

MIN-MAX OPTIMIZATION OVER SLOWLY TIME-VARYING GRAPHS

© 2023 г. Nhat Trung Nguyen^{1,*}, Alexander Rogozin^{1,**}, Dmitriy Metelev^{1,***},
Alexander Gasnikov^{1,2,3,4,†}

Presented at AiJourney 2023 conference

Поступило 30.08.2023

После доработки

Принято к публикации

Distributed optimization is an important direction of research in modern optimization theory. Its applications include large scale machine learning, distributed signal processing and many others. The paper studies decentralized min-max optimization for saddle point problems. Saddle point problems arise in training adversarial networks and in robust machine learning. The focus of the work is optimization over (slowly) time-varying networks. The topology of the network changes from time to time, and the velocity of changes is limited. We show that, analogically to decentralized optimization, it is sufficient to change only two edges per iteration in order to slow down convergence to the arbitrary time-varying case. At the same time, we investigate several classes of time-varying graphs for which the communication complexity can be reduced.

Ключевые слова и фразы: saddle point problem, decentralized optimization, time-varying graph, extragradient method

1. INTRODUCTION

This paper studies min-max optimization problems of type

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} f(x, y) := \frac{1}{M} \sum_{m=1}^M f_m(x, y), \quad (1)$$

where functions $f_m(x, y)$ are convex in x and concave in y and \mathcal{X}, \mathcal{Y} are closed convex sets. Each function $f_m(x, y)$ is held locally at some computational node. The nodes are connected to each other via a decentralized communication network. Each agent can perform local computations and exchange information with its immediate neighbors in the network. Additionally, the network is allowed to change with time. Due to some malfunctions or disturbances, the links in the network may fail or reappear from time to time. This type of networks is called *time-varying graphs*.

There are numerous applications for optimization over time-varying networks [15, 13]. They include distributed machine learning [16, 5, 14], distributed control of power systems [17, 6], vehicle control [18], distributed sensing [1].

Decentralized optimization and min-max optimization first-order methods use two types of steps: local gradient updates and inter-node communications. We consider the case when communications are done in synchronized rounds. Therefore, the complexity of the method is measured by two quantities: number of communication rounds and number of local oracle calls. These quantities depend on the problem characteristics, which include network condition number χ , function condition number L/μ and desired accuracy ε . Here L is the Lipschitz constant of the objective gradient and μ is the strong convexity constant.

¹Moscow Institute of Physics and Technology, Moscow, Russia

²Institute for Information Transportation Problems, Moscow, Russia

³Caucasus Mathematic Center of Adygh State University, Moscow, Russia

⁴Ivannikov Institute for System Programming of the Russian Academy of Sciences, Research Center for Trusted Artificial Intelligence, Moscow, Russia

*E-mail: ngnhtrg@phystech.edu

**E-mail: aleksandr.rogozin@phystech.edu

***E-mail: metelev.ds@phystech.edu

†E-mail: gasnikov@yandex.ru

This paper is devoted to *slowly time-varying graphs*. That means that only a limited number of edges is allowed to change after each communication round. We provide lower complexity bounds for the considered class of problems. Also we propose min-max optimization methods with better communication complexity for two particular classes of slowly time-varying networks.

For optimization over networks (not min-max) lower bounds are known as well as corresponding optimal algorithms. For static graphs, lower bounds were derived in [19] and in the same paper the optimal dual (i.e. using a dual oracle) algorithm was proposed. Optimal decentralized methods with primal oracle were proposed in [9]. Considering time-varying graphs (with arbitrary changes at each iteration), lower complexity bounds were proposed in [8]. Optimal primal algorithm was proposed in the same paper [8], and optimal dual method first appeared in [10]. After that, lower bounds for slowly time-varying graphs with different velocities of network changes were studied in [12]. In [11], it was shown that it is sufficient to change only two edges at each iteration to make communication complexity equal to arbitrary time-varying graph. The overview of lower bounds for decentralized optimization is presented in Table 1 (notation $\Omega(\cdot)$ is omitted).

The results for decentralized saddle-point problems are analogical to optimization. Lower bounds for min-max optimization over static graphs were given in [4]. The same paper [4] proposed algorithms optimal up to a logarithmic term. The case of (arbitrary) time-varying graphs was studied in [2] along with methods optimal up to a logarithmic factor. Optimal algorithms for sum-type variational inequalities (a generalization of saddle point problem) were proposed in [7]. Finally, this paper studies lower bounds for saddle-point problems over slowly time-varying graphs (only two edge changes per iteration). The corresponding results are presented in Table 2. It is worth noting that the lower complexity bounds are same as for optimization (Table 1), except for replacing $\sqrt{L/\mu}$ by L/μ .

This paper has the following organization. In Section 2, we introduce the needed assumptions and notation. After that, in Section 3, we show how to get an acceleration in communications using additional assumptions on the time-varying network. In Section 4, we provide lower bounds for slowly time-varying networks.

	static	time-var.	slowly time-var.
comm.	$\sqrt{\chi}\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$	$\chi\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$	$\chi\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$
oracle	$\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$	$\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$	$\sqrt{\frac{L}{\mu}}\log\frac{1}{\varepsilon}$
paper	[19]	[8]	[11]

Table 1: Lower bounds for optimization

Table 2: Lower bounds for saddle point problems

	static	time-var.	slowly time-var.
comm.	$\sqrt{\chi}\frac{L}{\mu}\log\frac{1}{\varepsilon}$	$\chi\frac{L}{\mu}\log\frac{1}{\varepsilon}$	$\chi\frac{L}{\mu}\log\frac{1}{\varepsilon}$
oracle	$\frac{L}{\mu}\log\frac{1}{\varepsilon}$	$\frac{L}{\mu}\log\frac{1}{\varepsilon}$	$\frac{L}{\mu}\log\frac{1}{\varepsilon}$
paper	[4]	[2]	This paper

Table 2: Lower bounds for saddle point problems

2. NOTATION AND ASSUMPTIONS

Smoothness and strong convexity.

