tilde{n}} [\Phi_b(k)]_{iu} [\tilde{B}(k)]_{uj} \geq n^{-b} \sum_{u=1}^{\tilde{n}} [\tilde{B}(k)]_{uj} \geq n^{-b}.$$

where the last inequality follows from the column-stochasticity of $\tilde{B}(k)$. The desired result is obtained by induction.

The following lemma is a direct result of Lemma 2, which specifies μ and \tilde{t} used later.

Lemma 3 Under assumptions of Lemma 2, let $\mu = \frac{1}{4} n^{-(b+1)}$ and $\tilde{t} = \left\lceil nb \log_{1-\theta} \left(\frac{\mu}{2} \right) \right\rceil + 1$. If $\lambda \geq (1 - \theta)^{\frac{1}{n(b+3)^2 \log_2 n}}$, then we have $\lambda^{\tilde{t}} \geq \frac{1}{2}$ and

$$\left\| \Phi_i^A(k) - \mathbf{1} \phi_i^A(k)^T \right\|_2 \leq \mu, \quad \left\| \Phi_i^B(k) - \phi_i^B(k) \mathbf{1}^T \right\|_2 \leq \mu. \quad (16)$$

Proof In view of Lemma 2(ii), it can be readily checked that (16) satisfies. We have

$$\lambda^{\tilde{t}} \geq (1 - \theta)^{\frac{\tilde{t}}{n(b+3)^2 \log_2 n}} \geq \left(\frac{\mu}{2} \right)^{\frac{b+1}{(b+3)^2 \log_2 n}} = n^{-\frac{(b+1)^2}{(b+3)^2 \log_2 n} - \frac{b+1}{(b+3)^2 \log_2 n}} = 2^{-\frac{(b+1)^2}{(b+3)^2} - \frac{3(b+1)}{(b+3)^2 \log_2 n}} = 2^{-\frac{(b+1)^2 \log_2 n + 3(b+1)}{(b+3)^2 \log_2 n}} \geq \frac{1}{2}.$$

Finally, we introduce the *absolute probability sequence* (Touri, 2012).

Lemma 4 (Theorem 4.2 in Touri (2012)) For a sequence of row-stochastic matrices $\{A(k)\}$, there exists a sequence of stochastic vectors $\{\pi(k)\}$ satisfying

$$\pi(k+1)^T A(k) = \pi(k)^T, \quad \forall k \in \mathbb{N}. \quad (17)$$

$\{\pi(k)\}$ is called an absolute probability sequence of $\{A(k)\}$.

In the sequel, we use $\pi(k) \in \mathbb{R}^{\tilde{n}}$ to denote an absolute probability sequence of $\tilde{A}(k)$, which implies that $\pi(k+t)^T \Phi_i^A(k) = \pi(k)^T, \forall k, t \in \mathbb{N}$.

6 Proof of Theorem 1 via LMIs

Under the augmented time-varying digraph, we are ready to prove Theorem 1.

6.1 Outline of the proof

As Nedić et al. (2017), for a nonnegative sequence $\{p(k)\}$, we define that

$$p^{\lambda, k} = \sup_{t \in \mathbb{N}, t \leq k} \frac{p(t)}{\lambda^t}, \quad (18)$$

for $\lambda \in (0, 1)$. We call $\{p^{\lambda, k}\}$ the λ -sequence of $p(k)$. Clearly, if $p^{\lambda, k}$ is uniformly bounded by some constant c , then $p(k) \leq c\lambda^k$ for all k . Our method to prove Theorem 1 is to show the boundedness of $p^{\lambda, k}, k \in \mathbb{N}$ in (18) for some nonnegative sequences $p(k)$. Under Assumptions 1-3 and (10), the proof of Theorem 1 relies on the following four lemmas, whose proofs are given in next subsections.

Lemma 5 Let

$$Q(k) = I_{\tilde{n}} - \mathbf{1}_{\tilde{n}} \pi(k-1)^T, \quad \left\| \tilde{X}(k) \right\|_Q = \left\| Q(k) \tilde{X}(k) \right\|_F,$$

and $\left\| \tilde{X} \right\|_Q^{\lambda, k}$ be the λ -sequence of $\left\| \tilde{X}(k) \right\|_Q$. If $\lambda^{\tilde{t}} \geq 1/2$, then

$$\left\| \tilde{X} \right\|_Q^{\lambda, k} \leq 4\tilde{\gamma}\tilde{n}\tilde{t} \left\| \tilde{Y} \right\|_Q^{\lambda, k} + c_1, \quad (19)$$

where $\left\| \tilde{Y} \right\|_F^{\lambda, k}$ is the λ -sequence of $\left\| \tilde{Y}(k) \right\|_F$, \tilde{t} is defined in Lemma 3, and c_1 is a constant given in (31).

$\left\| \tilde{X}(k) \right\|_Q$ in Lemma 5 is a weighted difference among nodes' states $x_i(k)$, and (19) bounds it by the gradient estimate $\left\| \tilde{Y}(k) \right\|_F$.

Lemma 6 Let

$$v(k+1) = \tilde{B}(k)v(k), \quad v(0) = [\mathbf{1}_n; \mathbf{0}_{\tilde{n}-n}], \quad V(k) = \text{diag}(v(k)), \quad Y_V(k) = V(k)^\dagger \tilde{Y}(k),$$

where $V(k)^\dagger$ is the pseudo inverse of $V(k)$, i. e. ,

$$[\mathbf{V}(k)^\dagger]_{ij} = \begin{cases} 1/[\mathbf{V}(k)]_{ii}, & \text{if } i = j \text{ and } [\mathbf{V}(k)]_{ii} > 0, \\ 0, & \text{otherwise.} \end{cases}$$

Let $\tilde{\mathbf{I}}(k) = \mathbf{V}(k)\mathbf{V}(k)^\dagger$, $\tilde{\mathbf{I}}(k) = \tilde{\mathbf{I}}(k)\mathbf{1}_n$ and

$$\mathbf{S}(k) = \tilde{\mathbf{I}}(k) - \frac{1}{n}\tilde{\mathbf{I}}(k)\mathbf{v}(k)^\top, \quad \|\mathbf{Y}_v(k)\|_s = \|\mathbf{S}(k)\mathbf{Y}_v(k)\|_F. \quad (20)$$

Define the corresponding λ -sequence $\|\mathbf{Y}_v\|_s^{\lambda,k}$ of $\{\|\mathbf{Y}_v(k)\|_s\}$. If $\lambda^i \geq 1/2$, then

$$\|\mathbf{Y}_v\|_s^{\lambda,k} \leq 4\beta n^{b+1}b\tilde{t} \left(2\|\tilde{\mathbf{X}}\|_Q^{\lambda,k} + \bar{\gamma}\|\tilde{\mathbf{Y}}\|_F^{\lambda,k} \right) + c_2, \quad (21)$$

where β is given in Assumption 1, b is defined in Lemma 1(ii), \tilde{t} is defined in Lemma 3, and c_2 is given by (38).

Similarly, $\|\mathbf{Y}_v(k)\|_s$ measures the differences between the weighted gradient estimates of different nodes, which is bounded by $\|\tilde{\mathbf{X}}(k)\|_Q$ and $\|\tilde{\mathbf{Y}}(k)\|_F$.

Lemma 7 Let

$$\mathbf{x}_\pi(k) = \boldsymbol{\pi}(k)^\top \tilde{\mathbf{X}}(k). \quad (22)$$

Define $\tilde{f}^{\lambda,k}$ be the λ -sequence of $\{\sqrt{f(\mathbf{x}_\pi(k)) - f^*}\}$. If $\bar{\gamma} \leq \frac{1}{4nb\beta}$ and $\lambda^b \geq 1 - \frac{1}{8}\alpha\bar{\gamma}n^{-b}$, then

$$\tilde{f}^{\lambda,k} \leq \frac{16n^{\frac{b}{2}}b}{\alpha\sqrt{\bar{\gamma}}} \left(n\beta\|\tilde{\mathbf{X}}\|_Q^{\lambda,k} + \|\mathbf{Y}_v\|_s^{\lambda,k} + n\beta\bar{\gamma}\|\tilde{\mathbf{Y}}\|_F^{\lambda,k} \right) + c_3, \quad (23)$$

where α, β are given in Assumption 1, b is defined in Lemma 1(ii), and c_3 is a constant given in (44).

