

Simultaneous Perturbation Stochastic Approximation-based Consensus for Tracking under Unknown—but-Bounded Disturbances

Oleg Granichin, *Senior Member, IEEE*, Victoria Erofeeva, Yury Ivanskiy,
and Yuming Jiang, *Senior Member, IEEE*

Abstract—We consider a setup where a distributed set of sensors working cooperatively can estimate an unknown signal of interest, whereas any individual sensor cannot fulfil the task due to lack of necessary information diversity. This paper deals with these kinds of estimation and tracking problems and focuses on a class of simultaneous perturbation stochastic approximation (SPSA)-based consensus algorithms for the cases when the corrupted observations of sensors are transmitted between sensors with communication noise and the communication protocol has to satisfy a prespecified cost constraints on the network topology. Sufficient conditions are introduced to guarantee the stability of estimates obtained in this way, without resorting to commonly used but stringent conventional statistical assumptions about the observation noise such as randomness, independence, and zero mean. We derive an upper bound of the mean-square error of the estimates in the problem of unknown time-varying parameters tracking under unknown—but-bounded observation errors and noisy communication channels. The result is illustrated by a practical application to the multi-sensor multi-target tracking problem.

Index Terms—Distributed tracking, multi-agent networks, consensus algorithm, simultaneous perturbation stochastic approximation, SPSA, randomized algorithm, arbitrary noise, unknown—but-bounded disturbances, stochastic stability, tracking performance.

I. INTRODUCTION

Distributed cooperative control of networked systems has been investigated and numerous potential applications to complex manufacturing, energy and social systems have been developed [1]–[3] over the past few decades. One of the fundamental concepts in multi-agent cooperative control is *consensus*. This approach aims to find an agreement between all agents in a network regarding a common value using only local information and communicating among neighboring agents.

The goal of distributed optimization is usually to find the minimum of some loss function $\bar{F}(\mathbf{x}) = \sum_{i=1}^n F^i(\mathbf{x})$ via interaction between agents. Here, $\mathbf{x} \in \mathbb{R}^d$ and $F^i(\mathbf{x}) :$

$\mathbb{R}^d \rightarrow \mathbb{R}$ is the loss function of agent i , which is typically known only to the agent itself. Studies of consensus and distributed optimization algorithms began from the 1970-80s [4], [5]. Distributed asynchronous stochastic approximation algorithms were studied in [6]. To date, there exist a number of approaches for the case when functions $F^i(\mathbf{x})$ are convex. In particular, the Alternating Direction Method of Multipliers [7], [8], as well as the subgradient method [9], [10] were proposed. For non-convex tasks, the works [11], [12] develop a large class of distributed algorithms based on various “functional-surrogate units”. The distributed tracking problem is considered when the estimated parameters vary over time.

Recently, for large-scale systems consisting of many individuals (components, targets), a distributed optimization is often used to estimate the unknown parameters which minimize a loss function, based on the information obtained by distributed sensors. So-called *multitarget-multisensor tracking problems* have been widely studied in many practical applications such as adaptive mobile networks, cognitive radio systems, target localization in biological networks, fish schooling, bee swarming, and bird flight (see, e.g., [13], [14]). It is well known that distributed tracking algorithms have some significant advantages over the centralized ones or the fusion methods. Centralized algorithms usually require the distributed sensor network to transmit the whole system information into a fusion center to estimate the unknown signals. This leads to the necessity of strong communication capabilities over sensor networks which is hard to provide in many practical situations when sensors may only have the capability to exchange information locally between their neighbors. An alternative approach for multitarget-multisensor tracking problems assumes only local interaction between sensors without the governing data fusion center. A detailed literature overview of the recent advances in the stability analysis of a consensus-based least squares algorithms is performed in [15] for distributed estimation and tracking problems.

The maximum likelihood estimator and stochastic approximation (SA) algorithms with decreasing to zero step-size are actively used to optimize mean-risk functionals. In gradient-free conventional stochastic approximation algorithms, two measurements are used to approximate each component of the vector-gradient of the cost function (implying $2d$ measurements for the d -dimension state space). *Simultaneous perturbation stochastic approximation* (SPSA) was proposed by Spall [16]. It can be used to solve optimization problems

The theoretical research in Sections I-VI of this work was supported by the Russian Fund for Basic Research (project no. 20-01-00619). The obtaining of experimental results in Section VII was supported by Russian Science Foundation (project no. 19-71-10012).

O. Granichin, V. Erofeeva and Y. Ivanskiy are with Saint Petersburg State University, 198504, Universitetskii pr. 28, St. Petersburg, Russia. e-mail: o.granichin@spbu.ru, victoria@grenka.net, ivanskiy.yuriy@gmail.com.

Y. Jiang is with Department of Information Security and Communication Technology, Norwegian University of Science and Technology (NTNU), NO-7491, Trondheim, Norway, e-mail: jiang@ntnu.no.

in the case when it is difficult or impossible to obtain a gradient of the loss function with respect to the parameters being optimized. In any multidimensional case ($d \gg 1$), SPSA requires only two measurements of a loss function at each iteration. The current estimate is updated along a randomly chosen direction Δ with ± 1 Bernoulli distributed components.

Traditionally, a stochastic optimization problem under uncertainties focused on finding a set of system parameters that deliver a minimum (or maximum) value to a certain mean-risk functional. In practical applications, these parameters may also vary over time. The problem of tracking changes in system parameters is considered in [17]–[19]. In this paper, such a problem is called the *minimum-point tracking of a nonstationary mean-risk functional*. In centralized (non-distributed) cases, SPSA-like algorithms for parameter tracking problems were considered in [20]–[22]. The stochastic approximation method with a constant step-size has also been used in [23] to achieve the approximate mean-squared consensus in multi-agent systems operating under noisy measurement conditions.

Contributions. In the case of differentiable time-varying loss functions and almost arbitrary external bounded noise, an upper bound of the mean square estimation error was derived in [20] for estimates of the SPSA type algorithms with constant step-size. This upper bound may be sufficiently small compared to the significant level of noise when the rate of change of parameters is low enough. One of the main conditions is a strong convexity property of the minimized mean-risk functional. In this paper we weaken this assumption by combining SPSA with the consensus algorithm from [23]. We propose a new SPSA-based consensus algorithm for distributed tracking under *unknown-but-bounded disturbances*. The preliminary concept of this paper is presented in [24]. In many practical applications, the network processes the data under certain constraints, and the data transmission is accompanied by noise. In this paper, compared with [24], we consider such noisy data transmission and a communication protocol with prespecified cost constraints on the network topology. Also, we study a more general type of simultaneous perturbation and we choose the current points of observations in a more general manner. We obtain an upper bound of the mean square error of estimates of unknown time-varying parameters tracking. Communication cost constraints are satisfied by exploiting a specific intentionally randomized topology of the network communication graph.

The paper is organized as follows. The preliminary information regarding concepts of the graph theory and network topology is given in Section II. A formal problem setting of a distributed non-constrained time-varying mean-risk optimization with noisy local communications is given in Section III. The main result including assumptions and the SPSA-based consensus algorithm for tracking is presented in Section IV. In Section V, the efficiency of the proposed algorithm is illustrated through the numerical simulation.

II. PRELIMINARIES

Let (Ω, \mathcal{F}, P) be the underlying probability space corresponding to sample space Ω , set of all events \mathcal{F} , and probability measure P . \mathbb{E} denotes mathematical expectation.

A. Concepts of Graph Theory

Given a network consisting of n sensors. Let the interaction between sensors be described by the directed graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where $\mathcal{N} = \{1, \dots, n\}$ is a set of vertices and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ is a set of edges. A subgraph of \mathcal{G} is a graph $\bar{\mathcal{G}} = (\mathcal{N}_{\bar{\mathcal{G}}}, \mathcal{E}_{\bar{\mathcal{G}}})$, where $\mathcal{N}_{\bar{\mathcal{G}}} \subseteq \mathcal{N}$ and $\mathcal{E}_{\bar{\mathcal{G}}} \subseteq \mathcal{E}$. Denote by $i \in \mathcal{N}$ an identifier of i -th sensor and $(j, i) \in \mathcal{E}$ if there is a directed edge from sensor j to sensor i . The latter means that sensor j is able to transmit data to sensor i . For a sensor $i \in \mathcal{N}$, the set of *neighbors* is defined as $\mathcal{N}^i = \{j \in \mathcal{N} : (j, i) \in \mathcal{E}\}$. The *in-degree* of $i \in \mathcal{N}$ equals $|\mathcal{N}^i|$. Here and after, $|\cdot|$ is the cardinality of a set, and the identifier of i -th sensor is used as a superscript and not as an exponent.

Let $c^{i,j} > 0$ be the weight associated with the edge $(j, i) \in \mathcal{E}$, and $c^{i,j} = 0$ whenever $(j, i) \notin \mathcal{E}$. Let $C = [c^{i,j}]$, be the *weighted adjacency matrix*, or simply *connectivity matrix*. Denote by $\mathcal{G}_C = (\mathcal{N}_C, \mathcal{E}_C)$ the weighted directed graph, where $\mathcal{N}_C \equiv \mathcal{N}$ and $\mathcal{E}_C \equiv \mathcal{E}$. We assume that weight $c^{i,j}$ is the cost of data transmission through the edge $(j, i) \in \mathcal{E}_C$. The *weighted in-degree* of $i \in \mathcal{N}_C$ is defined as $\deg_i^+(C) = \sum_{j=1}^n c^{i,j}$, the maximum in-degree among all nodes contained in the graph \mathcal{G}_C as $\deg_{\max}^+(C)$, and the diagonal matrix as $\mathcal{D}(C) = \text{diag}_n(\text{col}(\deg_1^+(C), \dots, \deg_n^+(C)))$. Then, $\mathcal{L}(C) = \mathcal{D}(C) - C$ is the *Laplacian* of graph \mathcal{G}_C .