We work with the problem (1), where the sets $\mathcal{X} \subseteq \mathbb{R}^{n_x}$ and $\mathcal{Y} \subseteq \mathbb{R}^{n_y}$ are closed convex sets. Additionally, we introduce the set $\mathcal{Z} = \mathcal{X} \times \mathcal{Y} \subseteq \mathbb{R}^{n_z}$, $z = (x, y)$, $n_z = n_x + n_y$, and the operator F :

$$F_m(z) = F_m(x, y) = \begin{pmatrix} \nabla_x f_m(x, y) \\ -\nabla_y f_m(x, y) \end{pmatrix}, \quad F(z) = \frac{1}{M} \sum_{m=1}^M F_m(z).$$

Assumption 1. Let functions $f(x, y)$ and $f_m(x, y)$ satisfy following properties:

1. Function $f(x, y)$ is L -smooth, i.e. for all $z_1, z_2 \in \mathcal{Z}$ it holds

$$\|F(z_1) - F(z_2)\| \leq L\|z_1 - z_2\|.$$

2. For all m , $f_m(x, y)$ is L_{\max} -smooth, i.e. for all $z_1, z_2 \in \mathcal{Z}$ it holds

$$\|F_m(z_1) - F_m(z_2)\| \leq L_{\max}\|z_1 - z_2\|.$$

3. Function $f(x, y)$ is strongly-convex-strongly-concave with constant μ , i.e. for all $z_1, z_2 \in \mathcal{Z}$ it holds

$$\langle F(z_1) - F(z_2), z_1 - z_2 \rangle \geq \mu\|z_1 - z_2\|^2.$$

Decentralized Communication.

At each communication round, we use a graph to represent the connection between computing nodes. Denote the network of communications over time by the sequence of graphs $\{\mathcal{G}^k\}_{k=0}^{\infty} = \{(\mathcal{V}, \mathcal{E}^k)\}_{k=0}^{\infty}$, where

$\mathcal{V} = \{1, \dots, M\}$ is the set of nodes and \mathcal{E}^k is the set of available connections at k -th communication round. For each node $m \in \mathcal{V}$, we use the notation $\mathcal{N}_m^k = \{i \in \mathcal{V} | (i, m) \in \mathcal{E}^k\}$ to indicate the set of its neighbors at round k and at that time it can only communicate with nodes in \mathcal{N}_m^k .

Gossip Matrices Each computational node m contains its own local vector $z_m = (x_m, y_m)$. It is required to satisfy the consensus constraints $z_1 = \dots = z_M$. For this purpose, we use a concept called *gossip matrix*. **Assumption 2.** Each graph in the time-varying network correspond to a gossip matrix $W^k \in \mathbb{R}^{M \times M}$ that satisfies the following properties.

1. W^k is positive semi-definite,
2. $W_{i,j}^k = 0$ if $i \neq j$ and $(i, j) \notin \mathcal{E}^k$,
3. $\ker W^k = \text{span}(\mathbf{1})$, where $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^M$.

The number $\chi(W) = \frac{\lambda_{\max}(W)}{\lambda_{\min}^+(W)}$ is called the *condition number* of gossip matrix W , where $\lambda_{\max}(W)$ and $\lambda_{\min}^+(W)$ denote the largest and smallest positive eigenvalue of W . For a time-varying network $\{\mathcal{G}^k\}_{k=1}^{\infty}$, the condition number is given by $\chi = \sup_{k \in \mathbb{N} \cup \{0\}} \frac{\lambda_{\max}(W^k)}{\lambda_{\min}^+(W^k)}$.

In this paper, we also consider network with Laplacian matrices $\mathbf{L}(\mathcal{G}^k)$, which is a typical example of gossip matrix.

We introduce the *consensus space* $\mathcal{L} \subseteq \mathbb{R}^{Mn_z}$, defined by

$$\mathcal{L} = \{\mathbf{z} = (z_1^T, \dots, z_M^T)^T \in \mathbb{R}^{Mn_z} : z_1 = \dots = z_M\}. \quad (2)$$

Consider also the space $\mathcal{L}^\perp \subseteq \mathbb{R}^{Mn_z}$ which is an orthogonal complement to the space \mathcal{L} , defined by

$$\mathcal{L}^\perp = \{\mathbf{z} = (z_1^T, \dots, z_M^T)^T \in \mathbb{R}^{Mn_z} : \sum_{m=1}^M z_m = 0\}. \quad (3)$$

3. UPPER BOUNDS

In this section, we cover two classes of time-varying networks: networks with connected skeleton and networks with small random Markovian changes. For both scenarios, we propose a decentralized optimization algorithm that uses an auxiliary consensus procedure. We show that for the considered classes of problems, the dependence of communication complexity on factor χ may be reduced from χ to $\sqrt{\chi}$ with additional terms. The overview of results is presented in Table 3 (the $O(\cdot)$ notation is omitted).

