

Convergence of a Distributed Least Squares

Siyu Xie, Yaqi Zhang, and Lei Guo, *Fellow, IEEE*

Abstract—In this paper, we consider a least-squares (LS)-based distributed algorithm build on a sensor network to estimate an unknown parameter vector of a dynamical system, where each sensor in the network has partial information only but is allowed to communicate with its neighbors. Our main task is to generalize the well-known theoretical results on the traditional LS to the current distributed case by establishing both the upper bound of the accumulated regrets of the adaptive predictor and the convergence of the distributed LS estimator, with the following key features compared with the existing literature on distributed estimation: Firstly, our theory does not need the previously imposed independence, stationarity or Gaussian property on the system signals, and hence is applicable to stochastic systems with feedback. Secondly, the cooperative excitation condition introduced and used in this paper for the convergence of the distributed LS estimate is the weakest possible one, which shows that even if any individual sensor cannot estimate the unknown parameter by the traditional LS, the whole network can still fulfill the estimation task by the distributed LS.

Index Terms—Least squares, distributed estimation, convergence, diffusion strategies, cooperative excitation, regret, martingale theory

I. INTRODUCTION

Distributed estimation over sensor networks has received increasing research attention recently, and has been studied and used in many areas widely, e.g., collaborative spectral sensing in cognitive radio systems, target localization in biological networks, environmental monitoring, military surveillance, and so on [1]. Note that different cooperation strategies will lead to different distributed estimation algorithms, for examples, incremental least mean squares (LMS) [2], consensus LMS [3], [4], diffusion LMS [5], [6], incremental LS [7], consensus LS [8], diffusion LS [9]–[15], and distributed Kalman filter (KF) [16]–[19]. In our recent work (see e.g. [3]–[5]), we have given the stability and performance results for the consensus and diffusion LMS filters, without imposing the usual independence and stationarity assumptions for the system signals.

Note that when the unknown parameter is time-invariant, the LS algorithm may generate more accurate estimates in the transient phase and have faster convergence speed compared with LMS algorithm. This is one of the main motivations for us to consider the LS-based distributed estimation algorithm in this paper. Another reason for us to study this problem is that the existing convergence theory in the literature [7]–[15] can hardly be applied to non-independent and non-

stationary signals coming from practical complex systems where feedback loops exist inevitably.

Fortunately, in the traditional single sensor case, there is a vast literature on the convergence theory of the classical LS, which is indeed applicable to stochastic systems with feedback. In fact, motivated by the need to establish a rigorous theory for the well-known LS-based self-tuning regulators proposed by Åström and Wittenmark [20] in stochastic adaptive control, the convergence study of LS with possible stochastic feedback signals had received a great deal of attention in the literature, see e.g., [21]–[29]. At the same time, much effort had also been devoted to stochastic adaptive control, see e.g., [30]–[32]. Among the many significant contributions in this direction, here we only mention that Lai and Wei [25] established a celebrated convergence result under a weakest possible decaying excitation condition on the system signals, and Guo and Chen [29] and Guo [21] finally resolved the longstanding problem concerning the global stability and convergence of the LS-based self-tuning regulators.

In this paper, we will provide a theoretical analysis for a distributed LS algorithm of diffusion type [17], [18], where the diffusion strategy is designed via the so called covariance intersection fusion rule [33], [34]. In such a diffusion strategy, each node is only allowed to communicate with its neighbors, and both the estimates of the unknown parameter and the inverse of the covariance matrices are diffused between neighboring nodes. We will generalize the well-known convergence results on the classical LS by establishing both the upper bound of the accumulated regrets of the adaptive predictor and the convergence of the distributed LS estimator, with the following key features compared with the related results in the existing literature:

- Our theory does not need the usually assumed independence, stationarity or Gaussian property on the system signals, and hence does not exclude the applications of the theory to stochastic feedback systems.
- Our theory for the convergence of the distributed LS is established under a weakest possible cooperative excitation condition which is a natural extension of the single sensor case. The cooperative excitation condition introduced in this paper implies that even if any individual sensor is not able to estimate the unknown parameter, the distributed LS can still accomplish the estimation task.