Definition 1. A directed graph \mathcal{G} is said to be strongly connected if for every pair of nodes $j, i \in \mathcal{N}$, there exists a path of directed edges that goes from j to i .

Denote the eigenvalues of Laplacian $\mathcal{L}(C)$ by $\lambda_1, \dots, \lambda_n$ and arrange them in ascending order of real parts: $0 \leq \text{Re}(\lambda_1) \leq \text{Re}(\lambda_2) \leq \dots \leq \text{Re}(\lambda_n)$. It is known, that if the graph is strongly connected then $\lambda_1 = 0$ and all other eigenvalues of \mathcal{L} are in the open right half of the complex plane (see, e.g., [3]). The eigenvalue of matrix C with maximum absolute magnitude is defined as $\lambda_{\max}(C)$.

B. Network Topology Constraints

In practice, we have constraints on the bandwidth of communication channels, network response time, hardware and financial requirements, *etc.* In this paper, we associate these constraints with matrix C , which characterizes the cost of data transmission in the network. In many practical applications, we may represent cost constraints of sensor $i \in \mathcal{N}$ as a predefined upper bound Q : $\deg_i^+(C) \leq Q$. This bound may be thought of as the total cost of communication with neighbors of sensor i . To satisfy this constraint, we may generate at each time instant t subgraph $\mathcal{G}_{B_t} \subset \mathcal{G}_C$ associated with the weighted connectivity matrix B_t such that $\deg_i^+(B_t) \leq Q$. Obviously, the cost constraint $\deg_i^+(B_t) \leq Q$ may not be satisfied for given $B_t = C$ and Q , e.g. when $n = 6$, \mathcal{G}_C is the complete graph with $c^{i,j} = 1$, $i \neq j$, $c^{i,i} = 0$, and $Q < 5$. One possible solution is to use a randomized topology, when we drop $5 - Q$ edges randomly. Such randomized strategy for $Q = 1$ is similar to the scheme used in *gossip* algorithms [25]. Moreover, random subgraphs naturally arise in many practical applications.

Next, we consider a communication protocol needed to satisfy a predefined averaged cost constraint.

Definition 2. Random subgraph \mathcal{G}_{B_t} satisfies the averaged cost constraints with level Q if

$$\mathbb{E} \deg_{\max}^+(B_t) \leq Q. \quad (1)$$

In the example considered above we are able to satisfy averaged cost constraints if each sensor i randomly selects its neighbors out of all $j \in \mathcal{N}^i$ with probability $\frac{Q}{\deg_i^+(B_t)} = 0.2Q$ at each time instant t .

III. DISTRIBUTED TRACKING

A. Non-stationary Mean-risk Functional

Let Ξ be a set, $\{f_{\xi}^i(\theta)\}_{\xi \in \Xi}$ be a family of differentiable functions: $\forall i \in \mathcal{N} \ f_{\xi}^i(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$. We assume that parameter θ cannot be directly measured. Hence, we introduce a sequence of measurement points $\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, i \in \mathcal{N}$ chosen according to an observation plan. The values y_1^i, y_2^i, \dots of the functions $f_{\xi_t}^i(\cdot)$ are observable at every time instant $t = 1, 2, \dots$ with additive external *unknown-but-bounded* noise v_t^i

$$y_t^i = f_{\xi_t}^i(\mathbf{x}_t^i) + v_t^i, \quad (2)$$

where $\{\xi_t\}$, $\xi_t \in \Xi$, is a non-controllable deterministic (e.g., $\Xi = \mathbb{N}$ and $\xi_t = t$) or random sequence. In the latter case we assume that a probability distribution of ξ_t exists and may be known or unknown.

Let \mathcal{F}_{t-1} be the σ -algebra of all probabilistic events which happened up to time instant t , $\mathbb{E}_{\mathcal{F}_{t-1}}$ denotes the conditional mathematical expectation with respect to the σ -algebra \mathcal{F}_{t-1} . We consider an optimization problem in which the cost function $\bar{F}_t(\theta)$ is expressed as the sum of local contributions $F_t^i(\theta) = \mathbb{E}_{\mathcal{F}_{t-1}} f_{\xi_t}^i(\theta)$ and all of them depend on a common optimization variable θ . Moreover, minimizer θ of $\bar{F}_t(\theta)$ may vary over time. Formally, the *non-stationary mean-risk optimization problem* is as follows: estimate the time-varying minimum point θ_t of the distributed function

$$\bar{F}_t(\theta) = \sum_{i \in \mathcal{N}} F_t^i(\theta) = \mathbb{E}_{\mathcal{F}_{t-1}} \sum_{i \in \mathcal{N}} f_{\xi_t}^i(\theta) \rightarrow \min_{\theta}. \quad (3)$$

More precisely, the problem is to obtain an estimate $\hat{\theta}_t$ of an unknown vector θ_t based on the observations $y_1^i, y_2^i, \dots, y_t^i$ and measurement points $\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_t^i$, $i \in \mathcal{N}$, through the minimization of time-varying *mean-risk functional* (3).

Remark. There exist two special cases of measurement model (2) related to the different types of noise v_t^i and disturbance ξ_t : (i) If drift θ_t is deterministic then $F_t^i(\theta) = f_{\xi_t}^i(\theta)$, $\xi_t = t$, and the measurement model may be defined in more conventional way as $y_t^i = F_t^i(\mathbf{x}_t^i) + v_t^i$, (ii) If noise v_t^i has a probability distribution then we may consider it as a part of disturbance ξ_t . The measurement model for this case is $y_t^i = f_{\xi_t}^i(\mathbf{x}_t^i)$.

B. Communication Network

In centralized networks, it is required to transmit all needed information such as $y_1^i, y_2^i, \dots, y_t^i$, $\mathbf{x}_1^i, \mathbf{x}_2^i, \dots, \mathbf{x}_t^i$, $i \in \mathcal{N}$, to a fusion center in order to estimate the unknown vector θ_t . In such networks, robustness of the fusion center and quality of the communication channels become a critical factor. In many

situations, sensors may only have the capability to exchange information with their neighbors. The communication with neighbors may be much cheaper or much faster than transmission to fusion center as well. Moreover, the information may be transmitted over the noisy communication channels and with delays, and the network topology may vary over time. Also, in practice, the network cost constraints naturally arise. For example, we don't have communication channels with infinite bandwidth and the response time of the network should be practically reasonable. All these factors motivate the development of distributed decentralized algorithms.

To formalize the distributed setting, we assume that at time instant t sensors are able to communicate with their neighbors through the network defined by graph $\mathcal{G}_{B_t} = (\mathcal{N}_{B_t}, \mathcal{E}_{B_t})$. The corresponding connectivity matrix B_t satisfies some averaged cost constraints (1) with level Q .

We also assume that sensor i obtains its current estimate $\hat{\theta}_t^i$ based on its noisy observation (2) and, if the set $\mathcal{N}_t^i = \{j \in \mathcal{N}_{B_t} : (j, i) \in \mathcal{E}_{B_t}\}$ is not empty, also on the current estimates transmitted by its neighbors through the noisy channels

$$\tilde{\theta}_t^{i,j} = \hat{\theta}_t^j + \mathbf{w}_t^{i,j}, \quad j \in \mathcal{N}_t^i, \quad (4)$$

where $\mathbf{w}_t^{i,j}$ is communication noise. If $j \notin \mathcal{N}_t^i$ we set $\tilde{\theta}_t^{i,j} = 0$.

C. Example

In this subsection, we present an example illustrating the considered problem statement. Given a distributed network consisting of $n = 6$ planar sensors identified by $i \in \mathcal{N} = \{1, 2, \dots, 6\}$. The state of sensor i is $\mathbf{s}^i \in \mathbb{R}^2$. We assume that the states are known and doesn't depend on time, i.e. the sensors are stationary. In the sensing range of the sensors, there are $m = 10$ moving planar targets identified by $l \in \mathcal{M} = \{1, 2, \dots, 10\}$. The goal of each sensor i is to estimate the states of all targets $\mathbf{r}_t^l \in \mathbb{R}^2$ at time instant t .

Let $\theta_t = \text{col}(\mathbf{r}_t^1, \dots, \mathbf{r}_t^{10}) \in \mathbb{R}^{20}$ be the common state vector of all targets, $\hat{\theta}_t^i = \text{col}(\hat{\mathbf{r}}_t^{i,1}, \dots, \hat{\mathbf{r}}_t^{i,10})$ be a common vector of i -th sensor current estimates. Each target $l \in \mathcal{M}$ changes the position according to the following dynamics:

$$\mathbf{r}_t^l = \mathbf{r}_{t-1}^l + \zeta_{t-1}^l, \quad l \in \mathcal{M}, \quad (5)$$

where ζ_{t-1}^l are random vectors uniformly distributed in a ball. We assume that at time instant t sensor i is able to measure the squared distance $\rho_t^{i,l} = \rho(\mathbf{s}^i, \mathbf{r}_t^l) = \|\mathbf{r}_t^l - \mathbf{s}^i\|^2$ to some moving target \mathbf{r}_t^l .