	arbitrary time-var.	slowly time-var.	connected skeleton	Markovian
comm.	$\chi \frac{L}{\mu} \log \frac{1}{\varepsilon}$	$\chi \frac{L}{\mu} \log \frac{1}{\varepsilon}$	$\sqrt{\chi} \log \chi \frac{L}{\mu} \log^2 \frac{1}{\varepsilon}$	$\tau \left(\sqrt{\chi} + \frac{\rho^2}{(\lambda_{\min}^+)^2} \right) \frac{L}{\mu} \log^2 \frac{1}{\varepsilon}$
oracle	$\frac{L}{\mu} \log \frac{1}{\varepsilon}$	$\frac{L}{\mu} \log \frac{1}{\varepsilon}$	$\frac{L}{\mu} \log \frac{1}{\varepsilon}$	$\frac{L}{\mu} \log \frac{1}{\varepsilon}$
paper	[7]	[7]	This paper	This paper

Table 3: Upper bounds for saddle point problems over arbitrary and slowly time-varying graphs

Our algorithms are based on extragradient method with consensus subroutine. They make several communication rounds after each extragradient step. After a sufficient number of communications, consensus is reached up to a desired accuracy. Such approximate averaging makes the trajectories of computational nodes almost synchronized. For each class of time-varying networks, we use a corresponding consensus subroutine and incorporate it into extragradient method to get a decentralized optimization algorithm.

Accelerated Gossip with Connected Skeleton and Non-Recoverable Links.

First, we focus on the type of graphs with connected skeleton. We assume that all graphs in the sequence have a common connected subgraph that we call a *skeleton*. The edges may still appear and disappear, but each node remembers which incident links have failed at least once and stops communicating by that links. This strategy we call *non-recoverable links*. Effectively the communication network only loses edges at each iteration, but remains connected. In other words, the graph of interest can be called "monotonically decreasing".

Assumption 3. Graph sequence $\{\mathcal{G}^k = (\mathcal{V}, \mathcal{E}^k)\}_{k=0}^{\infty}$ has a connected skeleton: there exists a connected graph $\hat{\mathcal{G}} = (\mathcal{V}, \hat{\mathcal{E}})$ such that for all $k \in \mathbb{N} \cup \{0\}$ we have $\hat{\mathcal{E}} \subset \mathcal{E}^k$, $\lambda_{\max}(\mathbf{L}(\mathcal{G}^k)) \leq \lambda_{\max}$ and $\lambda_{\min}^+ \leq \lambda_{\min}^+(\hat{\mathcal{G}})$.

With these properties of the network, we introduce the following algorithm for consensus.

Algorithm 1 Accelerated Gossip with Non-Recoverable Links (**AccGossipNonRecoverable**)

Require: Vectors z_1, \dots, z_M , number of iterations H , current communication round number k_0 . Step sizes $\eta, \beta > 0$.
Construct column vector $\mathbf{z} = (z_1^T \dots z_M^T)^T$.
Set $\mathbf{u}^0 = \mathbf{z}^0 = \mathbf{z}$.
Every node $i = 1, 2, \dots, M$ initializes set of neighbors $\mathcal{N}_i = \mathcal{N}_i^{k_0}$.
for $k = 0, 1, \dots, H - 1$ **do**
 for $i = 1, 2, \dots, M$ **do**
 Update the set of nodes to which the node communicates: $\mathcal{N}_i = \mathcal{N}_i \cap \mathcal{N}_i^{k_0+k}$.
 $u_i^{k+1} = z_i^k - \eta \left(|\mathcal{N}_i| z_i^k - \sum_{j \in \mathcal{N}_i} z_j^k \right)$
 $z_i^{k+1} = (1 + \beta) u_i^{k+1} - \beta u_i^k$
 end for
end for
return z_1^H, \dots, z_M^H .

Lemma 1. (From the proof of Theorem 4.3 in [12]) Let Assumptions 3 hold and $\{\hat{z}_m\}_{m=1}^M$ be output of Algorithm 1 with input $\{z_m\}_{m=1}^M$ and step sizes $\eta = 1/\lambda_{\max}$, $\beta = (\sqrt{\chi} - 1)/(\sqrt{\chi} + 1)$, where $\chi = \lambda_{\max}/\lambda_{\min}^+$. Then after H iterations, it holds that

$$\frac{1}{M} \sum_{m=1}^M \|\hat{z}_m - \bar{z}\|^2 \leq \frac{2\chi}{M} \left(1 - \frac{1}{\sqrt{\chi}}\right)^H \sum_{m=1}^M \|z_m - \bar{z}\|^2, \quad (4)$$

where $\bar{z} = \frac{1}{M} \sum_{m=1}^M z_m$.

Based on Algorithm 1, it is possible to develop a decentralized algorithm for solving the saddle point problem 1 with a sequence of graphs that have a connected skeleton and non-recoverable links

Algorithm 2 Time-Varying with Non-Recoverable Links Decentralized Extra Step Method

Require: Step size $\gamma > 0$; number of **AccGossipNonRecoverable** steps H , communication rounds K , number of iterations N .
Choose $(x^0, y^0) = z^0 \in \mathcal{Z}, z_m^0 = z^0$.
for $k = 0, 1, 2, \dots, N - 1$ **do**
 Each machine m computes $\hat{z}_m^{k+1/2} = z_m^k - \gamma \cdot F_m(z_m^k)$
 Communication: $\hat{z}_1^{k+1/2}, \dots, \hat{z}_M^{k+1/2} = \text{AccGossipNonRecoverable}(\hat{z}_1^{k+1/2}, \dots, \hat{z}_M^{k+1/2}, H)$
 Each machine m computes $z_m^{k+1/2} = \text{proj}_{\mathcal{Z}}(\hat{z}_m^{k+1/2})$
 Each machine m computes $\hat{z}_m^{k+1} = z_m^{k+1/2} - \gamma \cdot F_m(z_m^{k+1/2})$
 Communication: $\hat{z}_1^{k+1}, \dots, \hat{z}_M^{k+1} = \text{AccGossipNonRecoverable}(\hat{z}_1^{k+1}, \dots, \hat{z}_M^{k+1}, H)$
 Each machine m computes $z_m^{k+1} = \text{proj}_{\mathcal{Z}}(\hat{z}_m^{k+1})$
end for