The rest of the paper is organized as follows. In Section II, we present some preliminaries on notations and graph theory, the observation model, and the distributed LS algorithm studied in the paper. The main results are stated in Section III. In Section IV, we provide the proofs of the main results. Finally, some concluding remarks are given in Section V.

S. Y. Xie is with the Department of Electrical and Computer Engineering, Wayne State University, Detroit, MI 48202, USA. Email: syxie@wayne.edu.

Y. Q. Zhang and L. Guo are with Key Laboratory of Systems and Control, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, P. R. China. They are also with School of Mathematical Science, University of Chinese Academy of Sciences, Beijing 100049, P. R. China. Email : zhangyq@amss.ac.cn., lguo@amss.ac.cn.

This work was supported by the National Natural Science Foundation of China under grant 11688101.

II. PROBLEM FORMULATION

A. Basic Notations

Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times n}$ be two symmetric matrices with real entries, then $A \geq B$ ($A > B$) means $A - B$ is a positive semidefinite (definite) matrix. Also, let $\lambda_{\max}\{\cdot\}$ and $\lambda_{\min}\{\cdot\}$ denote the largest and the smallest eigenvalues of the corresponding matrix, respectively. For any matrix $X \in \mathbb{R}^{m \times n}$, $\|X\|$ denotes the operator norm induced by the Euclidean norm, i.e., $(\lambda_{\max}\{XX^\top\})^{\frac{1}{2}}$, where $(\cdot)^\top$ denotes the transpose operator. We use $\mathbb{E}[\cdot]$ to denote the mathematical expectation operator, and $\mathbb{E}[\cdot|\mathcal{F}_k]$ to denote the conditional mathematical expectation operator, where $\{\mathcal{F}_k\}$ is a sequence of nondecreasing σ -algebras [35]. We also use $\log(\cdot)$ to denote the natural logarithm function, and $\text{Tr}(\cdot)$ to denote the trace of the corresponding matrix. Through out the paper, $|\cdot|$ denotes the determinant of the corresponding matrix, which should not be confused with the absolute value of a scalar from the context.

Let $\{A_k, k \geq 0\}$ be a matrix sequence and $\{b_k, k \geq 0\}$ be a positive scalar sequence. Then by $A_k = O(b_k)$ we mean that there exists a constant $M > 0$ such that $\|A_k\| \leq Mb_k, \forall k \geq 0$, and by $A_k = o(b_k)$ we mean that $\lim_{k \rightarrow \infty} \|A_k\|/b_k = 0$.

B. Graph Theory

As usual, let the communication structure among sensors be represented by an undirected weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$, where $\mathcal{V} = \{1, 2, \dots, n\}$ is the set of sensors and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges. The structure of the graph \mathcal{G} is described by $\mathcal{A} = \{a_{ij}\}_{n \times n}$ which is called the weighted adjacency matrix, where $a_{ij} > 0$ if $(i, j) \in \mathcal{E}$ and $a_{ij} = 0$ otherwise. In this paper, we assume that the elements of the weighted matrix \mathcal{A} satisfy $a_{ij} = a_{ji}, \forall i, j = 1, \dots, n$, and $\sum_{j=1}^n a_{ij} = 1, \forall i = 1, \dots, n$. Thus the matrix \mathcal{A} is symmetric and doubly stochastic¹.