Intuitively, $\mathbf{x}_\pi(k)$ is a weighted average of $\mathbf{x}_i(k)$, $i \in \mathcal{V}$, and $f(\mathbf{x}_\pi(k)) - f^*$ is the optimality gap. Eq. (22) shows that the square root of the optimality gap can be bounded by $\|\tilde{\mathbf{X}}\|_Q^{\lambda,k}$, $\|\mathbf{Y}_v\|_s^{\lambda,k}$, and $\|\tilde{\mathbf{Y}}\|_F^{\lambda,k}$.

Lemma 8 With the above-defined $\|\tilde{\mathbf{X}}\|_Q^{\lambda,k}$, $\|\mathbf{Y}_v\|_s^{\lambda,k}$ and $\tilde{f}^{\lambda,k}$, it holds that

$$\|\tilde{\mathbf{Y}}\|_F^{\lambda,k} \leq n\sqrt{b}\beta\|\tilde{\mathbf{X}}\|_Q^{\lambda,k} + n\|\mathbf{Y}_v\|_s^{\lambda,k} + \frac{\beta n\sqrt{2b}}{\sqrt{\alpha}}\tilde{f}^{\lambda,k}, \quad (24)$$

where β and b are given in Assumption 1 and Lemma 1, respectively.

Proof of Theorem 1

Note that λ and γ_i , $i \in \mathcal{V}$ defined in Theorem 1 satisfy all the conditions on λ in the above four lemmas and the condition on γ_i in Lemma 7. Thus, (19), (21), (23) and (24) hold. Let

$$\mathbf{e}(k) = \left[\|\tilde{\mathbf{X}}\|_Q^{\lambda,k}, \|\mathbf{Y}_v\|_s^{\lambda,k}, \tilde{f}^{\lambda,k}, \|\tilde{\mathbf{Y}}\|_F^{\lambda,k} \right]^\top, \quad \mathbf{c} = [c_1, c_2, c_3, 0]^\top.$$

Combining (19), (21), (23), and (24), we obtain that for all $k \in \mathbb{N}$,

$$\mathbf{e}(k) \leq \mathbf{M}\mathbf{e}(k) + \mathbf{c}, \quad (25)$$

where \leq is the element-wise inequality and \mathbf{M} is a nonnegative matrix

$$\mathbf{M} = \begin{bmatrix} 0 & 0 & 0 & 4nb\tilde{t}\bar{\gamma} \\ 8\beta n^{b+1}b\tilde{t} & 0 & 0 & 4\beta n^{b+1}b\tilde{t}\bar{\gamma} \\ \frac{16n^{b/2+1}b\beta}{\alpha\sqrt{\bar{\gamma}}} & \frac{16n^{b/2}b}{\alpha\sqrt{\bar{\gamma}}} & 0 & \frac{16n^{b/2+1}b\beta\bar{\gamma}}{\alpha\sqrt{\bar{\gamma}}} \\ n\beta\sqrt{b} & n & \frac{n\beta\sqrt{2b}}{\sqrt{\alpha}} & 0 \end{bmatrix}.$$

It follows from (25) that if the spectral radius $\rho(\mathbf{M})$ of \mathbf{M} is strictly less than 1, then $\|\tilde{\mathbf{X}}\|_Q^{\lambda,k}$, $\|\mathbf{Y}_v\|_s^{\lambda,k}$ and $\tilde{f}^{\lambda,k}$ are all bounded for all $k \in \mathbb{N}$.

Define a transformation matrix $T = \text{diag}\left(\left[1, 1, \sqrt{\underline{\gamma}}, \sqrt{\underline{\gamma}}\right]\right)$. Then, TMT^{-1} can be arbitrarily close to a strictly lower triangular matrix by making $\bar{\gamma}/\sqrt{\underline{\gamma}}$ sufficiently small, which implies that $\rho(M) = \rho(TMT^{-1}) < 1$. To obtain an upper bound of $\bar{\gamma}$, we notice that $\rho(M^4) = \mathcal{O}\left(\bar{\gamma}/\sqrt{\underline{\gamma}}\right)$. By letting $\rho(M^4) < 1$, we have

$$\bar{\gamma} < \frac{\alpha^4 \sqrt{\underline{\gamma}}}{16^4 n^{4(b+1)} b^4 \beta^4 \tilde{t}^2}.$$

The bound in Theorem 1 is obtained by substituting \tilde{t} in Lemma 3.

Define $\bar{c} := \frac{\max\{c_1, c_2, c_3\}}{1 - \rho(M)} < \infty$, where c_1, c_2 and c_3 are given in (31), (38) and (44), respectively. It follows

from (25) that $\|\tilde{X}\|_Q^{\lambda, k} \leq \bar{c}$ and $\tilde{f}^{\lambda, k} \leq \bar{c}, \forall k$. We have

$$\begin{aligned} \|\mathbf{x}_i(k) - \Pi_{\mathcal{X}^*}(\mathbf{x}_\pi(k-1))\|_2 &\leq \|\mathbf{x}_i(k) - \mathbf{x}_\pi(k-1)\|_2 + \|\mathbf{x}_\pi(k-1) - \Pi_{\mathcal{X}^*}(\mathbf{x}_\pi(k-1))\|_2 \\ &\leq \|\tilde{X}(k)\|_Q + \sqrt{2/\alpha} \sqrt{f(\mathbf{x}_\pi(k-1)) - f^*}, \end{aligned}$$

where the last inequality used the equivalence between the Polyak-Łojasiewicz condition and the quadratic growth condition (Karimi et al., 2016, Theorem 2), i. e., $f(\mathbf{x}) - f^* \geq \frac{\alpha}{2} \|\mathbf{x} - \Pi_{\mathcal{X}^*}(\mathbf{x})\|_2^2, \forall \mathbf{x} \in \mathbb{R}^m$. Let $\mathbf{x}^*(k) = \Pi_{\mathcal{X}^*}(\mathbf{x}_\pi(k-1))$ and $\|\mathbf{x}_i - \mathbf{x}^*\|_2^{\lambda, k}$ be the λ -sequence of $\{\|\mathbf{x}_i(k) - \mathbf{x}^*(k)\|_2\}$. We have

$$\|\mathbf{x}_i - \mathbf{x}^*\|_2^{\lambda, k} \leq \|\tilde{X}\|_Q^{\lambda, k} + \sqrt{2/\alpha} \tilde{f}^{\lambda, k} \leq (1 + \sqrt{2/\alpha})\bar{c}, \forall i.$$

Combined with the boundedness of $\|\tilde{Y}\|_F^{\lambda, k}$, the result in Theorem 1 follows by the definition of λ -sequence.

6.2 Two useful propositions

We establish two important results in this subsection. The first one shows a property of λ -sequence and the second one recalls the contraction relation of gradient methods.

Proposition 1 Let $\{p(k)\}, \{q(k)\}$ be nonnegative sequences satisfying

$$p(t+j) \leq rp(t) + \sum_{i=0}^{j-1} q(t+i), \tag{26}$$

where $r \in [0, 1)$. If we choose λ such that $\lambda^j \in (r, 1)$, then the λ -sequences $p^{\lambda, k}$ and $q^{\lambda, k}$ in (18) satisfy

$$p^{\lambda, k} \leq \frac{j}{\lambda^j - r} q^{\lambda, k} + c_\lambda, \forall k \in \mathbb{N},$$

where $c_\lambda = \frac{\lambda^j}{\lambda^j - r} \sum_{t=1}^m \lambda^{-t} p(t)$ is a constant.

It follows from (26) that

$$\lambda^{-(t+j)} p(t+j) \leq \frac{r}{\lambda^j} \lambda^{-t} p(t) + \sum_{i=0}^{j-1} \frac{1}{\lambda^{j-i}} \lambda^{-(t+i)} q(t+i), \forall t = 1, \dots, k.$$

This gives k inequalities by selecting $t = 1, \dots, k$. On the other hand, we have $\lambda^{-t} p(t) \leq \lambda^{-t} p(t)$, which gives another j inequalities by selecting $t = 1, \dots, j$. Take the maximum on both sides of these $k+j$ inequalities and use the definition of λ -sequence, we obtain

$$p^{\lambda, k+j} \leq \frac{r}{\lambda^j} p^{\lambda, k} + \max\left\{q^{\lambda, k+j} \sum_{i=0}^{j-1} \frac{1}{\lambda^{j-i}}, \max_{t=1, \dots, j} \lambda^{-t} p(t)\right\} \leq \frac{r}{\lambda^j} p^{\lambda, k+j} + q^{\lambda, k+j} \sum_{i=0}^{j-1} \frac{1}{\lambda^{j-i}} + \sum_{t=1}^j \lambda^{-t} p(t). \tag{27}$$

If $\lambda^j \in (r, 1)$, then (27) implies

$$\begin{aligned} p^{\lambda, k+j} &\leq \frac{\lambda^j \sum_{i=0}^{j-1} \frac{1}{\lambda^{j-i}}}{\lambda^j - r} q^{\lambda, k+j} + \frac{\lambda^j}{\lambda^j - r} \sum_{t=1}^j \lambda^{-t} p(t) = \frac{1 - \lambda^j}{(\lambda^j - r)(1 - \lambda)} q^{\lambda, k+j} + \frac{\lambda^j}{\lambda^j - r} \sum_{t=1}^j \lambda^{-t} p(t) \\ &\leq \frac{j}{\lambda^j - r} q^{\lambda, k+j} + \frac{\lambda^j}{\lambda^j - r} \sum_{t=1}^j \lambda^{-t} p(t), \end{aligned}$$

for all $k + j \in \mathbb{N}$, where we have used that $1 - \lambda^j \leq j(1 - \lambda)$. The result is obtained immediately.