Theorem 1. Suppose Assumptions 1, 3 hold. Let problem (1) be solved by Algorithm 2 with $\gamma \leq \frac{1}{4L_{\max}}$. Then in order to achieve ε_0 -approximate consensus at each iteration, it takes

$$H = \left(\sqrt{\chi} \log \left[\chi \left(4 + \frac{\frac{1}{2} \|z^0 - z^*\|^2 + \frac{Q^2}{2L_{\max}^2}}{\varepsilon_0^2} \right) \right] \right) \text{ communications,}$$

where $Q^2 = \frac{1}{M} \sum_{m=1}^M \|F_m(z^*)\|^2$ and z^* is a solution of (1).

Proof. Let us have ε_0 -accuracy of consensus after k iterations, i.e.

$$\frac{1}{M} \sum_{m=1}^M \|z_m^k - z^k\|^2 \leq \varepsilon_0^2.$$

We introduce the following notation:

$$g_m^k = F_m(z_m^k), \quad g_m^{k+1/2} = F_m(z_m^{k+1/2}),$$

and

$$\begin{aligned} z^k &= \frac{1}{M} \sum_{m=1}^M z_m^k, & z^{k+1/2} &= \frac{1}{M} \sum_{m=1}^M z_m^{k+1/2}, & g^k &= \frac{1}{M} \sum_{m=1}^M g_m^k, & g^{k+1/2} &= \frac{1}{M} \sum_{m=1}^M g_m^{k+1/2}, \\ \hat{z}^k &= \frac{1}{M} \sum_{m=1}^M \hat{z}_m^k, & \hat{z}^{k+1/2} &= \frac{1}{M} \sum_{m=1}^M \hat{z}_m^{k+1/2}, & \tilde{z}^k &= \frac{1}{M} \sum_{m=1}^M \tilde{z}_m^k, & \tilde{z}^{k+1/2} &= \frac{1}{M} \sum_{m=1}^M \tilde{z}_m^{k+1/2}. \end{aligned}$$

We have

$$\frac{1}{M} \sum_{m=1}^M \|g_m^k - g^k\|^2 \leq \frac{1}{M} \sum_{m=1}^M \|g_m^k\|^2.$$

Let $T = 2\chi \left(1 - \frac{1}{\sqrt{\chi}}\right)^H$. Then

$$\begin{aligned} \frac{1}{M} \sum_{m=1}^M \|\hat{z}_m^{k+1/2} - \tilde{z}_m^{k+1/2}\|^2 &\leq \frac{T}{M} \sum_{m=1}^M \|\hat{z}_m^{k+1/2} - \tilde{z}_m^{k+1/2}\|^2 = \frac{T}{M} \sum_{m=1}^M \|z_m^k - \gamma g_m^k - z^k - \gamma g^k\|^2 \\ &\leq \frac{2T}{M} \sum_{m=1}^M \|z_m^k - z^k\|^2 + \frac{2T\gamma^2}{M} \sum_{m=1}^M \|g_m^k - g^k\|^2 \\ &\leq \frac{2T}{M} \sum_{m=1}^M \|z_m^k - z^k\|^2 + \frac{2T\gamma^2}{M} \sum_{m=1}^M \|g_m^k\|^2 = 2T\varepsilon_0^2 + \frac{2T\gamma^2}{M} \sum_{m=1}^M \|g_m^k\|^2. \end{aligned}$$

On the other side

$$\begin{aligned} \frac{1}{M} \sum_{m=1}^M \|g_m^k\|^2 &= \frac{1}{M} \sum_{m=1}^M \|F_m(z_m^k)\|^2 \leq \frac{2}{M} \sum_{m=1}^M \|F_m(z_m^k) - F_m(z^*)\|^2 + \frac{2}{M} \sum_{m=1}^M \|F_m(z^*)\|^2 \\ &\leq 2L_{\max}^2 \|z_m^k - z^*\|^2 + \frac{2}{M} \sum_{m=1}^M \|F_m(z^*)\|^2 \leq 2L_{\max}^2 \|z_m^0 - z^*\|^2 + \frac{2}{M} \sum_{m=1}^M \|F_m(z^*)\|^2. \end{aligned}$$

Let $Q^2 = \frac{1}{M} \sum_{m=1}^M \|F_m(z^*)\|^2$, then we have

$$\begin{aligned} \frac{1}{M} \sum_{m=1}^M \|z_m^{k+1/2} - z^{k+1/2}\|^2 &= \frac{1}{M} \sum_{m=1}^M \|\text{proj}_{\mathcal{Z}} \hat{z}_m^{k+1/2} - \text{proj}_{\mathcal{Z}} \tilde{z}_m^{k+1/2}\|^2 \leq \frac{1}{M} \sum_{m=1}^M \|\hat{z}_m^{k+1/2} - \tilde{z}_m^{k+1/2}\|^2 \\ &\leq 2T \left(\varepsilon_0^2 + 2\gamma^2 (L_{\max}^2 \|z_m^0 - z^*\|^2 + Q^2) \right) \leq 2T \left(\varepsilon_0^2 + \frac{1}{8} \|z_m^0 - z^*\|^2 + \frac{Q^2}{8L_{\max}^2} \right) \\ &= \chi \left(1 - \frac{1}{\sqrt{\chi}}\right)^H \left(4\varepsilon_0^2 + \frac{1}{2} \|z^0 - z^*\|^2 + \frac{Q^2}{2L_{\max}^2} \right). \end{aligned}$$