The network is modelled by complete graph \mathcal{G}_C , for which we have the following *topology constraints*: each sensor $i \in \mathcal{N}$ at each time instant t is able to measure the noisy squared distance to only one target $l \in \mathcal{M}$ and to receive estimates $\hat{\theta}_t^j$ and measurements $\rho_t^{j,l}$ only from one randomly chosen neighbor $j \in \mathcal{N}_t^i$. This leads to the communication protocol satisfying averaged cost constraints with level $Q = 1$ considered as example at the end of Section II.

Let sensor i receive the current estimate and measurement from some neighbor with identifier $j \in \mathcal{N}$. Denote by $\mathbf{u} = \text{col}(i, j, l)$ the vector, where the first element is the identifier of a sensor, the intermediate element is the identifier of a neighbor, which shares its information with sensor i ,

and the last element is the identifier of a target, which this sensor observe at time instant t . Note that in general there may be several intermediate elements. Also, denote by $\bar{\rho}_t(\mathbf{u}) = \rho(\mathbf{s}^i, \mathbf{r}_t^{l(\mathbf{u})}) - \rho(\mathbf{s}^j, \mathbf{r}_t^{l(\mathbf{u})})$ the residual between measurements of sensor i and its neighbor j . Here and after, $l(\mathbf{u}) : \mathbb{R}^{|\mathbf{u}|} \rightarrow \mathbb{R}$ gives the last element of \mathbf{u} . In this case, using the square difference formula we derive

$$C^{\mathbf{u}} \mathbf{r}_t^l = D_t^{\mathbf{u}} \Rightarrow C^{\mathbf{u}^T} C^{\mathbf{u}} \mathbf{r}_t^l = C^{\mathbf{u}^T} D_t^{\mathbf{u}} \Rightarrow I^{\mathbf{u}} \mathbf{r}_t^l = H_t^{\mathbf{u}}, \quad (6)$$

where $I^{\mathbf{u}} = [C^{\mathbf{u}^T} C^{\mathbf{u}}]' C^{\mathbf{u}^T} C^{\mathbf{u}}$, $H_t^{\mathbf{u}} = [C^{\mathbf{u}^T} C^{\mathbf{u}}]' C^{\mathbf{u}^T} D_t^{\mathbf{u}}$, $C^{\mathbf{u}} = 2(\mathbf{s}^j - \mathbf{s}^i)^T$, $D_t^{\mathbf{u}} = \bar{\rho}_t(\mathbf{u}) + \|\mathbf{s}^j\|^2 - \|\mathbf{s}^i\|^2$, and $[\cdot]'$ denotes a vector or matrix Moore–Penrose inverse.

Denote by U^i the set of all vectors \mathbf{u} with the first element i . Let $\mathbf{u}_t^i \in U^i$ be a random variable and input $\mathbf{x}^i = \hat{\theta}_t^i$ be fixed. We consider observation model (2) as follows

$$y_t^i = f_{\xi_t}^i(\mathbf{x}^i) = \|I^{\mathbf{u}_t^i} \hat{\mathbf{r}}^{i,l(\mathbf{u}_t^i)} - H_t^{\mathbf{u}_t^i}\|^2, \quad (7)$$

where ξ_t consists of all \mathbf{u}_t^i generated at time instant t , i.e. $\xi_t = \text{col}(\theta_t, \mathbf{u}_t^1, \mathbf{u}_t^2, \dots, \mathbf{u}_t^{\theta})$.

This leads us to following individual mean-risk sub-functionals $F_t^i(\mathbf{x}^i) = \mathbb{E}_{\mathcal{F}_{t-1}} f_{\xi_t}^i(\mathbf{x}^i)$, which are equal to $\frac{1}{|U^i|} \sum_{\mathbf{u}^i \in U^i} \|I^{\mathbf{u}^i} \hat{\mathbf{r}}^{i,l(\mathbf{u}^i)} - H_t^{\mathbf{u}^i}\|^2$ when positions of all targets do not evolve over time.

IV. MAIN RESULT

In this section, we present the main result of this paper. All proofs are relegated to Appendix.

A. SPSA-based Consensus Algorithm

Let Δ_k^i , $k = 1, 2, \dots$, $i \in \mathcal{N}$, be an observed sequence of independent random vectors in \mathbb{R}^d , called the *simultaneous test perturbation*, with symmetrical distribution functions $P_k^i(\cdot)$, and let $\mathbf{K}_k^i(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $k = 1, 2, \dots$, be a set of vector functions (kernels).

Let us take fixed nonrandom initial vectors $\hat{\theta}_0^i \in \mathbb{R}^d$, positive step-size α , gain coefficient γ , and choose sequences of such nonnegative numbers $\{\beta_k^+\}$ and $\{\beta_k^-\}$ that $\beta_k = \beta_k^+ + \beta_k^- > 0$. We consider the algorithm with two observations of distributed sub-functions $f_{\xi_t}^i(\theta)$ for each agent $i \in \mathcal{N}$ for constructing sequences of measurement points $\{\mathbf{x}_t^i\}$ and estimates $\{\hat{\theta}_t^i\}$:

$$\begin{cases} \mathbf{x}_{2k}^i = \hat{\theta}_{2k-2}^i + \beta_k^+ \Delta_k^i, & \mathbf{x}_{2k-1}^i = \hat{\theta}_{2k-2}^i - \beta_k^- \Delta_k^i, \\ \hat{\theta}_{2k-1}^i = \hat{\theta}_{2k-2}^i, \\ \hat{\theta}_{2k}^i = \hat{\theta}_{2k-1}^i - \alpha \left(\frac{y_{2k}^i - y_{2k-1}^i}{\beta_k} \mathbf{K}_k^i(\Delta_k^i) + \right. \\ \left. \gamma \sum_{j \in \mathcal{N}_{2k-1}^i} b_{2k-1}^{i,j} (\hat{\theta}_{2k-1}^{i,j} - \hat{\theta}_{2k-1}^i) \right). \end{cases} \quad (8)$$

Algorithm (8) consists of two parts: (i) The first one is similar to SPSA-like algorithm from [20]. This part represents a pseudo-gradient of sub-functions $f_{\xi_t}^i(\theta)$; (ii) The second one coincides with Local Voting Protocol (LVP) from [23], where it is used to solve load balancing problem in stochastic networks. This part is determined for each sensor i by the weighted sum of differences between the information about the current estimate $\hat{\theta}_{2k-1}^i$ of sensor i and noisy information about the estimates of its neighbors.

Further, we denote by $\bar{\theta}_t = \text{col}(\hat{\theta}_t^1, \dots, \hat{\theta}_t^n)$ the common vector of estimates of all sensors at time instant t and by $\tilde{\theta}_t = \text{col}(\hat{\theta}_t^{1,1}, \tilde{\theta}_t^{2,1}, \dots, \hat{\theta}_t^{n,1}, \hat{\theta}_t^{1,2}, \dots, \hat{\theta}_t^{n,n})$ the corresponding vector of data transmitted over the noisy channels. Also we introduce the following notations: $\bar{\mathbf{y}}_t = \text{diag}_n(\text{col}(y_t^1, \dots, y_t^n))$, $\bar{\Delta}_k = \text{col}(\mathbf{K}_k^1(\Delta_k^1), \dots, \mathbf{K}_k^n(\Delta_k^n))$. Using the notations introduced above, the algorithm (8) can be rewritten in the following form

$$\bar{\theta}_{2k} = \bar{\theta}_{2k-1} - \alpha \left[\left(\frac{\bar{\mathbf{y}}_{2k} - \bar{\mathbf{y}}_{2k-1}}{\beta_k} \otimes \mathbf{I}_d \right) \bar{\Delta}_k + \gamma \left(\bar{\mathcal{L}}_{2k-1} \otimes \mathbf{I}_d \right) \tilde{\theta}_{2k-1} \right] \quad (9)$$

where $(n \times n^2)$ matrix $\bar{\mathcal{L}}_{2k-1}$ is defined in such a way that its i -th row consists of zeros except the elements from the position $(j-1)n+1$ to jn which coincide with i -th row of $\mathcal{L}(B_{2k-1})$.

B. Main Assumptions

For any $i \in \mathcal{N}$ let us formulate assumptions about functions $F_t^i(\mathbf{x})$, $f_{\xi}^i(\mathbf{x})$, disturbances, network topology, randomized perturbations Δ_k^i , and noises.

1: The functions $F_t^i(\cdot)$ are convex, they have a common minimum point θ_t and

$$\forall \mathbf{x} \in \mathbb{R}^d \quad \langle \mathbf{x} - \theta_t, \mathbb{E}_{\mathcal{F}_{t-1}} \nabla f_{\xi_t}^i(\mathbf{x}) \rangle \geq 0.$$

2: $\forall \xi \in \Xi$, $\forall i \in \mathcal{N}$ the gradient $\nabla f_{\xi}^i(\mathbf{x})$ satisfies the Lipschitz condition: $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$

$$\|\nabla f_{\xi}^i(\mathbf{x}_1) - \nabla f_{\xi}^i(\mathbf{x}_2)\| \leq M \|\mathbf{x}_1 - \mathbf{x}_2\|$$

with the same constant $M > 0$.

3: The gradient $\nabla f_{\xi_t}^i$ is uniformly bounded in the mean-square sense at the minimum point θ_t : $\forall t \mathbb{E} \|\nabla f_{\xi_t}^i(\theta_t)\|^2 \leq g_2^2$, $\mathbb{E} \langle \nabla f_{\xi_t}^i(\theta_t), \nabla f_{\xi_{t-1}}^i(\theta_{t-1}) \rangle \leq g_2^2$ ($g_2 = 0$ if ξ_t is not a random parameter, i.e. $f_{\xi_t}^i(\mathbf{x}) = F_t^i(\mathbf{x})$).