For any nonnegative sequences $\{p(k)\}$, $\{q(k)\}$, letting $r(k) = p(k) + q(k)$, it holds that $r^{\lambda, k} \leq p^{\lambda, k} + q^{\lambda, k}$, which can be easily checked by definition.

The following proposition shows the convergence rate of a perturbed gradient descent method for minimizing functions satisfying the PL condition. As a special case, it recovers the linear convergence rate of the standard gradient descent method.

Proposition 2 Suppose that f is β -Lipschitz smooth and satisfies the Polyak-Łojasiewicz condition in Assumption 1(iii). Let $\eta \in \left(0, \frac{1}{2\beta}\right)$, $\sigma = 1 - \alpha\eta(1 - 2\eta\beta) < 1$, and $\mathbf{x}^+ = \mathbf{x} - \eta\nabla f(\mathbf{x}) + \boldsymbol{\varepsilon}$. Then

$$f(\mathbf{x}^+) - f^* \leq \sigma(f(\mathbf{x}) - f^*) + (2\eta^{-1} + \beta)\|\boldsymbol{\varepsilon}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

Proof It follows from the β -Lipschitz smoothness that

$$\begin{aligned} f(\mathbf{x}^+) &\leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(-\eta\nabla f(\mathbf{x}) + \boldsymbol{\varepsilon}) + \frac{\beta}{2}\|-\eta\nabla f(\mathbf{x}) + \boldsymbol{\varepsilon}\|_2^2 \\ &\leq f(\mathbf{x}) - \eta\|\nabla f(\mathbf{x})\|_2^2 + \frac{\eta}{2}\|\nabla f(\mathbf{x})\|_2^2 + \frac{2}{\eta}\|\boldsymbol{\varepsilon}\|_2^2 + \beta\eta^2\|\nabla f(\mathbf{x})\|_2^2 + \beta\|\boldsymbol{\varepsilon}\|_2^2 \\ &\leq f(\mathbf{x}) - \eta\left(\frac{1}{2} - \eta\beta\right)\|\nabla f(\mathbf{x})\|_2^2 + \left(\frac{2}{\eta} + \beta\right)\|\boldsymbol{\varepsilon}\|_2^2. \end{aligned}$$

Applying the Polyak-Łojasiewicz inequality on $\|\nabla f(\mathbf{x})\|_2^2$ implies the desired result:

$$\begin{aligned} f(\mathbf{x}^+) - f^* &\leq f(\mathbf{x}) - f^* - 2\alpha\eta\left(\frac{1}{2} - \eta\beta\right)(f(\mathbf{x}) - f^*) + \left(\frac{2}{\eta} + \beta\right)\|\boldsymbol{\varepsilon}\|_2^2 \\ &\leq (1 - \alpha\eta(1 - 2\eta\beta))(f(\mathbf{x}) - f^*) + (2\eta^{-1} + \beta)\|\boldsymbol{\varepsilon}\|_2^2. \end{aligned}$$

We end this subsection with two inequalities which will be frequently used later. For any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$,

$$\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_2\|\mathbf{B}\|_F, \quad \|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_\infty\|\mathbf{A}\|_1},$$

and $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{n}$ for any row-stochastic matrix \mathbf{A} .

6.3 Proof of Lemma 5

Let \tilde{t} be defined in Lemma 3. It follows from (10) and (15) that

$$\begin{aligned} \|\tilde{\mathbf{X}}(k + \tilde{t})\|_Q &= \|\mathbf{Q}(k + \tilde{t})\tilde{\mathbf{X}}(k + \tilde{t})\|_F \\ &\leq \|\mathbf{Q}(k + \tilde{t})\boldsymbol{\Phi}_{\tilde{t}}^A(k)\tilde{\mathbf{X}}(k)\|_F + \sum_{t=0}^{\tilde{t}-1} \|\mathbf{Q}(k + \tilde{t})\boldsymbol{\Phi}_{\tilde{t}-t}^A(k+t)\tilde{\Gamma}I_{k+t}^a\tilde{\mathbf{Y}}(k+t)\|_F \end{aligned} \quad (28)$$

By the stochastic vector $\boldsymbol{\Phi}_{\tilde{t}}^A(k)$ in Lemma 3, (28) implies that

$$\begin{aligned} \|\tilde{\mathbf{X}}(k + \tilde{t})\|_Q &\leq \|\mathbf{Q}(k + \tilde{t})(\boldsymbol{\Phi}_{\tilde{t}}^A(k) - \mathbf{1}\boldsymbol{\Phi}_{\tilde{t}-1}^A(k)^\top)\mathbf{Q}(k)\tilde{\mathbf{X}}(k)\|_F \\ &\quad + \bar{\gamma}\sum_{t=0}^{\tilde{t}-1} \|\mathbf{Q}(k + \tilde{t})\|_2\|\boldsymbol{\Phi}_{\tilde{t}-t}^A(k+t)\|_2\|\mathbf{I}_{k+t}^a\tilde{\mathbf{Y}}(k+t)\|_F \\ &\leq \|\mathbf{Q}(k + \tilde{t})\|_2\|\boldsymbol{\Phi}_{\tilde{t}}^A(k) - \mathbf{1}\boldsymbol{\Phi}_{\tilde{t}-1}^A(k)^\top\|_2\|\mathbf{Q}(k)\tilde{\mathbf{X}}(k)\|_F + \tilde{n}\bar{\gamma}\sum_{t=0}^{\tilde{t}-1} \|\tilde{\mathbf{Y}}(k+t)\|_F \\ &\leq \sqrt{\tilde{n}}\mu\|\tilde{\mathbf{X}}(k)\|_Q + \bar{\gamma}\tilde{n}\sum_{t=0}^{\tilde{t}-1} \|\tilde{\mathbf{Y}}(k+t)\|_F, \end{aligned} \quad (29)$$

where $\mu = \frac{1}{4}n^{-(b+1)} < 1$ is given in Lemma 3, the first inequality used the relation

$$\begin{aligned} \mathbf{Q}(k+t)(\mathbf{A} - \mathbf{1}\mathbf{r}^\top)\mathbf{Q}(k) &= \mathbf{Q}(k+t)\mathbf{A}\mathbf{Q}(k) - \mathbf{Q}(k+t)\mathbf{1}\mathbf{r}^\top\mathbf{Q}(k) \\ &= \mathbf{Q}(k+t)(\mathbf{A} - \mathbf{1}_{\tilde{n}}\boldsymbol{\pi}(k-1)^\top) - \mathbf{Q}(k+t)(\mathbf{1}_{\tilde{n}}\mathbf{r}^\top - \mathbf{1}_{\tilde{n}}\boldsymbol{\pi}(k-1)^\top) \\ &= \mathbf{Q}(k+t)\mathbf{A} - \mathbf{Q}(k+t)\mathbf{1}_{\tilde{n}}\mathbf{r}^\top = \mathbf{Q}(k+t)\mathbf{A}, \end{aligned} \quad (30)$$

for any row-stochastic matrix \mathbf{A} and stochastic vector \mathbf{r} , and the last inequality follows from Lemma 3, $\|\mathbf{Q}(k)\|_2 \leq \sqrt{\tilde{n}}$ and $\|\Phi_i^A(k)\|_2 \leq \sqrt{\tilde{n}}, \forall k$.