If we take

$$H \geq \sqrt{\chi} \log \left[\chi \left(4 + \frac{\frac{1}{2} \|z^0 - z^*\|^2 + \frac{Q^2}{2L_{\max}^2}}{\varepsilon_0^2} \right) \right],$$

then

$$\frac{1}{M} \sum_{m=1}^M \|z_m^{k+1/2} - z^{k+1/2}\|^2 \leq \varepsilon_0^2.$$

Analogously, we get the same estimate for H to ensure that

$$\frac{1}{M} \sum_{m=1}^M \|z_m^{k+1} - z^{k+1}\|^2 \leq \varepsilon_0^2.$$

Hence, to achieve accuracy ε_0 , we need to perform

$$H = \mathcal{O} \left(\sqrt{\chi} \log \left[\chi \left(4 + \frac{\frac{1}{2} \|z^0 - z^*\|^2 + \frac{Q^2}{2L_{\max}^2}}{\varepsilon_0^2} \right) \right] \right) \text{ communications.}$$

□

Theorem 2. (From Theorem 6 in [3]) Let $\{z_m^k\}_{k \geq 0}$ denote the iterates of Algorithm 2 for solving problem 1. Let Assumptions 1, 3 be satisfied. Then if $\gamma \leq \frac{1}{4L_{\max}}$, we have the following estimates

$$\|z^N - z^*\| = \mathcal{O} \left(\|z^0 - z^*\|^2 \exp \left(-\frac{\mu K}{8L \cdot H} \right) \right).$$

Corollary 2.1. In the setting of Theorem 1 and Theorem 2, if $H = \mathcal{O}(\sqrt{\chi} \log \chi \log(1/\varepsilon))$, then the number of communication rounds required for Algorithm 2 to obtain ε -solution is upper bounded by

$$\mathcal{O} \left(\sqrt{\chi} \log \chi \frac{L}{\mu} \log^2 \frac{1}{\varepsilon} \right),$$

and the number of local computations on each device is upper bounded by

$$\mathcal{O} \left(\frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$

Consensus for networks with Markovian changes

This subsection is devoted to slowly time-varying graphs with random changes satisfying Markovian law. At each iteration, several randomly chosen edges may appear or disappear. The choice of edges depends only on the current graph topology, so the sequence of graphs is a Markovian process.

Let introduce some requirements for time-varying network with Markovian changes.

Assumption 4. The communication network satisfy the following conditions

1. $\{W^k\}_{k=0}^{\infty}$ is a stationary Markov chain on (W_G, W_σ) , where W_G is the set of all possible gossip matrices for the network and W_σ is the σ -field on W_G and the chain $\{W^k\}_{k=0}^{\infty}$ has a Markov kernel Q and a unique stationary distribution π .
2. Q is uniformly geometrically ergodic with mixing time $\tau \in \mathbb{N}$, i.e., for every $m \in \mathbb{N}$,

$$\Delta(Q^m) = \sup_{W, W' \in W_G} (1/2) \|Q^m(W, \cdot) - Q^m(W', \cdot)\|_{TV} \leq (1/4)^{\lfloor m/\tau \rfloor}.$$

3. For all $k \in \mathbb{N} \cup \{0\}$, it holds $\mathbb{E}_\pi[W(q)] = \bar{W}$ and the \bar{W} satisfies Assumption 2.

Denote $\lambda_{\max} = \lambda_{\max}(\bar{W})$, $\lambda_{\min}^+ = \lambda_{\min}^+(\bar{W})$, $\chi = \frac{\lambda_{\max}}{\lambda_{\min}^+}$.

4. For any graph \mathcal{G} that can appear in the network it holds:

$$\|W(\mathcal{G}) - \bar{W}\| \leq \rho.$$

Consider the consensus search problem:

$$\begin{aligned} \min_{\mathbf{z} \in \mathbb{R}^{Mn_z}} & \left[r(\mathbf{z}) = \left\| \left(\sqrt{\bar{W}} \otimes \mathbf{I}_{n_z} \right) \mathbf{z} \right\|^2 \right] \\ \text{s.t.} & \sum_{m=1}^M z_m = \sum_{m=1}^M z_m^0 \end{aligned} \quad (5)$$

where $\mathbf{z} = (z_1^T, \dots, z_M^T)^T$.

Theorem 3. (Theorem 1 from [11]) Let Assumptions 4 hold. Let problem (5) be solved by Algorithm 3. Then for any $b \in \mathbb{N}$,

$$\gamma \in \left(0; \min \left\{ \frac{3}{4\lambda_{\max}}; \frac{\lambda_{\min}^3}{[1800\rho^2(\tau b^{-1} + \tau^2 b^{-2})]^2} \right\} \right),$$

and $\beta, \theta, \eta, p, S, B$ satisfying

$$p = \frac{1}{4}, \beta = \sqrt{\frac{4p^2 \mu \gamma}{3}}, \eta = \frac{3\beta}{p\mu\gamma} = \sqrt{\frac{12}{\mu\gamma}}, \theta = \frac{p\eta^{-1} - 1}{\beta p \eta^{-1} - 1}$$

Algorithm 3 Accelerated consensus over graphs with Markovian changes (ACOGWMC)

Require: stepsize $\gamma > 0$, momentums θ, η, β, p , number of iterations N , batchsize limit S .