4: The drift is bounded: a) $\|\theta_t - \theta_{t-1}\| \leq \delta_\theta < \infty$, or $\mathbb{E} \|\theta_t - \theta_{t-1}\|^2 \leq \delta_\theta^2$ and $\mathbb{E} \|\theta_t - \theta_{t-1}\| \|\theta_{t-1} - \theta_{t-2}\| \leq \delta_\theta^2$ if a sequence $\{\xi_t\}$ is random;

b) $\mathbb{E}_{\mathcal{F}_{2k-2}} f_{\xi_{2k}}^i(\mathbf{x}) - f_{\xi_{2k-1}}^i(\mathbf{x})|^q \leq \delta_\theta^q (g_0^q + g_1^q \|\mathbf{x} - \theta_{2k-2}\|^q)$ for $q = 1, 2$ and for any $i \in \mathcal{N}$.

5: a) Graphs \mathcal{G}_{B_t} , $t = 0, \dots$ are i.i.d., i.e. the random events of appearance of ‘‘time-varying’’ edge (j, i) in graph \mathcal{G}_{B_t} are independent and identically distributed for the fixed pair (j, i) , $i \in \mathcal{N}, j \in \mathcal{N}_{\max}^i = \cup_t \mathcal{N}_t^i$.

b) For all $i \in \mathcal{N}$, $j \in \mathcal{N}_t^i$ weights $b_t^{i,j}$ are independent random variables with mean values (mathematical expectations): $\mathbb{E} b_t^{i,j} = b_{av}^{i,j}$, and bounded variances: $\mathbb{E} \|B_t - B_{av}\|^2 \leq \sigma_B^2$ where $B_{av} = [b_{av}^{i,j}]$.

c) $\mathbb{E} \sum_{j \in \mathcal{N}_t^i} (b_t^{i,j})^2 \leq \frac{Q^2}{n-1}$.

d) Graph $\mathcal{G}_{B_{av}}$ is strongly connected.

6: For $k = 1, 2, \dots$, the successive differences $\tilde{v}_k^i = v_{2k}^i - v_{2k-1}^i$ of observation noise are bounded: $|\tilde{v}_k^i| \leq c_v < \infty$, or $\mathbb{E} (\tilde{v}_k^i)^2 \leq c_v^2$ if a sequence $\{\tilde{v}_k^i\}$ is random.

7: For $t = 1, 2, \dots$, $\forall i \in \mathcal{N}, \forall j \in \mathcal{N}$ the communication noise $\mathbf{w}_t^{i,j}$ is random i.i.d. (independent identically distributed) with zero-mean $\mathbb{E} \mathbf{w}_t^{i,j} = 0$ and bounded disturbances: $\mathbb{E} \|\mathbf{w}_t^{i,j}\|^2 \leq$

σ_w^2 . All random vectors and values $\mathbf{w}_t^{i,j}$, $b_t^{i,j}$, ξ_t , and ξ_{t+1} are mutually independent (if they are random).

8: For any $i, j \in \mathcal{N}$, $k = 1, 2, \dots$,

a) Vectors Δ_k^i are mutually independent.

b) Δ_k^i and ξ_{2k-1}, ξ_{2k} (if they are random) do not depend on the σ -algebra \mathcal{F}_{2k-2} .

c) If $\xi_{2k-1}, \xi_{2k}, \tilde{v}_k^i, \mathbf{w}_{2k-1}^{i,j}, b_{2k-1}^{i,j}$ are random, then random vectors Δ_k^i and elements $\xi_{2k-1}, \xi_{2k}, \tilde{v}_k^i, \mathbf{w}_{2k-1}^{i,j}, b_{2k-1}^{i,j}$ are independent.

d) For $k = 1, 2, \dots$, vectors Δ_k^i and vector functions $\mathbf{K}_k^i(\cdot)$ along with simultaneous perturbation symmetrical distribution functions $P_k(\cdot)$ satisfy the conditions

$$\int \mathbf{x} P_k(d\mathbf{x}) = \int \mathbf{x} \|\mathbf{K}_k^i(\mathbf{x})\|^2 P_k(d\mathbf{x}) = \int \mathbf{K}_k^i(\mathbf{x}) P_k(d\mathbf{x}) = 0,$$

$$\int \langle \mathbf{e}, \mathbf{x} \rangle \mathbf{K}_k^i(\mathbf{x}) P_k(d\mathbf{x}) = \langle \mathbf{e}, \mathbf{1}_d \rangle \mathbf{1}_d, \int \|\mathbf{x}\|^2 P_k(d\mathbf{x}) \leq c_\Delta^2, \quad (10)$$

$$\int \|\mathbf{K}_k^i(\mathbf{x})\|^2 P_k(d\mathbf{x}) \leq c_\Delta^2, \quad \int \|\mathbf{K}_k^i(\mathbf{x})\|^2 \|\mathbf{x}\|^2 P_k(d\mathbf{x}) \leq c_\Delta^4.$$

Note that all Assumptions 1–8 are standard for the considered problem.

Remark. Usually, it is practically reasonable to define $\{\Delta_k^i\}$ as a sequence of independent Bernoulli random vectors from \mathbb{R}^d with each component independently taking values $\pm \frac{1}{\sqrt{2}}$ with probabilities $\frac{1}{2}$ and $\mathbf{K}_k^i(\mathbf{x}) \equiv \mathbf{x}$ as kernel functions. For this case, we have $c_\Delta = 1$. The case, when $\beta_k^+ = \beta_k^-$ and decreasing to zero sequence α_k is used instead of constant step-size α , corresponds to the SPSA algorithm in [16]. The similar algorithm with randomly varying truncations and randomized difference was studied in [26] where the case $\beta_k^- = 0$ was additionally considered.

Example. Return back to the example from Section III-C and check Assumptions 1–5.

1. Using (6) and (7), we obtain for gradient

$$\langle \mathbf{x} - \theta_t, \mathbb{E}_{\mathcal{F}_{t-1}} \nabla f_{\xi_t}^i(\mathbf{x}) \rangle = \mathbb{E}_{\mathcal{F}_{t-1}} (\mathbf{x}^{i,l(\mathbf{u}_t^i)} - \mathbf{r}^{l(\mathbf{u}_t^i)})^T [I^{\mathbf{u}_t^i}]^T$$

$$I^{\mathbf{u}_t^i} (\mathbf{x}^{i,l(\mathbf{u}_t^i)} - \mathbf{r}^{l(\mathbf{u}_t^i)}) \geq 0.$$

2. Using (7), we obtain $\|\nabla f_{\xi}^i(\mathbf{x}_1) - \nabla f_{\xi}^i(\mathbf{x}_2)\| = \|2[I^{\mathbf{u}_t^i}]^T I^{\mathbf{u}_t^i} (\mathbf{x}_1^{l(\mathbf{u}_t^i)} - \mathbf{x}_2^{l(\mathbf{u}_t^i)})\| \leq M \|\mathbf{x}_2 - \mathbf{x}_1\|$, where $M = \max_i \|2[I^{\mathbf{u}_t^i}]^T I^{\mathbf{u}_t^i}\|$.

3. $\nabla f_{\xi_t}^i(\theta_t) = 0$. Hence, $g_2 = 0$.

4. Assumption about the drift holds for $\delta_\theta = n\delta_\zeta$ and by virtue of drift model (5) when ζ_t^i are random i.i.d. vectors with $\mathbb{E}\zeta_t^i = 0$, and $\mathbb{E}\|\zeta_t^i\|^2 \leq \delta_\zeta^2$, $g_0 = 4\sqrt{2}\bar{s}^2$, $g_1 = 8\sqrt{2}\bar{s}^2$, where $\bar{s} = \max_{i,j} \|s^i - s^j\|$.

5. a), c), d) hold by the construction; in b) $b_{av}^{i,j} = 0.2$, $i \neq j$, $b_{av}^{i,i} = 0$, $\sigma_B^2 = 4.8$.

C. Analysis of the Estimation Accuracy

To analyze the quality of estimates we apply the following definition for the problem of minimum tracking for mean-risk functional (3):

Definition 2. A sequence of estimates $\{\bar{\theta}_{2k}^i\}$ has an asymptotically efficient upper bound $\bar{L} > 0$ of residuals of estimation if $\forall \varepsilon > 0 \exists \bar{k}$ such that $\forall k > \bar{k}$

$$\sqrt{E\|\bar{\theta}_{2k}^i - \mathbf{1}_n \otimes \theta_{2k}\|^2} \leq \bar{L} + \varepsilon.$$

Denote $\bar{\lambda}_2 = \text{Re}(\lambda_2(\mathcal{L}(B_{av})))$, $\bar{\lambda}_m = \frac{1}{\lambda_{\max}(\mathcal{L}(B_{av})^T \mathcal{L}(B_{av}))}$, $c_+ = \max_k \frac{\beta_k^+}{\beta_k}$, $\bar{\beta} = \max_k \frac{1}{\beta_k}$, $\bar{c} = \max_k \left(\frac{\beta_k^+}{\beta_k} \right)^2 + \left(\frac{\beta_k^-}{\beta_k} \right)^2$, $\bar{\beta} = \max_k \frac{(\beta_k^+)^2}{\beta_k} + \frac{(\beta_k^-)^2}{\beta_k}$, $c_m = \bar{\lambda}_m^2 + \sigma_B^2$, $c_1 = c_\Delta \bar{\lambda}_m M (\delta_\theta g_1 \bar{\beta} + c_\Delta)$, $c_2 = \frac{2c_\Delta^2 (\delta_\theta^2 g_1^2 \bar{\beta}^2 + c_\Delta^2 M^2)}{c_\Delta}$, $c_\mu = (\bar{\lambda}_2 - \alpha c_1) / c_m$, $c_d = \sqrt{1 - \alpha^2 c_2 c_m} / (\bar{\lambda}_2 - \alpha c_1)^2$.