Note that $\sqrt{\tilde{n}} \mu \leq 1/4$ and $\lambda^i \geq 1/2$. By Proposition 1 and (29), we obtain that $\|\tilde{\mathbf{X}}\|_Q^{\lambda,k} \leq 4\bar{\gamma}\tilde{n}\tilde{\ell}\|\tilde{\mathbf{Y}}\|_Q^{\lambda,k} + c_1$ where

$$c_1 = 4\bar{\gamma}\tilde{n} \sum_{t=1}^i \lambda^{-t} \|\tilde{\mathbf{X}}(t)\|_Q. \quad (31)$$

6.4 Proof of Lemma 6

Let $\mathbf{v}(k) = [v_1(k), \dots, v_{\tilde{n}}(k)]^T$ and $\mathcal{I}_V(k) = \{i | v_i(k) = [\mathbf{V}(k)]_{ii} = 0\}$. Note that $v_i(k) \geq n^{-b}$ for all $i \in \mathcal{V}, k \in \mathbb{N}$ from Lemma 2(iii). It can be shown that the i -th row of $\tilde{\mathbf{Y}}(k)$ is $\mathbf{0}_m^T$ for all $i \in \mathcal{I}_V(k)$, and thus $\tilde{\mathbf{Y}}(k) = \mathbf{V}(k)\mathbf{Y}_V(k)$.

Let $\mathbf{R}(k) = \nabla(k+1) - \nabla(k)$. It then follows from (10) that

$$\mathbf{Y}_V(k+1) = \tilde{\mathbf{B}}_V(k)\mathbf{Y}_V(k) + \mathbf{V}(k+1)^\dagger \mathbf{R}(k), \quad (32)$$

where $\tilde{\mathbf{B}}_V(k) = \mathbf{V}(k+1)^\dagger \tilde{\mathbf{B}}(k)\mathbf{V}(k)$ and one can prove that each row except the i -th ($i \in \mathcal{I}_V(k)$) row of $\tilde{\mathbf{B}}_V(k)$ has row sum 1 using similar arguments as in Lemma 4 of Nedić et al. (2016;2017). By $\nabla(k)$ in (11) and Assumption 1, we have

$$\|\mathbf{R}(k)\|_F = \|\nabla(k+1) - \nabla(k)\|_F \leq \beta \|\mathbf{X}(k+1) - \mathbf{X}(k)\|_F \leq \beta \|\tilde{\mathbf{X}}(k+1) - \tilde{\mathbf{X}}(k)\|_F. \quad (33)$$

Notice that

$$\begin{aligned} \|\tilde{\mathbf{X}}(k+1) - \tilde{\mathbf{X}}(k)\|_F &= \|\tilde{\mathbf{A}}(k)(\tilde{\mathbf{X}}(k) - \tilde{\mathbf{I}}\mathbf{I}_k^a\tilde{\mathbf{Y}}(k)) - \tilde{\mathbf{X}}(k)\|_F \leq \|(\tilde{\mathbf{A}}(k) - \mathbf{I})\mathbf{Q}(k)\tilde{\mathbf{X}}(k)\|_F + \bar{\gamma}\sqrt{n}\|\tilde{\mathbf{Y}}(k)\|_F \\ &\leq 2\sqrt{n}\|\tilde{\mathbf{X}}(k)\|_Q + \bar{\gamma}\sqrt{n}\|\tilde{\mathbf{Y}}(k)\|_F, \end{aligned} \quad (34)$$

where the first inequality follows from the row-stochasticity of $\tilde{\mathbf{A}}(k)$, and the last inequality is from that $\|\tilde{\mathbf{A}}(k) - \mathbf{I}\|_2 \leq \sqrt{n} + 1 \leq 2\sqrt{n}$. By Combining (33) and (34), we obtain

$$\|\mathbf{R}(k)\|_F \leq 2\beta\sqrt{n}\|\tilde{\mathbf{X}}(k)\|_Q + \beta\bar{\gamma}\sqrt{n}\|\tilde{\mathbf{Y}}(k)\|_F. \quad (35)$$

To analyze the sequence $\{\mathbf{Y}_V(k)\}$, we define

$$\begin{aligned} \tilde{\Phi}_i(k) &:= \tilde{\mathbf{B}}_V(k+t-1)\tilde{\mathbf{B}}_V(k+t-2)\cdots\tilde{\mathbf{B}}_V(k+1)\tilde{\mathbf{B}}_V(k) = \mathbf{V}(k+t)^\dagger \left(\prod_{l=t-1}^k \tilde{\mathbf{B}}(k+l)\tilde{\mathbf{I}}(k+l) \right) \tilde{\mathbf{B}}(k)\mathbf{V}(k) \\ &= \mathbf{V}(k+t)^\dagger \Phi_i^B(k)\mathbf{V}(k). \end{aligned}$$

where

$$[\tilde{\mathbf{I}}(k)]_{ij} = [\mathbf{V}(k)\mathbf{V}(k)^\dagger]_{ij} = \begin{cases} 1, & \text{if } i=j, i \notin \mathcal{I}_V(k), \\ 0, & \text{otherwise.} \end{cases}$$

and the last equality follows from that $\tilde{\mathbf{I}}(k+1)\tilde{\mathbf{B}}(k)\mathbf{V}(k) = \tilde{\mathbf{B}}(k)\mathbf{V}(k), \forall k \in \mathbb{N}$, where we used the fact that $[\tilde{\mathbf{B}}(k)\mathbf{V}(k)\mathbf{1}_{\tilde{n}}] = v_i(k+1) = 0$ for any $i \in \mathcal{I}_V(k+1)$, and thus the i -th row of $\tilde{\mathbf{B}}(k)\mathbf{V}(k)$ is $\mathbf{0}_m^T$.

By Lemma 2 and Lemma 3, it holds that $\Phi_i^B(k)$ can be written as

$$\Phi_i^B(k) = \phi_i^B(k)\mathbf{1}^T + \Delta\Phi_i(k),$$

where $\|\Delta\Phi_i(k)\|_2 \leq 2\rho'$ and $\|\phi_i^B(k)\|_2 \leq \mu$. Hence,

$$\mathbf{v}(k+t) = \Phi_i^B(k)\mathbf{v}(k) = \phi_i^B(k)\mathbf{1}^T\mathbf{v}(k) + \Delta\Phi_i(k)\mathbf{v}(k) = \phi_i^B(k)\mathbf{1}^T\mathbf{v}(0) + \Delta\Phi_i(k)\mathbf{v}(k) = n\phi_i^B(k) + \Delta\Phi_i(k)\mathbf{v}(k),$$

which implies that

$$\boldsymbol{\phi}_i^B(k) = \frac{1}{n}(\mathbf{v}(k+t) - \Delta\boldsymbol{\Phi}_i(k)\mathbf{v}(k)). \quad (36)$$

Let $\tilde{\mathbf{I}}(k) = \tilde{\mathbf{I}}(k)\mathbf{1}$ and $\mathbf{C}_i(k) = \mathbf{V}(k+t)^\dagger \Delta\boldsymbol{\Phi}_i(k) \left(\mathbf{I} - \frac{1}{n}\mathbf{v}(k)\mathbf{1}_n^\top \right) \mathbf{V}(k)$. It follows from (35)-(36) that $\tilde{\boldsymbol{\Phi}}_i(k) = \frac{1}{n}\tilde{\mathbf{I}}(k+t)\mathbf{v}(k)^\top + \mathbf{C}_i(k)$.

Now we turn to the sequence $\{Y_v(k)\}$. It follows from (20) and (32) that

$$\begin{aligned} \|Y_v(k+\tilde{t})\|_S &= \|S(k+\tilde{t})Y_v(k+\tilde{t})\|_F \\ &\leq \|S(k+\tilde{t})\tilde{\boldsymbol{\Phi}}_i(k)Y_v(k)\|_F + \sum_{i=1}^{\tilde{t}} \|S(k+\tilde{t})\tilde{\boldsymbol{\Phi}}_{i-1}(k+t)\mathbf{V}(k+t)^\dagger \mathbf{R}(k+t-1)\|_F \\ &\leq \|S(k+\tilde{t})\tilde{\boldsymbol{\Phi}}_i(k)Y_v(k)\|_F + \sum_{i=1}^{\tilde{t}} \|S(k+\tilde{t})\mathbf{V}(k+\tilde{t})^\dagger \boldsymbol{\Phi}_{i-1}^B(k+t)\mathbf{R}(k+t-1)\|_F \\ &\leq \|S(k+\tilde{t})\tilde{\boldsymbol{\Phi}}_i(k)Y_v(k)\|_F + \sqrt{\tilde{n}}n^b \sum_{i=1}^{\tilde{t}} \|\mathbf{R}(k+t-1)\|_F, \end{aligned} \quad (37)$$

where we used in the last inequality that $\|S(k)\mathbf{V}(k)^\dagger\|_2 = \|\mathbf{V}(k)^\dagger - \frac{1}{n}\tilde{\mathbf{I}}(k)\tilde{\mathbf{I}}(k)^\top\|_2 \leq n^b, \forall k$ from Lemma 2(iii).