Choose $z_f^0 = z^0, T^0 = 0$, set the same random seed for generating $\{J_k\}$ on all devices

for $k = 0, 1, \dots, N - 1$ **do**

$$z_g^k = \theta z_f^k + (1 - \theta) z^k$$

Sample $J_k \sim \text{Geom}(1/2)$

Send z_g^k to neighbors in the networks $\{\mathcal{G}^{T^k+i}\}_{i=1}^{2^{J_k}B}$

$$\text{Compute } g^k = g_0^k + \begin{cases} 2^{J_k} (g_{J_k}^k - g_{J_k-1}^k), & \text{if } 2^{J_k} \leq S \\ 0, & \text{otherwise} \end{cases}$$

with $g_j^k = 2^{-j} B^{-1} \sum_{i=1}^{2^j B} W^{T^k+i} z_g^k$

$$z_f^{k+1} = z_g^k - p\gamma g^k$$

$$z^{k+1} = \eta z_f^{k+1} + (p - \eta) z_f^k + (1 - p)(1 - \beta) z^k + (1 - p)\beta z_g^k$$

$$T^{k+1} = T^k + 2^{J_k} B$$

end for

$$S = \max \left\{ 2; \sqrt{\frac{1}{4} \left(1 + \frac{2}{\beta} \right)} \right\}, \quad B = \lceil b \log_2 S \rceil,$$

it holds that

$$\begin{aligned} & \mathbb{E} \left[\|z^N - z^*\|^2 + \frac{24}{\lambda_{\min}} (r(z_f^N) - r(z^*)) \right] \\ &= \mathcal{O} \left(\exp \left(-N \sqrt{\frac{p^2 \lambda_{\min} \gamma}{3}} \right) \left[\|z^0 - z^*\|^2 + \frac{24}{\lambda_{\min}} (r(z^0) - r(z^*)) \right] \right), \end{aligned}$$

where $z_m^* = \frac{1}{M} \sum_{i=1}^M z_i$ for $m = 1, \dots, M$.

Corollary 3.1. In the setting of Theorem 3, if $b = \tau$ and $\gamma \simeq \min \left\{ \frac{1}{\lambda_{\max}}; \frac{\lambda_{\min}^3}{\rho^4} \right\}$, then in order to achieve ε -approximate solution (in terms of $\mathbb{E}[\|z - z^*\|^2] \lesssim \varepsilon$) it takes

$$\tilde{\mathcal{O}} \left(\tau \left[\sqrt{\chi} + \frac{\rho^2}{\lambda_{\min}^2} \log \frac{1}{\varepsilon} \right] \right) \text{ communications.}$$

Algorithm 4 Time-Varying Decentralized Extra Step Method with Markovian changes

Require: Step size $\gamma \leq \frac{1}{4L}$, number of `AccGossipNonRecoverable` steps H , number of iterations N .

Choose $(x^0, y^0) = z^0 \in \mathcal{Z}, z_m^0 = z^0$.

for $k = 0, 1, 2, \dots, N$ **do**

Each machine m computes $\hat{z}_m^{k+1/2} = z_m^k - \gamma \cdot F_m(z_m^k)$

Communication: $\hat{z}_1^{k+1/2}, \dots, \hat{z}_M^{k+1/2} = \text{ACOGWMC}(\hat{z}_1^{k+1/2}, \dots, \hat{z}_M^{k+1/2}, H)$

Each machine m computes $z_m^{k+1/2} = \text{proj}_{\mathcal{Z}}(\hat{z}_m^{k+1/2})$

Each machine m computes $\hat{z}_m^{k+1} = z_m^{k+1/2} - \gamma \cdot F_m(z_m^{k+1/2})$

Communication: $\hat{z}_1^{k+1}, \dots, \hat{z}_M^{k+1} = \text{ACOGWMC}(\hat{z}_1^{k+1}, \dots, \hat{z}_M^{k+1}, H)$

Each machine m computes $z_m^{k+1} = \text{proj}_{\mathcal{Z}}(\hat{z}_m^{k+1})$

end for

Theorem 4. Let Assumptions 1, 4 hold. Let problem (1) be solved by Algorithm 4. Then, if $\gamma \leq \frac{1}{4L_{\max}}$ and $H = \mathcal{O} \left(\tau \left[\sqrt{\chi} + \frac{\rho^2}{\lambda_{\min}^2} \log \frac{1}{\varepsilon} \right] \right)$ it holds that to achieve ε -solution (in terms of $\mathbb{E}[f(z) - f(z^*)] \lesssim \varepsilon$) it takes

$$\tilde{\mathcal{O}} \left(\tau \left[\sqrt{\chi} + \frac{\rho^2}{\lambda_{\min}^2} \right] \frac{L}{\mu} \log^2 \frac{1}{\varepsilon} \right) \text{ communications and}$$

$$\mathcal{O} \left(\frac{L}{\mu} \log \frac{1}{\varepsilon} \right) \text{ local computations on each node.}$$

4. LOWER BOUNDS

The previous section was devoted to applying new consensus algorithms for specific types of time-varying networks: networks with connected skeleton and graphs changing according to Markovian law. In this section, we show that without specific constraints on the network change, such as connected skeleton or Markovian changes, acceleration is not obtained on generalised slowly changing graphs, even if the constraints on the rate of change are very high. That is, under constant network constraints, when at most a constant number of edges in the network per iteration is changed, the worst-case dependence on χ cannot be improved in comparison to arbitrary changing networks. In particular, we show that it sufficient to change only two edges per iteration.