The following theorem shows the asymptotically efficient upper bound of estimation residuals provided by algorithm (8).

Theorem 1: If Assumptions 1–8 hold, $\bar{\beta} = \min_k \beta_k > 0$, positive constant α is sufficiently small:

$$\alpha < \frac{\bar{\lambda}_2}{c_1 + \sqrt{c_2 c_m}} \quad (11)$$

and

$$c_\mu(1 - c_d) < \alpha\gamma < c_\mu(1 + c_d) \quad (12)$$

then the averaged cost constraint (1) holds and the sequence of estimates provided by algorithm (8) has an asymptotically efficient upper bound which equals to

$$\bar{L} = \frac{1}{\mu} \left(h + \sqrt{h^2 + l\mu} \right), \quad (13)$$

where $\mu = 2\gamma\bar{\lambda}_2 - \alpha(c_m\gamma^2 + \alpha(2\gamma c_1 + c_2))$, $h = \gamma c_3 + c_4$, $l = \alpha\gamma^2 Q^2 \sigma_w^2 + c_5$, $c_3 = 2\sqrt{n}\bar{\lambda}_m \delta_\theta + \alpha\bar{\lambda}_m c_\Delta M (\delta_\theta g_0 \bar{\beta} + \bar{\beta} c_\Delta^2)$, $c_4 = M(c_+ + c_\Delta^4 g_1 \bar{c} + 2c_\Delta^2(1 + c_+))\delta_\theta + c_\Delta^5 M^2 \bar{\beta}$, $c_5 = \frac{4n\delta_\theta^2}{\alpha} + 2\alpha c_\Delta^2 (\bar{\beta}^2 n(c_v^2 + \delta_\theta^2 g_0^2) + c_\Delta^2 \bar{c} n M(c_v + \delta_\theta g_0) + c_\Delta^3 n M \bar{\beta} (M\delta_\theta + \delta_\theta c_+ + g_2)) + 2Mn(\delta_\theta^2 c_+ + c_\Delta^3 \bar{\beta}) + c_\Delta^2 n (\delta_\theta^2 (1 + c_+)^2 + g_2^2 + M^2 \bar{\beta}^2 c_\Delta^4)$.

See the proof of Theorem 1 in Appendix.

Remarks. 1. The bound L in the Theorem 1 is tight, so there exists no $L' < L$ such that the statement of the Theorem 1 still holds if all inequalities from the Assumptions 1–8 are replaced by equation.

2. The observation noise v_t^i in Theorem 1 can be said to be almost arbitrary since it may either be nonrandom but bounded or it may also be a realization of some stochastic process with arbitrary internal dependencies. In particular, to prove the results of Theorem 1, there is no need to assume that v_t^i and \mathcal{F}_{t-1} are not dependent.

3. The proof of Theorem 1 allows for consideration of random sequences $\{\beta_k^+\}$ and $\{\beta_k^-\}$ whose values at iteration k are measurable under the corresponding σ -algebra \mathcal{F}_{2k-2} . This fact is sometimes useful from a practical point of view.

4. The result of the Theorem 1 shows that for the case without drift ($\delta_\theta = 0$) and $g_2 = 0$ under any noise level c_v the asymptotic upper bound can be made infinitely small by choosing sufficiently small α and β_k^\pm . At the same time, in the case of drift, the bigger drift norm δ_θ can be compensated by choosing a bigger step-size α and β_k^\pm . This leads to a tradeoff between making α smaller because of noisy observations and making α bigger due to the drift of optimal points.

V. SIMULATION

In this section, we present the numerical experiments, which illustrate the performance of the suggested algorithm. We apply the algorithm to the problem described in Section III-C.

The starting positions of the targets are chosen randomly from the interval $[0; 100]$. The states of the targets evolve over time as follows: $\mathbf{r}_t^l = \mathbf{r}_{t-1}^l + \chi_{t-1}^l$. Let χ_{t-1}^l be a random vector uniformly distributed on the ball of radius equal to 0.2. The sensors don't move and their coordinates are random values uniformly distributed in interval $[100; 120]$. We consider observation model (2) defined as $y_t^i = \|I^{u_t^i} \hat{\mathbf{r}}_t^{i,l}(\mathbf{u}_t^i) - H_t^{u_t^i}\|^2 + v_t^i$, where v_t^i is modelled as unknown-but-bounded disturbance falling within interval $[0.6; -0.6]$.

Algorithm (8) working on each node has the following parameters: $\alpha = 0.03$, $\beta = 1.5$, $\gamma = 1.5$. The initial estimate on each sensor for each target coordinate was chosen randomly from the interval $[50; 100]$. Fig. 1 shows how the residuals evolve over time. Figures show that there exists time instant t starting with which the estimations converge to the actual value and move next to it.

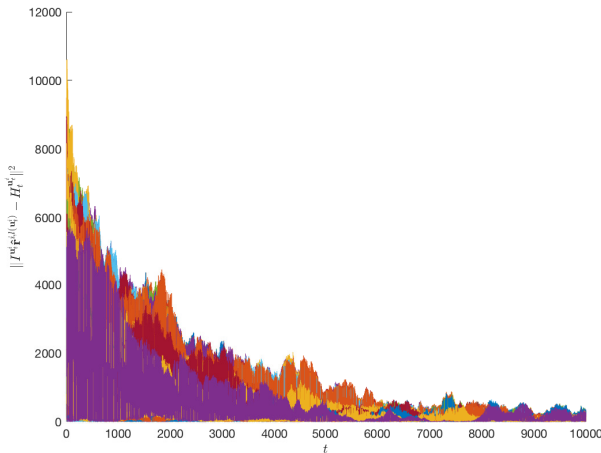


Fig. 1. Residuals $\|I^{u_t^i} \hat{\mathbf{r}}_t^{i,l}(\mathbf{u}_t^i) - H_t^{u_t^i}\|^2$ obtained by nodes.

VI. CONCLUSION

In this paper, we proposed a new SPSA-based consensus algorithm for distributed tracking under unknown-but-bounded disturbances. Compared to the SPSA algorithm, this algorithm is suitable for distributed problems due to the relaxed assumption regarding the strong convexity of the minimized mean-risk functional. In many practical applications, the network processes the data under certain constraints, and the data transmission is accompanied by noise. In this paper, we consider such noisy data transmission and the case where a communication protocol has to satisfy prespecified cost constraints. Communication cost constraints are satisfied by exploiting a specific intentionally randomized topology of the network communication graph. We obtain an upper bound on the mean square error of estimates of tracking unknown time-varying parameters under unknown-but-bounded observation errors and noisy communication channels.

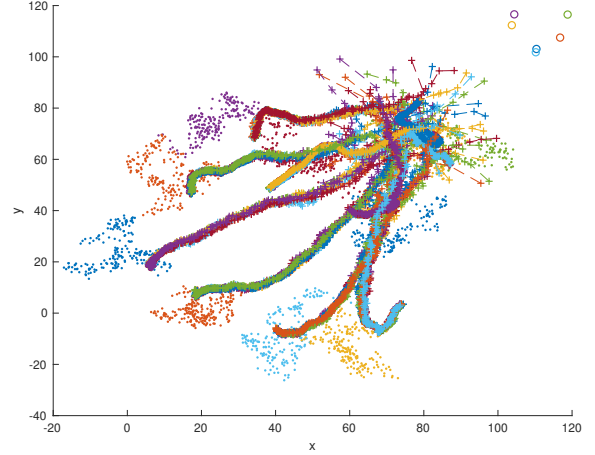


Fig. 2. The estimates $\hat{\mathbf{r}}_t^{i,l}$ obtained by nodes and actual targets positions $\mathbf{r}_t^{i,l}$. (Empty circles denote sensor positions, targets movement is depicted as a series of shaded circles and plus signs show the estimated target positions.) The figure shows sparse data for clarity: each 50th position of targets and the estimates.