Similar to (30), we have $S(k+\tilde{t})\tilde{\boldsymbol{\Phi}}_i(k) = S(k+\tilde{t})\left(\tilde{\boldsymbol{\Phi}}_i(k) - \frac{1}{n}\tilde{\mathbf{I}}(k+t)\mathbf{v}(k)^\top\right)S(k)$, and

$$\begin{aligned} \left\| S(k+\tilde{t})\left(\tilde{\boldsymbol{\Phi}}_i(k) - \frac{1}{n}\tilde{\mathbf{I}}(k+t)\right) \right\|_2 &= \|S(k+\tilde{t})\mathbf{C}_i(k)\|_2 \\ &\leq \left\| \mathbf{V}(k+t)^\dagger - \frac{1}{n}\tilde{\mathbf{I}}^\top \right\|_2 \|\Delta\boldsymbol{\Phi}_i(k)\|_2 \left\| \left(\mathbf{I} - \frac{1}{n}\mathbf{v}(k)\mathbf{1}_n^\top \right) \mathbf{V}(k) \right\|_2 \\ &\leq n^{b+1}\mu = \frac{1}{4}. \end{aligned}$$

Following a similar argument as (29), it follows from (35) and (37) that

$$\|Y_v(k+\tilde{t})\|_S \leq \frac{1}{4}\|Y_v(k)\|_S + \beta n^{b+1}b \sum_{i=0}^{\tilde{t}-1} \left(2\|\tilde{\mathbf{X}}(k+i)\|_Q + \bar{\gamma}\|\tilde{\mathbf{Y}}(k+i)\|_F \right).$$

By Proposition 1, we have for any $\lambda^i \geq 1/2$ that

$$\|Y_v\|_S^{\lambda,k} \leq 4\beta n^{b+1}b\tilde{t} \left(2\|\tilde{\mathbf{X}}\|_Q^{\lambda,k} + \bar{\gamma}\|\tilde{\mathbf{Y}}\|_F^{\lambda,k} \right) + c_2,$$

where

$$c_2 = 4\beta n^{b+1}b \sum_{i=1}^{\tilde{t}} \lambda^{-i} \left(2\|\tilde{\mathbf{X}}(t)\|_Q + \bar{\gamma}\|\tilde{\mathbf{X}}(t)\|_Q \right). \quad (38)$$

6.5 Proof of Lemma 7

It follows from (10), (17) and (22) that

$$\begin{aligned} \mathbf{x}_\pi(k+b) &= \boldsymbol{\pi}(k+b)^\top \tilde{\mathbf{X}}(k+b) \\ &= \boldsymbol{\pi}(k+b)^\top \boldsymbol{\Phi}_b^A(k) \tilde{\mathbf{X}}(k) - \boldsymbol{\pi}(k+b)^\top \sum_{i=0}^{b-1} \boldsymbol{\Phi}_{b-i}^A(k+t) \mathbf{I}_{k+i}^a \tilde{\mathbf{I}}^\top \mathbf{V}(k+t) Y_v(k+t) \\ &= \boldsymbol{\pi}(k)^\top \tilde{\mathbf{X}}(k) - \frac{1}{n} \sum_{i=0}^{b-1} \boldsymbol{\eta}_k(t) \mathbf{1}_n^\top \nabla(k+t) - \sum_{i=0}^{b-1} \mathbf{r}_k(t) \left(Y_v(k+t) - \frac{1}{n} \tilde{\mathbf{I}}(k+t) \mathbf{1}_n^\top \nabla(k+t) \right), \end{aligned} \quad (39)$$

where $\mathbf{V}(k)$ and $\tilde{\mathbf{I}}(k)$ are defined in Lemma 6, $\mathbf{r}_k(t) = \boldsymbol{\pi}(k+b)^\top \boldsymbol{\Phi}_{b-i}^A(k+t+1) \mathbf{I}_{k+i}^a \tilde{\mathbf{I}}^\top \mathbf{V}(k+t)$ and $\boldsymbol{\eta}_k(t) = \mathbf{r}_k(t) \tilde{\mathbf{I}}(k+t)$. Let $\bar{\boldsymbol{\eta}}_k := \sum_{i=0}^{b-1} \boldsymbol{\eta}_k(t)$. We have

$$\underline{\gamma} n^{-b} \leq \boldsymbol{\pi}(k+b)^\top \sum_{t=0}^{b-1} \boldsymbol{\Phi}_{b-t}^A(k+t) \mathbf{I}_{k+t}^a \tilde{\mathbf{T}} \mathbf{v}(k+t) \leq \bar{\eta}_k \leq \bar{\gamma} n b < \frac{1}{4\beta}, \forall k. \quad (40)$$

where we used the relation $v_i(k) \geq n^{-b}, \forall i \in \mathcal{V}, k \in \mathbb{N}$ from Lemma 2(iii) and $\bar{\gamma} \leq \frac{1}{4nb\beta}$ from Lemma 7.

By introducing an auxiliary term $\bar{\eta}_k \nabla f(\mathbf{x}_\pi(k))^\top$, Eq. (39) becomes

$$\begin{aligned} \mathbf{x}_\pi(k+b) &= \mathbf{x}_\pi(k) - \bar{\eta}_k \nabla f(\mathbf{x}_\pi(k))^\top + \underbrace{\sum_{t=0}^{b-1} \eta_k(t) \left(\nabla f(\mathbf{x}_\pi(k))^\top - \frac{1}{n} \mathbf{1}_n^\top \nabla(k+t) \right)}_{\mathbf{h}(k)} \\ &\quad - \sum_{t=0}^{b-1} \mathbf{r}_k(t) \left(\mathbf{Y}_v(k+t) - \frac{1}{n} \tilde{\mathbf{I}}(k+t) \mathbf{v}(k+t)^\top \mathbf{Y}_v(k+t) \right) \\ &= \mathbf{x}_\pi(k) - \bar{\eta}_k \nabla f(\mathbf{x}_\pi(k))^\top + \mathbf{h}(k) - \sum_{t=0}^{b-1} \mathbf{r}_k(t) \mathbf{S}(k+t) \mathbf{Y}_v(k+t), \end{aligned} \quad (41)$$

where we have used the relation $\mathbf{1}_n^\top \nabla(k) = \mathbf{1}_n^\top \tilde{\mathbf{Y}}(k) = \mathbf{v}(k)^\top \mathbf{Y}_v(k)$ in (14).

We now bound $\mathbf{h}(k)$ in (41). Recall that $\nabla f(\mathbf{x}) = \mathbf{1}_n^\top \nabla \mathbf{f}(\mathbf{1}_n \mathbf{x}^\top), \forall \mathbf{x} \in \mathbb{R}^m$ and $\mathbf{1}_n^\top \nabla(k) = [\mathbf{1}_n; \mathbf{0}]^\top \nabla(k) = \mathbf{1}_n^\top \nabla \mathbf{f}(\mathbf{X}(k)), \forall k$. Let $\tilde{\mathbf{I}} = [\mathbf{1}_n; \mathbf{0}_{\bar{n}-n}]$, we have

$$\begin{aligned} \|\mathbf{h}(k)\|_{\mathbb{F}} &= \left\| \sum_{t=0}^{b-1} \eta_k(t) \left(\mathbf{1}_n^\top \nabla \mathbf{f}(\mathbf{1}_n \mathbf{x}_\pi(k)) - \mathbf{1}_n^\top \nabla \mathbf{f}(\mathbf{X}(k+t)) \right) \right\|_{\mathbb{F}} \\ &\leq \sqrt{n} \beta \bar{\eta}_k \max_{0 \leq t \leq b-1} \left\| \mathbf{1}_n \mathbf{x}_\pi(k)^\top - \mathbf{X}(k+t) \right\|_{\mathbb{F}} \\ &\leq \sqrt{n} \beta \bar{\eta}_k \left(\max_{0 \leq t \leq b-1} \left\| \mathbf{1}_n \mathbf{x}_\pi(k)^\top - \mathbf{1}_n \mathbf{x}_\pi(k+t-1) \right\|_{\mathbb{F}} \right. \\ &\quad \left. + \max_{0 \leq t \leq b-1} \left\| \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \left(\mathbf{1}_n \boldsymbol{\pi}_{k+t-1}^\top \tilde{\mathbf{X}}(k+t) - \tilde{\mathbf{X}}(k+t) \right) \right\|_{\mathbb{F}} \right) \\ &\leq n\beta \bar{\eta}_k \sum_{t=0}^{b-1} \left(\left\| \tilde{\mathbf{X}}(k+t) \right\|_Q + \bar{\gamma} \left\| \tilde{\mathbf{Y}}(k+t) \right\|_{\mathbb{F}} \right), \end{aligned} \quad (42)$$

where we defined $\boldsymbol{\pi}_k := \boldsymbol{\pi}(k)$ for notional simplicity and used $\boldsymbol{\pi}_{k+t}^\top \tilde{\mathbf{X}}(k+t) = \boldsymbol{\pi}_k^\top \tilde{\mathbf{X}}(k) + \sum_{l=0}^{t-1} \boldsymbol{\pi}_{k+l}^\top \tilde{\mathbf{T}} \mathbf{I}_{k+l}^a \tilde{\mathbf{Y}}(k+l)$

from (10) and (17) to obtain the last inequality.