We start with the definition of black-box procedure class of algorithms on which we will evaluate the lower bound.

Definition 1. An algorithm with T local iterations and K communication rounds that satisfies following properties is called a *black-box procedure*, denoted by $\mathbf{BBP}(T, K)$.

Each node m maintains a local memory with \mathcal{M}_m^x and \mathcal{M}_m^y for the x - and y -variables, which are initialized as $\mathcal{M}_m^x = \mathcal{M}_m^y = \{0\}$. \mathcal{M}_m^x and \mathcal{M}_m^y are updated as follows:

- **Local computation:** Each node m computes and adds to its \mathcal{M}_m^x and \mathcal{M}_m^y a finite number of points x, y , each satisfying

$$x \in \text{span}\{x', \nabla_x f_m(x'', y'')\}, \quad y \in \text{span}\{y', \nabla_y f_m(x'', y'')\},$$

for given $x', x'' \in \mathcal{M}_m^x$ and $y', y'' \in \mathcal{M}_m^y$.

- **Communication:** \mathcal{M}_m^x and \mathcal{M}_m^y are updated according to

$$\mathcal{M}_m^x := \text{span} \left\{ \bigcup_{(i,m) \in \mathcal{E}^k} \mathcal{M}_i^x \right\}, \quad \mathcal{M}_m^y := \text{span} \left\{ \bigcup_{(i,m) \in \mathcal{E}^k} \mathcal{M}_i^y \right\},$$

where $\mathcal{G}^k = (\mathcal{V}, \mathcal{E}^k)$ is the current state of network.

- **Output:** The final global output at current moment of time is calculated as:

$$\hat{x} \in \text{span} \left\{ \bigcup_{m=1}^M \mathcal{M}_m^x \right\}, \quad \hat{y} \in \text{span} \left\{ \bigcup_{m=1}^M \mathcal{M}_m^y \right\}.$$

To estimate the lower bound of distributed saddle point problem (1), we need to provide a "bad function" and a "bad sequence of graphs" such that any black-box procedure cannot solve it using less than a given number of rounds. Using the time-varying network from [11] and the objective function from [20] we can prove the following theorem.

Theorem 5. For any $L > \mu > 0$ and any $\chi \geq 1$, there exists a decentralized distributed saddle point problem on $\mathcal{X} \times \mathcal{Y} = \mathbb{R}^n \times \mathbb{R}^n$ (where n is sufficiently large) with $x^*, y^* \neq 0$, a sequence of graphs $\{\mathcal{G}^k = (\mathcal{V}, \mathcal{E}^k)\}_{k=0}^\infty$, where consecutive graphs differ in no more than two edges, and a sequence of corresponding gossip matrices $\{W^k\}_{k=0}^\infty$ with characteristic number χ , such that for any output \hat{x}, \hat{y} after K communication rounds of any \mathbf{BBP} , the following estimate hold:

$$\|\hat{x} - x^*\|^2 + \|\hat{y} - y^*\|^2 \geq \Omega \left(\exp \left(-\frac{32\mu}{L - \mu} \cdot \frac{K}{\chi} \right) \|y_0 - y^*\|^2 \right)$$

Proof. Let us introduce the graph $T_{a,b}$ from [11]. This graph contains two partitions: left and right, these partitions consist of a and b vertices respectively. Each vertex in partitions is connected to its root. These roots are connected to the another vertex called central root.

Consider the network with $|\mathcal{V}| = 2d + 3$ nodes ($d \geq 2$). We select a node to be the central root and two other nodes to be the left and right roots. Central root can be changed over time but partition roots are fixed. We also select $\lfloor \frac{d}{2} \rfloor$ fixed vertices for each partition and denote them by \mathcal{V}_1 (left side) and \mathcal{V}_2 (right side). At each communication round k , the graph \mathcal{G}^k has the form $T_{a,b}$, where $a + b = 2d$ and $a, b \geq \lfloor \frac{d}{2} \rfloor$.

The communication network changes alternatively in two phases. The first phase starts with graph $T_{2d - \lfloor \frac{d}{2} \rfloor, \lfloor \frac{d}{2} \rfloor}$ and ends with graph $T_{\lfloor \frac{d}{2} \rfloor, 2d - \lfloor \frac{d}{2} \rfloor}$. At each iteration, the central root is moved to the right partition and one vertex from left partition, but not in \mathcal{V}_1 , becomes the central root. The second phase has the same procedure, but from right to left.

We modify the objective function from section B.1 in [3] based on our graph type:

Using result from the proof of Proposition 3.6 in [20], we have:

$$\ln(q) \geq \frac{q-1}{q} = \frac{2}{1 - \sqrt{\frac{L^2}{\mu^2} + 1}} \geq \frac{2}{1 - \frac{L}{\mu}} = \frac{-2\mu}{L - \mu}.$$

For each graph in our sequence we map a weighted Laplacian from Lemma 8 in [11], so $\chi \leq 8d$.