REFERENCES

- [1] R. Olfati-Saber and R. M. Murray, "Consensus problems in networks of agents with switching topology and time-delays," *IEEE Transactions on Automatic Control*, vol. 49, no. 9, pp. 1520–1533, 2004.
- [2] W. Ren and R. W. Beard, *Distributed consensus in multi-vehicle cooperative control*. Springer, 2008.
- [3] F. L. Lewis, H. Zhang, K. Hengster-Movric, and A. Das, *Cooperative control of multi-agent systems: optimal and adaptive design approaches*. Springer Science & Business Media, 2013.
- [4] M. DeGroot, "Reaching a consensus," *J. Am. Stat. Assoc.*, vol. 69, pp. 118–121, 1974.
- [5] V. Borkar and P. Varaiya, "Asymptotic agreement in distributed estimation," *IEEE Trans. Autom. Control*, vol. 27, pp. 650–655, 1982.
- [6] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," in *1984 American Control Conference*, 1984, pp. 484–489.
- [7] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [8] A. Falsone, I. Notarnicola, G. Notarstefano, and M. Prandini, "Tracking-admm for distributed constraint-coupled optimization," *Automatica*, vol. 117, p. 108962, 2020.
- [9] M. Rabbat and R. Nowak, "Distributed optimization in sensor networks," in *Proceedings of the 3rd international symposium on Information processing in sensor networks*. ACM, 2004, pp. 20–27.
- [10] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, p. 48, 2009.
- [11] M. Zhu and S. Martínez, "Discrete-time dynamic average consensus," *Automatica*, vol. 46, no. 2, pp. 322–329, 2010.
- [12] P. Di Lorenzo and G. Scutari, "Next: In-network nonconvex optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.
- [13] S. Aeron, V. Saligrama, and D. A. Castanon, "Efficient sensor management policies for distributed target tracking in multihop sensor networks," *IEEE Transactions on Signal Processing*, vol. 56, no. 6, pp. 2562–2574, 2008.
- [14] A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. J. Towfic, "Diffusion strategies for adaptation and learning over networks: An examination of distributed strategies and network behavior," *IEEE Signal Process. Mag.*, vol. 30, pp. 155–171, 2013.
- [15] S. Xie and L. Guo, "Analysis of normalized least mean squares-based consensus adaptive filters under a general information condition," *SIAM J. Control Optim.*, vol. 56, pp. 3404–3431, 2018.

- [16] J. C. Spall, "Multivariate stochastic approximation using a simultaneous perturbation gradient approximation," *IEEE Transactions on Automatic Control*, vol. 37, no. 3, pp. 332–341, 1992.
- [17] B. Polyak, *Introduction to optimization*. Optimization Software, Publications Division (New York), 1987.
- [18] H. J. Kushner and G. Yin, *Stochastic Approximation and Recursive Algorithms and Applications*. New York, Springer-Verlag, 2003.
- [19] V. S. Borkar, *Stochastic approximation: a dynamical systems viewpoint*. Springer, 2009.
- [20] O. Granichin and N. Amelina, "Simultaneous perturbation stochastic approximation for tracking under unknown but bounded disturbances," *IEEE Transactions on Automatic Control*, vol. 60, no. 6, pp. 1653–1658, 2015.
- [21] J. Zhu and J. C. Spall, "Tracking capability of stochastic gradient algorithm with constant gain," in *Decision and Control (CDC), 2016 IEEE 55th Conference on*. IEEE, 2016, pp. 4522–4527.
- [22] —, "Probabilistic bounds in tracking a discrete-time varying process," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 4849–4854.
- [23] N. Amelina, A. Fradkov, Y. Jiang, and D. J. Vergados, "Approximate consensus in stochastic networks with application to load balancing," *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1739–1752, 2015.
- [24] V. Erofeeva, O. Granichin, N. Amelina, Y. Ivanskiy, and Y. Jiang, "Distributed tracking via simultaneous perturbation stochastic approximation-based consensus algorithm," in *2019 IEEE 58th Conference on Decision and Control (CDC)*. IEEE, 2019, pp. 6050–6055.
- [25] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE/ACM Transactions on Networking (TON)*, vol. 14, no. SI, pp. 2508–2530, 2006.
- [26] H.-F. Chen, T. E. Duncan, and B. Pasik-Duncan, "A kiefer-wolfowitz algorithm with randomized differences," *IEEE Transactions on Automatic Control*, vol. 44, no. 3, pp. 442–453, 1999.

APPENDIX

Proof of Theorem 1. At first, we prove that the averaged graph $\mathcal{G}_{B_{av}}$ corresponds the average cost constraint (1).

By virtue of Assumption 5c and Cauchy-Bunyakovsky-Schwarz inequality, we get

$$\mathbb{E}Cost(B_t) = \mathbb{E} \sum_{i \in \mathcal{N}} b_t^{i,j} \leq \sqrt{n \mathbb{E} \sum_{i \in \mathcal{N}} (b_t^{i,j})^2} \leq \sqrt{n \frac{Q^2}{n}} = Q.$$

Hence, the average cost constraint (1) holds.

At second, we study the asymptotic mean square properties of residuals $\eta_k = \|\bar{\theta}_{2k} - \mathbf{1}_n \otimes \theta_{2k}\|$.

Denote $\bar{\mathbf{s}}_k = \frac{\alpha}{\beta_k} ((\bar{\mathbf{y}}_{2k} - \bar{\mathbf{y}}_{2k-1}) \otimes \mathbf{I}_d) \bar{\Delta}_k$, $\mathbf{d}_t^i = \hat{\theta}_{2\lceil \frac{t-1}{2} \rceil}^i - \theta_t$, $\bar{\mathbf{d}}_t = \text{col}\{\mathbf{d}_t^1, \dots, \mathbf{d}_t^n\}$, where $\lceil \cdot \rceil$ is a ceiling function, $\bar{\mathbf{w}}_t = \text{col}\{\mathbf{w}_t^{1,1}, \mathbf{w}_t^{2,1}, \dots, \mathbf{w}_t^{n,1}, \mathbf{w}_t^{1,2}, \dots, \mathbf{w}_t^{n,n}\}$, $\bar{\mathbf{v}}_t = \text{col}\{\tilde{v}_t^1, \dots, \tilde{v}_t^n\}$,

Let $\tilde{\mathcal{F}}_{k-1} = \sigma\{\mathcal{F}_{k-1}, \bar{\mathbf{v}}_{2k-1}, \bar{\mathbf{v}}_{2k}, \xi_{2k-1}, \xi_{2k}, \bar{\Delta}_k, B_{2k-1}\}$ be the σ -algebra of probabilistic events generated by $\mathcal{F}_{k-1}, \bar{\mathbf{v}}_{2k-1}, \bar{\mathbf{v}}_{2k}, \xi_{2k-1}, \xi_{2k}, \bar{\Delta}_k, B_{2k-1}$, and $\tilde{\mathcal{F}}_{k-1} = \sigma\{\mathcal{F}_{k-1}, \bar{\mathbf{v}}_{2k-1}, \bar{\mathbf{v}}_{2k}, \xi_{2k-1}, \xi_{2k}, \bar{\Delta}_k\}$, $\tilde{\mathcal{F}}_{k-1} = \sigma\{\mathcal{F}_{k-1}, \bar{\mathbf{v}}_{2k-1}, \bar{\mathbf{v}}_{2k}, \xi_{2k-1}, \xi_{2k}\}$,

$$\mathcal{F}_{k-1} \subset \tilde{\mathcal{F}}_{k-1} \subset \bar{\mathcal{F}}_{k-1} \subset \tilde{\tilde{\mathcal{F}}}_{k-1} \subset \mathcal{F}_k.$$

By virtue of communication model (4), we obtain $\tilde{\theta}_t = \mathbf{1}_n \otimes \bar{\theta}_t + \bar{\mathbf{w}}_t$ and, according to the algorithm (9), we have $\eta_k =$

$$\begin{aligned} & \|\bar{\theta}_{2k-2} - \mathbf{1}_n \otimes \theta_{2k} - \bar{\mathbf{s}}_k - \alpha\gamma \bar{\mathcal{L}}_{2k-1} (\mathbf{1}_n \otimes \bar{\theta}_{2k-2} + \bar{\mathbf{w}}_t)\| = \\ & = \|\bar{\mathbf{g}}_k - \bar{\mathbf{s}}_k + \alpha\gamma \bar{\mathcal{L}}_{2k-1} \bar{\mathbf{w}}_{2k-1}\| \end{aligned}$$

where $\bar{\mathbf{g}}_k = (\mathbf{I}_{nd} - \alpha\gamma \mathcal{L}(B_{2k-1}) \otimes \mathbf{I}_d) \bar{\mathbf{d}}_{2k-2} + \mathbf{1}_n \otimes (\theta_{2k-2} - \theta_{2k})$ since it is not so hard to prove that $(\mathcal{L}(B_{2k-1}) \otimes \mathbf{I}_d) \mathbf{1}_n \otimes \theta_{2k-2} = 0$ based on the properties of Laplacian matrix $\mathcal{L}(B_{2k-1})$. Taking the conditional expectation over σ -algebra $\tilde{\mathcal{F}}_{k-1}$, by virtue of Assumption 6, we derive

$$\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \eta_k^2 = \|\bar{\mathbf{g}}_k - \bar{\mathbf{s}}_k\|^2 + \alpha^2 \gamma^2 \|B_{2k-1}\|^2 \sigma_w^2$$

since $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \bar{\mathbf{w}}_{2k-1} = 0$.

Assumption 5c gives the bound: $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \|B_{2k-1}\|^2 \leq Q^2$. Taking the conditional expectation over σ -algebra $\tilde{\mathcal{F}}_{k-1}$, by virtue of Assumption 5b, we get

$$\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \eta_k^2 = \|\bar{\mathbf{g}}_k - \bar{\mathbf{s}}_k\|^2 + \alpha^2 \gamma^2 (Q^2 \sigma_w^2 + \sigma_B^2 \eta_{k-1}^2)$$

since $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} (\mathcal{L}(B_{2k-1}) - \mathcal{L}(B_{av})) \bar{\mathbf{d}}_{2k-2} = 0$.