A path of length ℓ in the graph \mathcal{G} is a sequence of nodes $\{i_1, \dots, i_\ell\}$ subject to $(i_j, i_{j+1}) \in \mathcal{E}$, for $1 \leq j \leq \ell - 1$. The maximum value of the distances between any two nodes in the graph \mathcal{G} is called the diameter of \mathcal{G} . Here in this paper, we assume that the graph is connected, and denote the diameter of the graph \mathcal{G} as $D_{\mathcal{G}}$. Then $1 \leq D_{\mathcal{G}} < \infty$ holds. The set of neighbors of the sensor i is denoted as $\mathcal{N}_i = \{j \in \mathcal{V} | (j, i) \in \mathcal{E}\}$, and the sensor i can only share information with its neighboring sensors from \mathcal{N}_i .

C. Observation Model

Let us consider a sensor network consisting of n sensors. Assume that at each time instant k , each sensor $i \in \{1, \dots, n\}$ in the sensor network receives a noisy scalar measurement $y_{k+1,i}$ and an m -dimensional regressor $\varphi_{k,i} \in \mathbb{R}^m$. They are related by a typical linear stochastic regression model

$$y_{k+1,i} = \varphi_{k,i}^\top \theta + w_{k+1,i}, \quad k \geq 0, \quad (1)$$

where $w_{k+1,i}$ is a random noise process, and $\theta \in \mathbb{R}^m$ is an unknown parameter vector which needs to be estimated. We

¹A matrix is called doubly stochastic, if all elements are nonnegative, both the sum of each row and the sum of each column equal to 1.

assume that at any time instant $k \geq 1$, each sensor i uses both the observations $y_{j+1,i}$ and the regressors $\varphi_{j,i}$ ($j \leq k$) to estimate the unknown parameter θ .

D. Distributed LS Algorithm

We now consider the following basic class of distributed LS algorithms of diffusion type:

Algorithm 1 A Distributed LS algorithm

For any given sensor $i \in \{1, \dots, n\}$, begin with an initial estimate $\theta_{0,i} \in \mathbb{R}^m$, and an initial positive definite matrix $P_{0,i} \in \mathbb{R}^{m \times m}$. The algorithm is recursively defined at any time instant $k \geq 0$ as follows:

- 1: Adapt (generate $\bar{\theta}_{k+1,i}$ and $\bar{P}_{k+1,i}$ on the bases of $\theta_{k,i}, P_{k,i}$ and $\varphi_{k,i}, y_{k+1,i}$):

$$\bar{\theta}_{k+1,i} = \theta_{k,i} + b_{k,i} P_{k,i} \varphi_{k,i} (y_{k+1,i} - \varphi_{k,i}^\top \theta_{k,i}), \quad (2)$$

$$\bar{P}_{k+1,i} = P_{k,i} - b_{k,i} P_{k,i} \varphi_{k,i} \varphi_{k,i}^\top P_{k,i}, \quad (3)$$

$$b_{k,i} = (1 + \varphi_{k,i}^\top P_{k,i} \varphi_{k,i})^{-1}, \quad (4)$$

- 2: Combine (generate $P_{k+1,i}^{-1}$ and $\theta_{k+1,i}$ by a convex combination of $\bar{P}_{k+1,j}^{-1}$ and $\bar{\theta}_{k+1,j}$):

$$P_{k+1,i}^{-1} = \sum_{j \in \mathcal{N}_i} a_{ji} \bar{P}_{k+1,j}^{-1}, \quad (5)$$

$$\theta_{k+1,i} = P_{k+1,i} \sum_{j \in \mathcal{N}_i} a_{ji} \bar{P}_{k+1,j}^{-1} \bar{\theta}_{k+1,j}. \quad (6)$$

Remark 2.1: When $\mathcal{A} = I_n$, the distributed LS will degenerate to the classical LS at any sensor i . Note that for stochastic gradient-based [28] and LMS-based [36] distributed estimation algorithms, the communication complexity may be reduced. However, for those algorithms, the estimation error either converges slowly to zero or does not converge to zero at all. Therefore, there is a tradeoff between the complexity and the convergence rate of the distributed estimation algorithms. Moreover, the convergence rate would be ‘‘optimal’’ when $\bar{P}_{k,i}$ is chosen to be the form in the paper. Furthermore, some existing methods may be used to reduce the communication complexity and to make the algorithm suitable for higher dimensional signals, for examples, event-driven methods [37], partial diffusion methods [12], [13], and compressed methods [38] and so on.