Hence

$$\ln(q) \cdot \frac{2K}{d} \geq \frac{-4\mu}{L - \mu} \cdot \frac{K}{d} \geq \frac{-32\mu}{L - \mu} \cdot \frac{K}{\chi}.$$

We get

$$q^{\frac{2K}{d}} = \exp\left(\ln(q) \cdot \frac{2K}{d}\right) \geq \exp\left(\frac{-32\mu}{L - \mu} \cdot \frac{K}{\chi}\right).$$

So we obtain

$$\|\hat{x} - x^*\|^2 + \|\hat{y} - y^*\|^2 = \Omega\left(\exp\left(-\frac{32\mu}{L - \mu} \cdot \frac{K}{\chi}\right) \|y_0 - y^*\|^2\right).$$

□

Corollary 5.1. In the setting of Theorem 1, the number of communication rounds required to obtain a ε -solution is lower bounded by

$$\Omega\left(\chi \frac{L}{\mu} \cdot \log\left(\frac{\|y^*\|^2}{\varepsilon}\right)\right),$$

and the number of local calculations on each node is lower bounded by:

$$\Omega\left(\frac{L}{\mu} \cdot \log\left(\frac{\|y^*\|^2}{\varepsilon}\right)\right).$$

5. CONCLUSION

In this paper, we studied min-max optimization over slowly time-varying graphs. We showed that if the graph changes in an adversarial manner and a constant number of edges is changed at each iteration, then lower complexity bounds coincide with those for arbitrary changing networks. Moreover, we showed that for particular time-varying graphs – networks with connected skeleton and networks with Markovian changes – acceleration of communication procedures is possible. We proposed the corresponding algorithms for saddle point problems for these two classes of problems.

The research was supported by Russian Science Foundation (project No. 23-11-00229), <https://rscf.ru/en/project/23-11-00229/>.

REFERENCES

- [1] J. A. Bazerque and G. B. Giannakis. Distributed spectrum sensing for cognitive radio networks by exploiting sparsity. *IEEE Transactions on Signal Processing*, 58(3):1847–1862, 2009.
- [2] A. Beznosikov, A. Rogozin, D. Kovalev, and A. Gasnikov. Near-optimal decentralized algorithms for saddle point problems over time-varying networks. In *International Conference on Optimization and Applications*, pages 246–257. Springer, 2021.
- [3] A. Beznosikov, V. Samokhin, and A. Gasnikov. Distributed saddle-point problems: Lower bounds, near-optimal and robust algorithms. *arXiv preprint arXiv:2010.13112*, 2020.
- [4] A. Beznosikov, V. Samokhin, and A. Gasnikov. Distributed saddle-point problems: Lower bounds, optimal algorithms and federated gans. *arXiv preprint arXiv:2010.13112*, 2021.
- [5] P. A. Forero, A. Cano, and G. B. Giannakis. Consensus-based distributed support vector machines. *Journal of Machine Learning Research*, 11(5), 2010.
- [6] L. Gan, U. Topcu, and S. H. Low. Optimal decentralized protocol for electric vehicle charging. *IEEE Transactions on Power Systems*, 28(2):940–951, 2012.
- [7] D. Kovalev, A. Beznosikov, A. Sadiev, M. Persiiyanov, P. Richtárik, and A. Gasnikov. Optimal algorithms for decentralized stochastic variational inequalities. *Advances in Neural Information Processing Systems*, 35:31073–31088, 2022.
- [8] D. Kovalev, E. Gasanov, A. Gasnikov, and P. Richtarik. Lower bounds and optimal algorithms for smooth and strongly convex decentralized optimization over time-varying networks. *Advances in Neural Information Processing Systems*, 34, 2021.

- [9] D. Kovalev, A. Salim, and P. Richtárik. Optimal and practical algorithms for smooth and strongly convex decentralized optimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- [10] D. Kovalev, E. Shulgin, P. Richtárik, A. Rogozin, and A. Gasnikov. Adom: Accelerated decentralized optimization method for time-varying networks. *arXiv preprint arXiv:2102.09234*, 2021.
- [11] D. Metev, A. Beznosikov, A. Rogozin, A. Gasnikov, and A. Proskurnikov. Decentralized optimization over slowly time-varying graphs: Algorithms and lower bounds. *arXiv preprint arXiv:2307.12562*, 2023.
- [12] D. Metev, A. Rogozin, D. Kovalev, and A. Gasnikov. Is consensus acceleration possible in decentralized optimization over slowly time-varying networks? In *International Conference on Machine Learning*, pages 24532–24554. PMLR, 2023.
- [13] A. Nedic. Distributed gradient methods for convex machine learning problems in networks: Distributed optimization. *IEEE Signal Processing Magazine*, 37(3):92–101, 2020.
- [14] A. Nedić, A. Olshevsky, and C. A. Uribe. Fast convergence rates for distributed non-bayesian learning. *IEEE Transactions on Automatic Control*, 62(11):5538–5553, 2017.
- [15] A. Nedic and A. Ozdaglar. Cooperative distributed multi-agent optimization. *Convex optimization in signal processing and communications*, 340, 2010.
- [16] M. Rabbat and R. Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pages 20–27, 2004.
- [17] S. S. Ram, V. V. Veeravalli, and A. Nedic. Distributed non-autonomous power control through distributed convex optimization. In *IEEE INFOCOM 2009*, pages 3001–3005. IEEE, 2009.
- [18] W. Ren and R. W. Beard. *Distributed consensus in multi-vehicle cooperative control*, volume 27. Springer, 2008.
- [19] K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3027–3036. JMLR. org, 2017.
- [20] J. Zhang, M. Hong, and S. Zhang. On lower iteration complexity bounds for the saddle point problems. *arXiv preprint arXiv:1912.07481*, 2019.