Then, it follows from Proposition 2 and (41) that

$$f(\mathbf{x}_\pi(k+b)) - f^* \leq \sigma^2 (f(\mathbf{x}_\pi(k)) - f^*) + \left(\frac{2}{\eta_k} + \beta \right) \left\| \mathbf{h}(k) - \sum_{t=0}^{b-1} \mathbf{r}_k(t) \mathbf{S}(k+t) \mathbf{Y}_v(k+t) \right\|_2^2,$$

where $\sigma = \sqrt{1 - \alpha \bar{\eta}_k (1 - 2\beta \bar{\eta}_k)}$, and $\sigma \leq 1 - \frac{1}{2} \alpha \bar{\eta}_k (1 - 2\beta \bar{\eta}_k) < 1 - \frac{1}{4} \alpha \bar{\eta}_k < 1$ as $\sqrt{1-a} \leq 1 - \frac{1}{2} a$,

$\forall a \in (0,1)$, and (40). Using $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and (42), we obtain

$$\begin{aligned} \tilde{f}(k+b) &= \sqrt{f(\mathbf{x}_\pi(k+b)) - f^*} \leq \sigma \sqrt{f(\mathbf{x}_\pi(k)) - f^*} \\ &\quad + \left(\sqrt{\beta} + \sqrt{2} \bar{\eta}_k^{-\frac{1}{2}} \right) \left(\left\| \mathbf{h}(k) \right\|_{\mathbb{F}} + \sum_{t=0}^{b-1} \left\| \mathbf{r}_k(t) \right\|_2 \left\| \mathbf{Y}_v(k+t) \right\|_s \right) \\ &\leq \tilde{\eta} \sum_{t=0}^{b-1} n\beta \left(\left\| \tilde{\mathbf{X}}(k+t) \right\|_Q + \bar{\gamma} \left\| \tilde{\mathbf{Y}}(k+t) \right\|_{\mathbb{F}} \right) + \left\| \mathbf{Y}_v(k+t) \right\|_s + \sigma \tilde{f}(k), \end{aligned} \quad (43)$$

where $\tilde{\eta} := \left(\bar{\eta}_k \sqrt{\beta} + \sqrt{2} \bar{\eta}_k^{-\frac{1}{2}} \right)$ and the two inequalities follow from (42) and $\left\| \mathbf{r}_k(t) \right\|_2 \leq \bar{\eta}_k, \forall k, t$.

Since $\lambda^b > 1 - \frac{1}{8} \alpha \underline{\gamma} n^{-b} > 1 - \frac{1}{8} \alpha \bar{\eta}_k, \sigma < 1 - \frac{1}{4} \alpha \bar{\eta}_k$, we obtain from Proposition 1 and (43) that

$$\begin{aligned}\tilde{f}^{\lambda,k} &\leq \frac{16b}{\alpha\sqrt{\bar{\eta}_k}} \left(n\beta \|\tilde{X}\|_Q^{\lambda,k} + \|\mathbf{Y}_v\|_S^{\lambda,k} + n\beta\bar{\gamma} \|\tilde{Y}\|_F^{\lambda,k} \right) + c_3 \\ &\leq \frac{16n^{\frac{b}{2}}b}{\alpha\sqrt{\underline{\gamma}}} \left(n\beta \|\tilde{X}\|_Q^{\lambda,k} + \|\mathbf{Y}_v\|_S^{\lambda,k} + n\beta\bar{\gamma} \|\tilde{Y}\|_F^{\lambda,k} \right) + c_3,\end{aligned}$$

where the first inequality used (40) and $\bar{\eta}_k \leq 1/(4\beta)$, and

$$c_3 = \frac{16n^{\frac{b}{2}}nb}{\alpha\sqrt{\underline{\gamma}}} \sum_{t=0}^{b-1} \lambda^{-t} \left(\beta \|\tilde{X}(t)\|_Q + \|\mathbf{Y}_v(t)\|_S + \beta\bar{\gamma} \|\tilde{Y}(t)\|_F \right). \quad (44)$$

6.6 Proof of Lemma 8

Since $\tilde{Y}(k) = \mathbf{V}(k)\mathbf{Y}_v(k) = \mathbf{V}(k)\mathbf{S}(k)\mathbf{Y}_v(k) + \frac{1}{n}\mathbf{V}(k)\tilde{\mathbf{I}}(k)\mathbf{v}(k)^\top\mathbf{Y}_v(k)$, we have

$$\|\tilde{Y}(k)\|_F \leq \|\mathbf{V}(k)\|_2 \left(\|\mathbf{Y}_v(k)\|_S + \left\| \frac{1}{n}\tilde{\mathbf{I}}(k)\mathbf{v}(k)^\top\mathbf{Y}_v(k) \right\|_F \right) \leq n\|\mathbf{Y}_v(k)\|_S + \|\tilde{\mathbf{I}}(k)\mathbf{1}_n^\top\nabla(k)\|_F. \quad (45)$$

Note that

$$\begin{aligned}\|\tilde{\mathbf{I}}(k)\mathbf{1}_n^\top\nabla(k)\|_F &= \left\| \tilde{\mathbf{I}}(k)\mathbf{1}_n^\top \left(\nabla f(\mathbf{X}(k)) - \nabla f(\mathbf{1}_n(\mathbf{x}^*)^\top) \right) \right\|_F \leq \|\tilde{\mathbf{I}}(k)\mathbf{1}_n^\top\|_2 \left\| \nabla f(\mathbf{X}(k)) - \nabla f(\mathbf{1}_n(\mathbf{x}^*)^\top) \right\|_F \\ &\leq \sqrt{\tilde{n}}\beta \left\| \mathbf{X}(k) - \mathbf{1}_n(\mathbf{x}^*)^\top \right\|_F = \sqrt{\tilde{n}}\beta \left\| \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \tilde{X}(k) - \begin{bmatrix} \mathbf{1}_n \\ \mathbf{0} \end{bmatrix} (\mathbf{x}^*)^\top \right\|_F \\ &\leq \sqrt{\tilde{n}}\beta \left(\left\| \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \tilde{X}(k) - \begin{bmatrix} \mathbf{1}_n \\ \mathbf{0} \end{bmatrix} \boldsymbol{\pi}(k-1)^\top \tilde{X}(k) \right\|_F + \left\| \begin{bmatrix} \mathbf{1}_n(\boldsymbol{\pi}(k-1)^\top \tilde{X}(k) - (\mathbf{x}^*)^\top) \\ \mathbf{0} \end{bmatrix} \right\|_F \right) \\ &\leq \sqrt{\tilde{n}}\beta \left(\left\| \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{1}_n \\ \mathbf{0} \end{bmatrix} \boldsymbol{\pi}(k-1)^\top \right\|_F Q \tilde{X}(k) \right\|_F + \sqrt{\tilde{n}}\beta \left\| \begin{bmatrix} \mathbf{1}_n(\mathbf{x}_\pi(k) - (\mathbf{x}^*)^\top) \\ \mathbf{0} \end{bmatrix} \right\|_F \\ &\leq \sqrt{\tilde{n}}\beta \left(\sqrt{n} \|\tilde{X}(k)\|_Q + \sqrt{n} \|\mathbf{x}_\pi(k) - \mathbf{x}^*\|_F \right) \\ &\leq \sqrt{\tilde{n}}\beta \left(\sqrt{n} \|\tilde{X}(k)\|_Q + \frac{\sqrt{2n}}{\sqrt{\alpha}} \sqrt{f(\mathbf{x}_\pi(k)) - f^*} \right),\end{aligned} \quad (46)$$

where $\mathbf{x}^* = \Pi_{\mathcal{X}^*}(\mathbf{x}_\pi(k))$ and the last inequality follows from that the Polyak-Łojasiewicz condition implies the quadratic growth condition (Theorem 2 (Karimi et al. , 2016)), i. e. , $f(\mathbf{x}) - f^* \geq \frac{\alpha}{2} \|\mathbf{x} - \Pi_{\mathcal{X}^*}(\mathbf{x})\|_2^2, \forall \mathbf{x}$.