III. THE MAIN RESULTS

A. Some Preliminaries

For the theoretical analysis, we need some standard assumptions on noise processes, regressors, and network topology.

Assumption 3.1: For each $i \in \{1, \dots, n\}$, the noise sequence $\{w_{k,i}, \mathcal{F}_k\}$ is a martingale difference (where $\{\mathcal{F}_k\}$ is a sequence of nondecreasing σ -algebras), and there exists a constant $\beta > 2$ such that $\sup_{k \geq 0} \mathbb{E}[|w_{k+1,i}|^\beta | \mathcal{F}_k] < \infty$, a.s.

Assumption 3.2: For each $i \in \{1, \dots, n\}$, the regressor sequence $\{\varphi_{k,i}, \mathcal{F}_k\}$ is an adapted sequence.

Assumption 3.3: The graph \mathcal{G} is connected.

Remark 3.1: From *Lemma 8.1.2* in [39], it is not difficult to see that for any two nodes i and j , there exists a path from i to j with length not larger than ℓ if and only if the (i, j) th entry of the matrix \mathcal{A}^ℓ is positive. From this, it is easy to see that each entry of the matrix \mathcal{A}^ℓ will be positive when ℓ is not smaller than the diameter of the graph \mathcal{G} , i.e., $D_{\mathcal{G}}$, see also [18].

B. Theoretical Results

Here we first introduce the following notations:

$$\begin{aligned} \mathbf{Y}_k &\triangleq \text{col}\{y_{k,1}, \dots, y_{k,n}\}, \quad \Phi_k \triangleq \text{diag}\{\varphi_{k,1}, \dots, \varphi_{k,n}\}, \\ \mathbf{W}_k &\triangleq \text{col}\{w_{k,1}, \dots, w_{k,n}\}, \quad \Theta \triangleq \text{col}\{\underbrace{\theta, \dots, \theta}_n\}, \\ \Theta_k &\triangleq \text{col}\{\theta_{k,1}, \dots, \theta_{k,n}\}, \quad \bar{\Theta}_k \triangleq \text{col}\{\bar{\theta}_{k,1}, \dots, \bar{\theta}_{k,n}\}, \\ \tilde{\Theta}_k &\triangleq \text{col}\{\tilde{\theta}_{k,1}, \dots, \tilde{\theta}_{k,n}\}, \quad \text{where } \tilde{\theta}_{k,i} = \theta - \theta_{k,i}, \\ \tilde{\bar{\Theta}}_k &\triangleq \text{col}\{\tilde{\bar{\theta}}_{k,1}, \dots, \tilde{\bar{\theta}}_{k,n}\}, \quad \text{where } \tilde{\bar{\theta}}_{k,i} = \theta - \bar{\theta}_{k,i}, \\ P_k &\triangleq \text{diag}\{P_{k,1}, \dots, P_{k,n}\}, \quad \bar{P}_k \triangleq \text{diag}\{\bar{P}_{k,1}, \dots, \bar{P}_{k,n}\}, \\ b_k &\triangleq \text{diag}\{b_{k,1}, \dots, b_{k,n}\}, \quad c_k \triangleq b_k \otimes I_m, \quad \mathcal{A} \triangleq \mathcal{A} \otimes I_m, \end{aligned}$$

where $\text{col}\{\dots\}$ denotes a vector by stacking the specified vectors, $\text{diag}\{\dots\}$ is used in a non-standard manner which means that $m \times 1$ column vectors are combined “in a diagonal manner” resulting in a $mn \times n$ matrix, and \otimes is the Kronecker product. Then (1) can be rewritten in the following compact form:

$$\mathbf{Y}_{k+1} = \Phi_k^\top \Theta + \mathbf{W}_{k+1}, \quad (7)$$

Similarly, for the distributed LS algorithm we have

$$\left\{ \begin{array}{l} \bar{\Theta}_{k+1} = \Theta_k + c_k P_k \Phi_k (\mathbf{Y}_{k+1} - \Phi_k^\top \Theta_k), \\ \bar{P}_{k+1} = P_k - c_k P_k \Phi_k \Phi_k^\top P_k, \\ b_k = (I_n + \Phi_k^\top P_k \Phi_k)^{-1}, \\ c_k = b_k \otimes I_m, \\ \text{vec}\{P_{k+1}^{-1}\} = \mathcal{A} \text{vec}\{\bar{P}_{k+1}^{-1}\}, \\ \Theta_{k+1} = P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} \bar{\Theta}_{k+1}, \end{array} \right. \quad (8)$$

where $\text{vec}\{\cdot\}$ denotes the operator that stacks the blocks of a block diagonal matrix on top of each other.

Since $\tilde{\Theta}_k = \Theta - \Theta_k$ and $\tilde{\bar{\Theta}}_k = \Theta - \bar{\Theta}_k$ by definition, substituting (7) into (8), we can get

$$\begin{aligned} \tilde{\Theta}_{k+1} &= \Theta - \Theta_k - c_k P_k \Phi_k (\Phi_k^\top \Theta + \mathbf{W}_{k+1} - \Phi_k^\top \Theta_k) \\ &= (I_{mn} - c_k P_k \Phi_k \Phi_k^\top) \tilde{\Theta}_k - c_k P_k \Phi_k \mathbf{W}_{k+1} \\ &= \bar{P}_{k+1} P_k^{-1} \tilde{\Theta}_k - c_k P_k \Phi_k \mathbf{W}_{k+1}. \end{aligned}$$

Note also that

$$\begin{aligned} &P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} \Theta \\ &= \text{col}\left\{ P_{k+1,1} \sum_{j \in \mathcal{N}_1} a_{j1} \bar{P}_{k+1,j}^{-1} \theta, \dots, P_{k+1,n} \sum_{j \in \mathcal{N}_n} a_{jn} \bar{P}_{k+1,j}^{-1} \theta \right\} \\ &= \Theta. \end{aligned}$$

Then we have

$$\begin{aligned} \tilde{\Theta}_{k+1} &= \Theta - P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} \bar{\Theta}_{k+1} = P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} \tilde{\Theta}_{k+1} \\ &= P_{k+1} \mathcal{A} P_k^{-1} \tilde{\Theta}_k \\ &\quad - P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} c_k P_k \Phi_k \mathbf{W}_{k+1}. \end{aligned} \quad (9)$$

Before establishing the convergence of the distributed LS, we first present a critical theorem, which requires no excitation conditions on the regression process $\varphi_{k,i}$.

Theorem 3.1: Let *Assumptions 3.1* and *3.2* be satisfied, we have as $t \rightarrow \infty$,

$$\begin{aligned} 1) \quad &\sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k b_k \Phi_k^\top \tilde{\Theta}_k = O(\log(r_t)), \quad \text{a.s.}, \\ 2) \quad &\tilde{\Theta}_{t+1}^\top P_{t+1}^{-1} \tilde{\Theta}_{t+1} = O(\log(r_t)), \quad \text{a.s.}, \end{aligned}$$

where

$$r_t = \lambda_{\max}\{P_0^{-1}\} + \sum_{i=1}^n \sum_{k=0}^t \|\varphi_{k,i}\|^2. \quad (10)$$