So, we obtain $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \eta_k^2 =$

$$\|\bar{\mathbf{g}}_k\|^2 + \|\bar{\mathbf{s}}_k\|^2 - 2\langle \bar{\mathbf{g}}_k, \bar{\mathbf{s}}_k \rangle + \alpha^2 \gamma^2 (Q^2 \sigma_w^2 + \sigma_B^2 \eta_{k-1}^2). \quad (14)$$

By virtue of Assumption 8c,d we have $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{v}_k \mathbf{K}_k^i(\Delta_k^i) = \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{v}_k \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \mathbf{K}_k^i(\Delta_k^i) = \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{v}_k \cdot 0 = 0$. Hence, taking the conditional expectation over σ -algebra $\tilde{\mathcal{F}}_{k-1}$ of both sides of the (14) and using observation model (2), we can assert the bound for $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \eta_k^2$ as follows:

$$\begin{aligned} \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \eta_k^2 & \leq \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \|\bar{\mathbf{g}}_k\|^2 - 2 \frac{\alpha}{\beta_k} \sum_{i \in \mathcal{N}} \langle \mathbf{d}_{2k}^i, \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{f}_k^i \mathbf{K}_k^i(\Delta_k^i) \rangle + \\ & + 2 \frac{\alpha}{\beta_k} \sum_{i \in \mathcal{N}} \langle \alpha\gamma (\mathcal{L}(B_{av}) \mathbf{d}_{2k-2}^i, \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{f}_k^i \mathbf{K}_k^i(\Delta_k^i)) \rangle + \\ & + \frac{\alpha^2}{\beta_k^2} \sum_{i \in \mathcal{N}} \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \left(\tilde{v}_k^i + \tilde{f}_k^i \right)^2 \|\mathbf{K}_k^i(\Delta_k^i)\|^2 + \\ & \alpha^2 \gamma^2 (Q^2 \sigma_w^2 + \sigma_B^2 \eta_{k-1}^2) \end{aligned} \quad (15)$$

where $\tilde{f}_k^i = f_{\xi_{2k}}^i(\mathbf{x}_{2k}) - f_{\xi_{2k-1}}^i(\mathbf{x}_{2k-1})$.

Under fulfilment of Assumption 5d, we have $\bar{\lambda}_2 > 0$ (see [1]). Hence, for the first term in (14) we derive

$$\begin{aligned} \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \|\bar{\mathbf{g}}_k\|^2 & \leq \mathbf{d}_{2k-2}^T (\mathbf{I}_{nd} - \alpha\gamma (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d))^T \times \\ & (\mathbf{I}_{nd} - \alpha\gamma (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d)) \mathbf{d}_{2k-2} + \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} 2\alpha\gamma \times \\ & \mathbf{d}_{2k-2}^T (\mathbf{I}_{nd} - \alpha\gamma (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d))^T \mathbf{1}_n \otimes (\theta_{2k-2} - \theta_{2k}) + \\ & \|\mathbf{1}_n \otimes (\theta_{2k-2} - \theta_{2k})\|^2 \leq \eta_{k-1}^2 - \mathbf{d}_{2k-2}^T \alpha\gamma \times \\ & (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d)^T \mathbf{d}_{2k-2} - \mathbf{d}_{2k-2}^T \alpha\gamma (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d) \mathbf{d}_{2k-2} + \\ & \alpha^2 \gamma^2 \mathbf{d}_{2k-2}^T (\mathbf{I}_{nd} - \mathcal{L}(B_{av}) \otimes \mathbf{I}_d)^T (\mathcal{L}(B_{av}) \otimes \mathbf{I}_d) \mathbf{d}_{2k-2} + \\ & \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} 2\alpha\gamma \eta_{k-1} \sqrt{n} \|(\mathbf{I}_d - \alpha\gamma (\mathcal{L}(B_{av}))\| \|\theta_{2k-2} - \theta_{2k}\| + 4n\delta_\theta^2 \\ & \leq (1 - 2\alpha\gamma \bar{\lambda}_2 + \alpha^2 \gamma^2 \bar{\lambda}_m^2) \eta_{k-1}^2 + 4\alpha\gamma \sqrt{n} \bar{\lambda}_m \delta_\theta \eta_{k-1} + 4n\delta_\theta^2. \end{aligned} \quad (16)$$

For any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$, by virtue of Taylor representation of $f_{\xi_t}^i(\mathbf{x})$ for $t^\pm = 2k - \frac{1}{2} \pm \frac{1}{2}$, we have

$$f_{\xi_{t^\pm}}^i(\mathbf{x}) = f_{\xi_{t^\pm}}^i(\mathbf{z}) + \langle \nabla f_{\xi_{t^\pm}}^i(\mathbf{z} + \rho_{\xi_{t^\pm}}^\pm (\mathbf{x} - \mathbf{z})), \mathbf{x} - \mathbf{z} \rangle, \quad (17)$$

where $\rho_{\xi_{t^\pm}}^\pm \in (0, 1)$.

For difference \tilde{f}_k^i , adding and subtracting $\langle \nabla f_{\xi_{t^\pm}}^i(\mathbf{z}), \mathbf{x}_{t^\pm}^i - \mathbf{z} \rangle$, we derive:

$$\tilde{f}_k^i = \sum_{t^\pm} \pm f_{t^\pm}^i(\mathbf{z}) \pm \langle \nabla f_{\xi_{t^\pm}}^i(\mathbf{z}), \mathbf{x}_{t^\pm}^i - \mathbf{z} \rangle \pm \bar{M}_{t^\pm}^i(\mathbf{z}) \quad (18)$$

where $\bar{M}_{t^\pm}^i(\mathbf{z}) = \langle \nabla f_{\xi_{t^\pm}}^i(\mathbf{z} + \rho_{\xi_{t^\pm}}^\pm(\mathbf{x}_{t^\pm}^i - \mathbf{z})) - \nabla f_{\xi_{t^\pm}}^i(\mathbf{z}, \mathbf{x}_{t^\pm}^i - \mathbf{z}) \rangle$. Hence, for $\mathbf{z} = \hat{\theta}_{2k-2}^i$, by virtue of Assumption 8c, we have $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{f}_k^i \mathbf{K}_k^i(\Delta_k^i) =$

$$\sum_{t^\pm} \pm \nabla f_{\xi_{t^\pm}}^i(\hat{\theta}_{2k-2}^i) \beta_k^\pm \pm \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \bar{M}_{t^\pm}^i(\hat{\theta}_{2k-2}^i) \mathbf{K}_k^i(\Delta_k^i),$$

since $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} f_{t^\pm}^i(\mathbf{z}) \mathbf{K}_k^i(\Delta_k^i) = 0$.

According to the Assumption 2, we have

$$\|\bar{M}_{t^\pm}^i(\hat{\theta}_{2k-2}^i)\| \leq M \|\rho_{\xi_{t^\pm}}^\pm(\mathbf{x}_{t^\pm}^i - \hat{\theta}_{2k-2}^i)\| \beta_k^\pm \|\Delta_k^i\| \leq M(\beta_k^\pm)^2 \|\Delta_k^i\|^2. \quad (19)$$

We can evaluate the second term in (15), using formula (19) and applying Assumptions 2,

$$\begin{aligned} \dots &\leq -2 \frac{\alpha}{\beta_k} \sum_{i \in \mathcal{N}} \sum_{t^\pm} \langle \hat{\theta}_{2k-2}^i - \theta_{t^\pm}, \nabla f_{\xi_{t^\pm}}^i(\hat{\theta}_{2k-2}^i) \beta_k^\pm \rangle - 2 \frac{\alpha}{\beta_k} \times \\ &\sum_{i \in \mathcal{N}} \langle \theta_{2k} - \theta_{2k-1}, \nabla f_{\xi_{2k-1}}^i(\hat{\theta}_{2k-2}^i) \beta_k^+ \rangle + 2 \frac{\alpha}{\beta_k} M c_\Delta^3 \sum_{t^\pm} (\beta_k^\pm)^2. \end{aligned}$$

Here the conditional expectation over σ -algebra \mathcal{F}_{k-1} for first terms with minus is not above zero by Assumption 1. By virtue of definition we have $\mathbb{E}_{\mathcal{F}_{k-1}} \nabla f_{\xi_{2k-1}}^i(\theta_{2k-1}) = 0$. Hence, applying the first part of Assumption 4, we get

$$\begin{aligned} \dots &\leq 2 \frac{\alpha}{\beta_k} M \mathbb{E}_{\mathcal{F}_{k-1}} \sum_{i \in \mathcal{N}} \delta_\theta \beta_k^+ \|\mathbf{d}_{2k-1}^i\| + c_\Delta^3 ((\beta_k^+)^2 + (\beta_k^-)^2) \\ &\leq 2\alpha M (\delta_\theta c_+ (\eta_{k-1} + n\delta_\theta) + n c_\Delta^3 \bar{\beta}). \end{aligned}$$