Substituting (46) into (45) yields

$$\|\tilde{Y}(k)\|_F \leq n\|\mathbf{Y}_v(k)\|_S + \sqrt{\tilde{n}}\beta \left(\sqrt{n} \|\tilde{X}(k)\|_Q + \frac{\sqrt{2n}}{\sqrt{\alpha}} \sqrt{f(\mathbf{x}_\pi(k)) - f^*} \right).$$

The desired result is obtained by the definition of $\|\tilde{Y}\|_F^{\lambda,k}$.

7 Numerical examples

We use APPG to train a multi-class logistic regression classifier in a distributed manner on the *Coverttype* dataset (Dheeru et al. , 2017), where the objective function takes the following form

$$f(\mathbf{X}) = -\sum_{i=1}^{n_s} \sum_{j=1}^{n_c} l_j^i \log \left(\frac{\exp(\mathbf{x}_j^\top \mathbf{s}^i)}{\sum_{j=1}^{n_c} \exp(\mathbf{x}_j^\top \mathbf{s}^i)} \right) + \frac{\rho}{2} \|\mathbf{X}\|_F^2.$$

Here $n_s = 581012$ is the number of training instances, $n_c = 7$ is the number of classes, $n_f = 55$ is the number of

features, $s^i \in \mathbb{R}^{55}$ is the feature vector of the i -th instance, $l^i = [l_1^i, \dots, l_7^i]^\top$ is the label vector of the i -th instance using the one-hot encoding, $X = [x_1, \dots, x_7] \in \mathbb{R}^{n_7 \times n_s}$ is the parameters to be optimized, $\rho = 20$ is a regularization factor.

Environment: APPG is implemented in Python with OpenMPI 1.10 on Ubuntu 14.04. The hardware is a server with 28 Xeon E5-2660 cores. Each core serves as a computing node.

Distributed Data: We first normalize non-categorical features by subtracting the mean and dividing by the standard deviation in the whole dataset. Then, we *sort* the data by digit label, and sequentially partition it into n parts (with different sizes), where each node (core) only has *exclusive* access to one part. Thus, we are dealing with distributed datasets.

Topology: The directed network among nodes is as follows: Each node i sends messages to node $\text{mod}(2^j + i, n)$, where $j \in \mathbb{N} \cap [0, \log_2(n))$ and $\text{mod}(a, b)$ returns the remainder after division of a by b . Thus, each node has $\mathcal{O}(\log(n))$ out-neighbors, which results in a relatively sparse directed networks. Note that gossip-based asynchronous algorithms are generally not applicable over this directed network.

Stepsize: The stepsize of each algorithm is tuned via a grid search around $0.5/n_s$.

Local Termination Criteria: Node i stops locally if the values of y_i in last n consecutive iterations are less than $300/n_s$.

7.1 Convergence performance and linear speedup

We implement APPG over $n = 1, 6, 12, 18, 24$ nodes ($n = 1$ is used as a baseline since APPG reduces to standard centralized gradient descent method). The training loss w. r. t. running time is plotted in Fig. 6(a), which validates the linear convergence of APPG, and shows that the training time is significantly reduced with the increase of number of nodes.

Fig. 6(b) depicts the training loss w. r. t. the number of iterations. We find that the number of iterations required to achieve the same accuracy is close to each other for different number of nodes. The time for a node to finish an iteration is proportional to the size of its local dataset, and hence is roughly inversely proportional to the number of nodes, which suggests that using n nodes may reduce $\mathcal{O}(n)$ times of training time than that of one node.

To further illustrate this property, we study the speedup of APPG defined as $S_n := T_n/T_1$, where T_n is the running time of the APPG with n node (s) when the training loss decays to 0.005. Fig. 7(a) shows that the APPG achieves a roughly linear speedup in convergence rate w. r. t. the number of nodes. Then, we test a synchronized version of APPG, which is done by adding a barrier after each update and is mathematically equivalent to the push-pull method in Pu et al. (2021). One can find that the synchronous version has an approximately linear speedup when the number of cores is small, but it decreases fast when the number of cores is relatively large.

Ideally, the speedup would be n when using n nodes. However, the communication among nodes introduces delays and staleness to the algorithm, which degrades the convergence rate. In practice, a higher speedup can be achieved by using a lower latency network with larger bandwidth.

7.2 Robustness of APPG to slow cores

We evaluate the robustness of APPG by forcing one core in the network to slow down. This is achieved by adding an artificial waiting time (20 ms, a normal iteration takes about 15 ms with 24 cores) after each local iteration of a node, which simulates either the slow computation or slow communication.

Fig. 7(b) shows the speedup of APPG and the synchronous implementation of APPG in this scenario. It indicates that the synchronous counterpart of APPG has a sharp reduction in convergence rate even when only 1 core slows down. In contrast, APPG still keeps an almost linear speedup. This result is also consistent with that in Lian et al. (2018); Zhang et al. (2020a).

Introducing the slowing core also brings an essential problem of asynchronous algorithms, that is, the cores

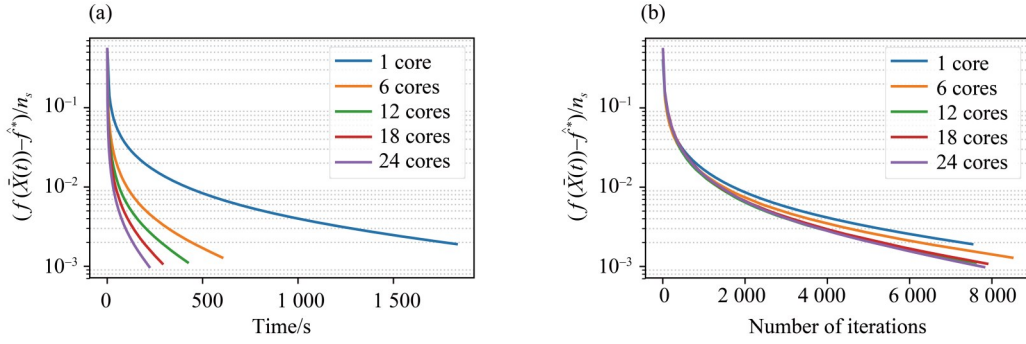


Fig. 6 Convergence performance with different number of nodes. (a) Training loss w. r. t. running time of APPG; (b) Training loss w. r. t. number of iterations of APPG

have *uneven* update rates. To show its effect to the performance, we compare APPG to a gossip-based asynchronous algorithm AD-PSGD (Lian et al., 2018) with full local gradients. Note that APPG can only work over *undirected* networks, and hence we modify the network by adding a reversed edge to each edge in the directed network, while APPG is still implemented over the directed network. Fig. 8 shows the result over 12 nodes and 24 nodes, where AD-PSGD fails to converge to the exact optimum. In contrast, APPG converges exactly despite a bit slower convergence rate.

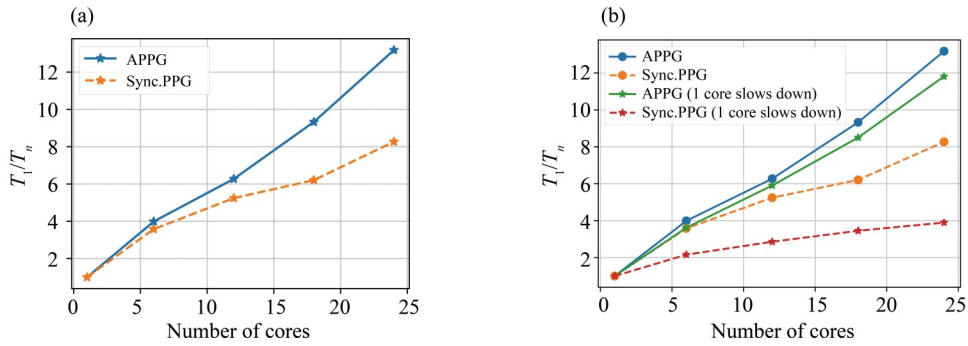


Fig. 7 (a) Speedup of APPG and the synchronous implementation of APPG; (b) Speedup of APPG and the 'synchronized' APPG with one slow core

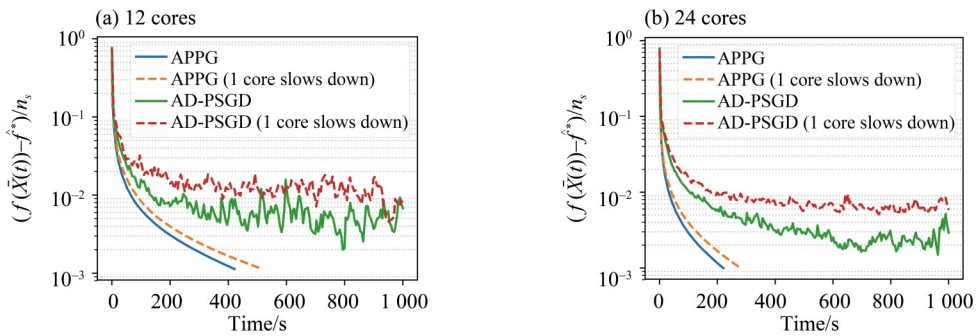


Fig. 8 Convergence of APPG and AD-PSGD with full local gradient when one node is artificially slowed down.