From this, we can obtain the upper bound of the accumulated regrets for the distributed LS-based adaptive predictor. For any $i \in \{1, \dots, n\}$, and at any time instant $k \geq 1$, the best prediction to the future observation $y_{k+1,i}$ is the conditional mathematical expectation $\mathbb{E}[y_{k+1,i} | \mathcal{F}_k] = \varphi_{k,i}^\top \theta$, since the noise is a martingale difference sequence with second moment. Unfortunately, this optimal predictor is unavailable because θ is unknown. A natural way is to construct an adaptive predictor $\hat{y}_{k+1,i}$ by using the online distributed LS estimate $\theta_{k,i}$, i.e., $\hat{y}_{k+1,i} = \varphi_{k,i}^\top \theta_{k,i}$. The error between the best predictor and the adaptive predictor is referred to as the regret denoted by

$$R_{k,i} = (\mathbb{E}[y_{k+1,i} | \mathcal{F}_k] - \hat{y}_{k+1,i})^2, \quad (11)$$

which may not be zero and even may not be small in sample paths due to the persistent disturbance of the unpredictable noises in the model.

However, one may evaluate the averaged regrets defined as follows:

$$\frac{1}{nt} \sum_{i=1}^n \sum_{k=0}^t R_{k,i},$$

which we are going to show tends to zero as t increases to infinity under essentially no excitation conditions on the regressors, see *Theorem 3.2* below. This is a celebrated property that has been widely studied in distributed online learning and optimization problems [40], [41], but under rather restrictive assumptions such as boundedness, stationarity or independence on the system signals.

Theorem 3.2: Let *Assumptions 3.1* and *3.2* be satisfied. Then the sample paths of the accumulated regrets have the following bound as $t \rightarrow \infty$:

$$\sum_{i=1}^n \sum_{k=0}^t R_{k,i} = O(\log(r_t)), \quad \text{a.s.}, \quad (12)$$

provided that $\Phi_t^\top P_t \Phi_t = O(1)$, a.s.

Remark 3.2: We remark that the order $O(\log(r_t))$ for the accumulated regrets may be shown to be the best possible among all adaptive predictors in a certain sense, as is already known in the traditional single sensor case, see [42]. The precise constant in $O(\cdot)$ may also be determined if we have

further conditions on the regressors, see *Corollary 3.3* in [21] in the single sensor case.

From *Theorem 3.1*, we can also obtain the strong consistency of the distributed LS to guarantee the generalization ability of learning.

Theorem 3.3: Let *Assumptions 3.1-3.3* be satisfied, we have as $t \rightarrow \infty$,

$$\|\tilde{\Theta}_{t+1}\|^2 = O\left(\frac{\log(r_t)}{\lambda_{\min}^{n,t}}\right), \quad \text{a.s.}, \quad (13)$$

where r_t is defined by (10) and

$$\lambda_{\min}^{n,t} = \lambda_{\min} \left\{ \sum_{j=1}^n P_{0,j}^{-1} + \sum_{j=1}^n \sum_{k=0}^{t-D_G+1} \varphi_{k,j} \varphi_{k,j}^\top \right\}. \quad (14)$$

Remark 3.3: *Theorem 3.3* shows that if

$$\lim_{t \rightarrow \infty} \frac{\log(r_t)}{\lambda_{\min}^{n,t}} = 0, \quad \text{a.s.}, \quad (15)$$

then the distributed LS estimate Θ_t will converge to the true unknown parameter. We may name (15) as a cooperative excitation condition. In the traditional single sensor case (where $n = 1$, $D_G = 1$), (15) reduces to the well-known Lai-Wei excitation condition, which is known to be the weakest possible data condition for the convergence of the classical LS estimates [25], and is much weaker than the well-known persistence of excitation (PE) condition usually used in the parameter estimation of finite-dimensional linear control systems.

Moreover, it is easy to convince oneself that the cooperative excitation condition (15) will make it possible for the distributed LS to consistently estimate the unknown parameter, even if any individual sensor cannot due to lack of suitable excitation, thanks to the cooperative nature of the excitation condition (15). Finally, we remark that the verification of (15) is straightforward in the ergodic case. For more general correlated non-stationary signals from control systems, the verification of (15) may be conducted following a similar way as that for the traditional single sensor case (see, [28]).