To evaluate the conditional expectation over σ -algebra $\tilde{\mathcal{F}}_{k-1}$ of the third term in (15) we use the following representation for the difference \tilde{f}_k^i

$$\begin{aligned} \tilde{f}_k^i &= f_{\xi_{2k}}^i(\mathbf{x}_{2k}) - f_{\xi_{2k-1}}^i(\mathbf{x}_{2k}) + f_{\xi_{2k-1}}^i(\mathbf{x}_{2k}) - f_{\xi_{2k-1}}^i(\mathbf{x}_{2k-1}) \\ &= \sum_{t^\pm} \pm f_{\xi_{t^\pm}}^i(\mathbf{x}_{2k}) + \langle \nabla f_{\xi_{2k-1}}^i(\hat{\theta}_{2k-2}^i \pm \rho_{\xi_{t^\pm}}^\pm \beta_k^\pm \Delta_k^i), \beta_k^\pm \Delta_k^i \rangle \end{aligned}$$

which is based on Taylor formula (17). By adding and subtraction $\sum_{t^\pm} \langle \nabla f_{\xi_{2k-1}}^i(\theta_{2k-1}), \beta_k^\pm \Delta_k^i \rangle$, using the first part of Assumption 9, we derive $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \tilde{f}_k^i \mathbf{K}_k^i(\Delta_k^i) =$

$$\begin{aligned} \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} \sum_{t^\pm} (\pm f_{\xi_{t^\pm}}^i(\mathbf{x}_{2k}) + \langle \nabla f_{\xi_{2k-1}}^i(\hat{\theta}_{2k-2}^i \pm \rho_{\xi_{t^\pm}}^\pm \beta_k^\pm \Delta_k^i), \beta_k^\pm \Delta_k^i \rangle) \\ \times \mathbf{K}_k^i(\Delta_k^i) + \langle \nabla f_{\xi_{2k-1}}^i(\theta_{2k-1}), \mathbf{1}_d \rangle \mathbf{1}_d. \end{aligned}$$

Taking the conditional expectation over σ -algebra \mathcal{F}_{k-1} , by virtue of properties $\mathbb{E}_{\mathcal{F}_{k-1}} \nabla f_{\xi_{2k-1}}^i(\theta_{2k-1}) = 0$ and the Assumptions 2,4,9, we get

$$\begin{aligned} \mathbb{E}_{\mathcal{F}_{k-1}} \|\tilde{f}_k^i \mathbf{K}_k^i(\Delta_k^i)\| &\leq (\delta_\theta (g_0 + g_1 \|\mathbf{d}_{2k-2}^i\|) + \\ &\sum_{t^\pm} M (\mathbb{E}_{\mathcal{F}_{k-1}} \|\mathbf{d}_{2k-1}^i\| + \beta_k^\pm c_\Delta) \beta_k^\pm c_\Delta) c_\Delta \quad (20) \end{aligned}$$

Hence, for the third term in (15) we have $\dots \leq 2 \frac{\alpha^2 \gamma}{\beta_k} \times$

$$\begin{aligned} \bar{\lambda}_m M c_\Delta \eta_{k-1} (\delta_\theta (g_0 + g_1 \eta_{k-1}) + \beta_k c_\Delta (\eta_{k-1} + \bar{\beta} c_\Delta)) \leq \\ 2 \frac{\alpha^2 \gamma}{\beta_k} \bar{\lambda}_m M c_\Delta ((\delta_\theta g_1 + \beta_k c_\Delta) \eta_{k-1}^2 + (\delta_\theta g_0 + \beta_k \bar{\beta} c_\Delta^2) \eta_{k-1}). \end{aligned}$$

Summing up the conditional expectations over σ -algebra \mathcal{F}_{k-1} of the second and third terms in (15) we derive

$$\begin{aligned} \dots &\leq 2\alpha^2 \gamma \bar{\lambda}_m M c_\Delta \left(\frac{\delta_\theta g_1}{\beta_k} + c_\Delta \right) \eta_{k-1}^2 + 2\alpha M (\delta_\theta c_+ + \\ &\alpha \gamma \bar{\lambda}_m \left(\frac{\delta_\theta g_0}{\beta_k} c_\Delta + \bar{\beta} c_\Delta^3 \right)) \eta_{k-1} + 2\alpha M n (\delta_\theta^2 c_+ + c_\Delta^3 \bar{\beta}). \quad (21) \end{aligned}$$

Consider the squared difference $(\tilde{v}_k^i + \tilde{f}_k^i)^2$. Using formula (18) with $\mathbf{z} = \hat{\theta}_{2k-2}$, the sum $(\tilde{v}_k^i + \tilde{f}_k^i)$ can be represented as a sum of five terms $\tilde{v}_k^i + \tilde{f}_k^i = a_1 + a_2 + a_3 + a_4$, where $a_1 = \tilde{v}_k^i$, $a_2 = \sum_{t^\pm} \pm f_{t^\pm}^i(\hat{\theta}_{2k-2})$, $a_3 = \sum_{t^\pm} \langle \nabla f_{\xi_{t^\pm}}^i(\hat{\theta}_{2k-2}), \Delta_k^i \beta_k^\pm \rangle$, and $a_4 = \sum_{t^\pm} \pm M_{t^\pm}^i(\theta_{2k-2})$.

The first two terms do not depend on Δ_k^i and $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} a_q \Delta_k^i \|\mathbf{K}_k^i(\Delta_k^i)\|^2 = 0$, $q = 1, 2$, by virtue of Assumption 8. Hence, we derive $\mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} (\tilde{v}_k^i + \tilde{f}_k^i)^2 \|\mathbf{K}_k^i(\Delta_k^i)\|^2 \leq$

$$\begin{aligned} c_\Delta^2 \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} (a_1 + a_2)^2 + 2(a_1 + a_2 + a_3)a_4 + a_3^2 + a_4^2 \leq \\ c_\Delta^2 \mathbb{E}_{\tilde{\mathcal{F}}_{k-1}} 2(a_1^2 + a_2^2 + (|a_1| + |a_2| + |a_3|)|a_4|) + a_3^2 + a_4^2. \end{aligned}$$

We need to estimate $\mathbb{E}_{\mathcal{F}_{k-1}} a_q^2$, $q = 1, \dots, 4$ and we can use formula $\mathbb{E}_{\mathcal{F}_{k-1}} |a_q| \leq \sqrt{\mathbb{E}_{\mathcal{F}_{k-1}} a_q^2}$, $q = 1, \dots, 4$ for the rest terms. Taking the conditional expectation over σ -algebra \mathcal{F}_{k-1} , by virtue of Assumptions 2–4 and (19), we evaluate

$$\begin{aligned} \mathbb{E}_{\mathcal{F}_{k-1}} a_1^2 &\leq c_v^2, \quad \mathbb{E}_{\mathcal{F}_{k-1}} |a_2|^q \leq \delta_\theta^q (g_0^q + g_1^q \|\mathbf{d}_{2k-2}^i\|^q), \quad q = 1, 2, \\ \mathbb{E}_{\mathcal{F}_{k-1}} a_3^q &\leq q \mathbb{E}_{\mathcal{F}_{k-1}} \left(\sum_{t^\pm} \langle \nabla f_{\xi_{t^\pm}}^i(\hat{\theta}_{2k-2}) - \nabla f_{\xi_{t^\pm}}^i(\theta_{t^\pm}), \Delta_k^i \beta_k^\pm \rangle \right)^q \\ &+ q \left(\sum_{t^\pm} \langle \nabla f_{\xi_{t^\pm}}^i(\theta_{t^\pm}), \Delta_k^i \beta_k^\pm \rangle \right)^q \leq 2c_\Delta^q ((M\beta_k (\|\mathbf{d}_{2k-2}^i\| + \delta_\theta) \\ &+ \delta_\theta \beta_k^+))^q + \beta_k^q g_2^q, \quad q = 1, 2, \\ \mathbb{E}_{\mathcal{F}_{k-1}} a_4^2 &\leq M^2 ((\beta_k^+)^2 + (\beta_k^-)^2) c_\Delta^4. \end{aligned}$$

Taking the conditional expectation over σ -algebra \mathcal{F}_{k-1} for the fourth term in (15) we get using Assumptions 2–5

$$\begin{aligned} 2 \frac{\alpha^2}{\beta_k^2} \mathbb{E}_{\mathcal{F}_{k-1}} \sum_{i \in \mathcal{N}} (\tilde{v}_k^i + \tilde{f}_k^i)^2 \|\mathbf{K}_k^i(\Delta_k^i)\|^2 &\leq c_\Delta^2 \frac{\alpha^2}{\beta_k^2} (2(nc_v^2 + \\ &n\delta_\theta^2 g_0^2 + \delta_\theta^2 g_1^2 \eta_{k-1}^2 + (nc_v + n\delta_\theta g_0 + \delta_\theta g_1 \eta_{k-1} + \\ &c_\Delta M \beta_k \eta_{k-1} + nc_\Delta (M\delta_\theta \beta_k + \delta_\theta \beta_k^+ + \beta_k g_2)) \times \\ &M((\beta_k^+)^2 + (\beta_k^-)^2) c_\Delta^2 + c_\Delta^2 (M^2 \eta_{k-1}^2 \beta_k^2 + 2\delta_\theta (\beta_k + \beta_k^+) \times \\ &M \beta_k \eta_{k-1} + n(\delta_\theta^2 (\beta_k + \beta_k^+)^2 + \beta_k^2 g_2^2))) + \\ &nM^2 ((\beta_k^+)^2 + (\beta_k^-)^2) c_\Delta^4). \quad (22) \end{aligned}$$

Summing up the findings bounds (16), (21), (22) and taking the conditional expectation over σ -algebra \mathcal{F}_{k-1} , we derive the following from (15)

$$\mathbb{E}_{\mathcal{F}_{k-1}} \eta_k^2 \leq (1 - \mu\alpha) \eta_{k-1}^2 + 2\alpha h \eta_{k-1} + \alpha l. \quad (23)$$

Consider the condition $0 < \mu\alpha < 1$ of Lemma 2 from [20]. The right part holds since $\bar{\lambda}_2 \leq c_m$. The left part is satisfied by virtue of condition (11)–(12). Hence, taking the unconditional expectation of both sides of (23), we see that all conditions of Lemma 2 from [20] hold for $e_k = \sqrt{\mathbb{E} \eta_k^2}$.

This completes the proof of Theorem 1. \square