8 Conclusion

This paper has proposed APPG for distributed optimization which allows nodes to connect via a directed communication network and update with uncoordinated computation and stale information from neighbors. Future works may focus on accelerating APPG and extending it to stochastic optimization.

References:

- ASSRAN M S, RABBAT M G, 2021. Asynchronous gradient push[J]. *IEEE Trans Autom Control*, 66(1): 168–183.
- ASSRAN M, LOIZOU N, BALLAS N, et al, 2019. Stochastic gradient push for distributed deep learning[C]//*Proc Int Conf Mach Learn*, 97: 344–353.
- ASSRAN M, AYTEKIN A, FEYZMAHDAVIAN H R, et al, 2020. Advances in asynchronous parallel and distributed optimization [J]. *Proc IEEE*, 108(11): 2013–2031.
- CHANG T H, HONG M Y, WANG X F, 2015. Multi-agent distributed optimization via inexact consensus ADMM[J]. *IEEE Trans Signal Process*, 63(2): 482–497.
- CHANG T H, HONG M Y, LIAO W C, et al, 2016. Asynchronous distributed ADMM for large-scale optimization—Part I: Algorithm and convergence analysis[J]. *IEEE Trans Signal Process*, 64(12): 3118–3130.
- DHEERU D, KARRA TANISKIDOU E, 2017. UCI machine learning repository [EB/OL]. University of California, Irvine, School of Information and Computer Sciences. <http://archive.ics.uci.edu/ml>.
- FAZEL M, GE R, KAKADE S M, et al, 2018. Global convergence of policy gradient methods for the linear quadratic regulator [C]//*Proc Int Conf Mach Learn*, 80: 1467–1476.
- JAKOVETIĆ D, 2019. A unification and generalization of exact distributed first-order methods[J]. *IEEE Trans Signal Inf Process Netw*, 5(1): 31–46.
- KARIMI H, NUTINI J, SCHMIDT M, 2016. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition[C]//*Joint Eur Conf Mach Learn Knowl Discovery Databases*, 795–811.
- LI H Q, HUANG C C, WANG Z, et al, 2020. Computation-efficient distributed algorithm for convex optimization over time-varying networks with limited bandwidth communication[J]. *IEEE Trans Signal Inf Process Netw*, 6: 140–151.
- LI S, BASAR T, 1987. Asymptotic agreement and convergence of asynchronous stochastic algorithms[J]. *IEEE Trans Autom Control*, 32(7): 612–618.
- LI Z, SHI W, YAN M, 2019. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates[J]. *IEEE Trans Signal Process*, 67(17): 4494–4506.
- LIAN X R, ZHANG W, ZHANG C, et al, 2018. Asynchronous decentralized parallel stochastic gradient descent[C]//*Proc Int Conf Mach Learn*, 80: 3043–3052.
- NEDIĆ A, 2011. Asynchronous broadcast-based convex optimization over a network[J]. *IEEE Trans Autom Control*, 56(6): 1337–1351.
- NEDIĆ A, OLSHEVSKY A, 2015. Distributed optimization over time-varying directed graphs[J]. *IEEE Trans Autom Control*, 60(3): 601–615.
- NEDIĆ A, OLSHEVSKY A, 2016. Stochastic gradient-push for strongly convex functions on time-varying directed graphs[J]. *IEEE Trans Autom Control*, 61(12): 3936–3947.
- NEDIĆ A, OZDAGLAR A, 2009. Distributed subgradient methods for multi-agent optimization[J]. *IEEE Trans Autom Control*, 54(1): 48–61.
- NEDIĆ A, OZDAGLAR A, 2010. Convergence rate for consensus with delays[J]. *J Glob Optim*, 47(3): 437–456.
- NEDIĆ A, OLSHEVSKY A, SHI W, 2017. Achieving geometric convergence for distributed optimization over time-varying graphs[J]. *SIAM J Optim*, 27(4): 2597–2633.
- OLSHEVSKY A, TSITSIKLIS J N, 2011. Convergence speed in distributed consensus and averaging[J]. *SIAM Rev*, 53(4): 747–772.
- PU S, SHI W, XU J M, et al, 2021. Push-pull gradient methods for distributed optimization in networks[J]. *IEEE Trans Autom Control*, 66(1): 1–16.
- SAADATNIAKI F, XIN R, KHAN U A, 2020. Decentralized optimization over time-varying directed graphs with row and column-stochastic matrices[J]. *IEEE Trans Autom Control*, 65(11): 4769–4780.
- SCUTARI G, SUN Y, 2019. Distributed nonconvex constrained optimization over time-varying digraphs[J]. *Math Program*, 176(1/2): 497–544.
- SPIRIDONOFF A, OLSHEVSKY A, PASCHALIDIS I C, 2020. Robust asynchronous stochastic gradient-push: Asymptotically optimal and network-independent performance for strongly convex functions[J]. *J Mach Learn Res*, 21(58): 1–47.

- SUN H, HONG M Y, 2019. Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms[J]. *IEEE Trans Signal Process*, 67(22): 5912–5928.
- TIAN Y, SUN Y, SCUTARI G, 2020. Achieving linear convergence in distributed asynchronous multi-agent optimization[J]. *IEEE Trans Autom Control*, 65(12): 5264–5279.
- TOURI B, 2012. *Product of random stochastic matrices and distributed averaging*[M]. Berlin: Springer-Verlag.
- TSITSIKLIS J, BERTSEKAS D, ATHANS M, 1986. Distributed asynchronous deterministic and stochastic gradient optimization algorithms[J]. *IEEE Trans Autom Control*, 31(9): 803–812.
- WU T Y, YUAN K, LING Q, et al, 2018. Decentralized consensus optimization with asynchrony and delays[J]. *IEEE Trans Signal Inf Process Netw*, 4(2): 293–307.
- XIE P, YOU K Y, TEMPO R, et al, 2018. Distributed convex optimization with inequality constraints over time-varying unbalanced digraphs[J]. *IEEE Trans Autom Control*, 63(12): 4331–4337.
- XIN R, KHAN U A, 2018. A linear algorithm for optimization over directed graphs with geometric convergence[J]. *IEEE Control Syst Lett*, 2(3): 315–320.
- XIN R, KAR S, KHAN U A, 2020. Decentralized stochastic optimization and machine learning: A unified variance-reduction framework for robust performance and fast convergence[J]. *IEEE Signal Process Mag*, 37(3): 102–113.
- XU J M, ZHU S Y, SOH Y C, et al, 2018. Convergence of asynchronous distributed gradient methods over stochastic networks[J]. *IEEE Trans Autom Control*, 63(2): 434–448.
- XU J M, TIAN Y, SUN Y, et al, 2021. Distributed algorithms for composite optimization: Unified framework and convergence analysis[J]. *IEEE Trans Signal Process*, 69: 3555–3570.
- YI X L, LI X X, XIE L H, et al, 2020. Distributed online convex optimization with time-varying coupled inequality constraints[J]. *IEEE Trans Signal Process*, 68: 731–746.
- YUAN K, YING B C, ZHAO X C, et al, 2019. Exact diffusion for distributed optimization and learning—Part I: Algorithm development[J]. *IEEE Trans Signal Process*, 67(3): 708–723.
- ZHANG J Q, YOU K Y, 2020a. AsySPA: An exact asynchronous algorithm for convex optimization over digraphs[J]. *IEEE Trans Autom Control*, 65(6): 2494–2509.
- ZHANG J Q, YOU K Y, CAI K, 2020b. Distributed dual gradient tracking for resource allocation in unbalanced networks[J]. *IEEE Trans Signal Process*, 68: 2186–2198.
- ZHAO X C, SAYED A H, 2015. Asynchronous adaptation and learning over networks—Part I: Modeling and stability analysis[J]. *IEEE Trans Signal Process*, 63(4): 811–826.

(责任编辑 冯兆永)