Furthermore, the convergence rate established in *Theorem 3.3* is essentially in terms of the increase of the number of observations rather than the number of iterations in computation.

Remark 3.4: Let us now compare the above distributed LS algorithm with centralized methods whereby, at each time instant k , all the n sensors transmit their raw data $\{y_{k+1,i}, \varphi_{k,i}\}$ to a fusion center for processing to obtain a centralized estimate. Although the centralized algorithm may have some advantages over the distributed algorithm in terms of communication complexity, it also has some drawbacks compared with the distributed case. Firstly, the distributed methods may have stronger structural robustness compared with the centralized ones. This is because the centralized algorithm will fail once the fusion center is broken down by outside attacks or some sensors lost the connection to the fusion center, while the distributed algorithm can still estimate the unknown parameters even if the communications among some sensors are interrupted, as long as the network connectivity is maintained. Secondly, if the fusion center

is far away from some sensors, the communications with the fusion center may not be feasible, and the transmission of observations and regression vectors may compromise the safety and privacy of the system even if the communication is possible. Hence, our distributed estimation problem is not a purely computational problem.

IV. PROOFS OF THE MAIN RESULTS

A. Proof of Theorem 3.1

To prove *Theorem 3.1*, we need to establish several lemmas first. The first lemma below is a key inequality on convex combination of nonnegative definite matrices.

Lemma 4.1: For any adjacency matrix $\mathcal{A} = \{a_{ij}\} \in \mathbb{R}^{n \times n}$, denote $\mathcal{A} = \mathcal{A} \otimes I_m$, and for any nonnegative definite matrices $Q_i \in \mathbb{R}^{m \times m}$, $i = 1, \dots, n$, denote $Q = \text{diag}\{Q_1, \dots, Q_n\}$, and $Q' = \text{diag}\{Q'_1, \dots, Q'_n\}$, where $Q'_i = \sum_{j=1}^n a_{ji} Q_j$. Then the following inequality holds:

$$\mathcal{A} Q \mathcal{A} \leq Q'. \quad (16)$$

Proof: By the definition of \mathcal{A} and Q , we can get that

$$\mathcal{A} Q \mathcal{A} = \begin{pmatrix} \sum_{j=1}^n a_{1j} a_{j1} Q_j & \cdots & \sum_{j=1}^n a_{1j} a_{jn} Q_j \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n a_{nj} a_{j1} Q_j & \cdots & \sum_{j=1}^n a_{nj} a_{jn} Q_j \end{pmatrix}.$$

In order to prove the inequality, we only need to prove that for any unit column vector $x \in \mathbb{R}^{mn}$ with $\|x\| = 1$, $x^\top \mathcal{A} Q \mathcal{A} x \leq x^\top Q' x$ holds. Denote $x = \text{col}\{x_1, x_2, \dots, x_n\}$ with $x_i \in \mathbb{R}^m$, then by the Schwarz inequality and noticing that $Q_j \geq 0$, $\sum_{j=1}^n a_{ij} = 1$, and $a_{ji} = a_{ij}$, ($i, j = 1, \dots, n$), we have

$$\begin{aligned} x^\top \mathcal{A} Q \mathcal{A} x &= \sum_{p=1}^n \sum_{q=1}^n \sum_{j=1}^n a_{pj} a_{jq} x_p^\top Q_j x_q \\ &= \sum_{p=1}^n \sum_{q=1}^n \sum_{j=1}^n \sqrt{a_{pj} a_{jq}} x_p^\top Q_j^{\frac{1}{2}} \cdot \sqrt{a_{pj} a_{jq}} Q_j^{\frac{1}{2}} x_q \\ &\leq \left\{ \sum_{p=1}^n \sum_{q=1}^n \sum_{j=1}^n a_{pj} a_{jq} x_p^\top Q_j x_p \right\}^{\frac{1}{2}} \\ &\quad \cdot \left\{ \sum_{p=1}^n \sum_{q=1}^n \sum_{j=1}^n a_{pj} a_{jq} x_q^\top Q_j x_q \right\}^{\frac{1}{2}} \\ &= \left\{ \sum_{p=1}^n \sum_{j=1}^n a_{pj} x_p^\top Q_j x_p \right\}^{\frac{1}{2}} \left\{ \sum_{q=1}^n \sum_{j=1}^n a_{jq} x_q^\top Q_j x_q \right\}^{\frac{1}{2}} \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ji} x_i^\top Q_j x_i = x^\top Q' x, \end{aligned}$$

which completes the proof. \blacksquare

By *Lemma 4.1*, we can obtain the following result:

Lemma 4.2: For any adjacency matrix $\mathcal{A} = \{a_{ij}\} \in \mathbb{R}^{n \times n}$, denote $\mathcal{A} = \mathcal{A} \otimes I_m$. Then for any $k \geq 1$,

$$\mathcal{A} \bar{P}_{k+1}^{-1} \mathcal{A} \leq P_{k+1}^{-1}, \quad (17)$$

and

$$\mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \leq \bar{\mathbf{P}}_{k+1}, \quad (18)$$

holds, where $\bar{\mathbf{P}}_{k+1}$ and \mathbf{P}_{k+1} are defined in (8).

Proof: By taking $Q_i = \bar{P}_{k+1,i}^{-1} \geq 0$ and noticing $P_{k+1,i}^{-1} = \sum_{j=1}^n a_{ji} \bar{P}_{k+1,j}^{-1} = Q'_i$, we know from *Lemma 4.1* that

$$\mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \mathcal{A} \leq \mathbf{P}_{k+1}^{-1},$$

holds. To prove (18), we first assume that \mathcal{A} is invertible. Then by *Lemma A.1* in Appendix A, it is easy to see that

$$\mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \leq \bar{\mathbf{P}}_{k+1}.$$

Next, we consider the case where \mathcal{A} is not invertible. Since the number of eigenvalues of the matrix \mathcal{A} is finite, then exists a constant $\varepsilon^* \in (0, 1)$ such that the perturbed adjacency matrix $\mathcal{A}^\varepsilon = \mathcal{A} + \varepsilon I_{mn} = \{a_{ij}^\varepsilon\}$ will be invertible for any $0 < \varepsilon < \varepsilon^*$. By the definition of \mathcal{A}^ε , we know that \mathcal{A}^ε is symmetric and the sums of each columns and rows of the matrix \mathcal{A}^ε are all $1 + \varepsilon$. Then we define $(P_{k+1,i}^\varepsilon)^{-1} = \sum_{j=1}^n a_{ji}^\varepsilon \bar{P}_{k+1,j}^{-1}$, and we can denote $\mathbf{P}_{k+1}^\varepsilon = \text{diag}\{P_{k+1,1}^\varepsilon, \dots, P_{k+1,n}^\varepsilon\}$ since $(P_{k+1,i}^\varepsilon)^{-1}$ defined above is invertible. Similar to the proof of *Lemma 4.1*, for any unit column vector $x \in \mathbb{R}^{mn}$, we have

$$x^\top \mathcal{A}^\varepsilon \bar{\mathbf{P}}_{k+1}^{-1} \mathcal{A}^\varepsilon x \leq (1 + \varepsilon) x^\top (\mathbf{P}_{k+1}^\varepsilon)^{-1} x.$$

Consequently, we have $\mathcal{A}^\varepsilon \bar{\mathbf{P}}_{k+1}^{-1} \mathcal{A}^\varepsilon \leq (1 + \varepsilon) (\mathbf{P}_{k+1}^\varepsilon)^{-1}$. Since \mathcal{A}^ε is invertible, we know from *Lemma A.1* in Appendix A that $\mathcal{A}^\varepsilon \mathbf{P}_{k+1}^\varepsilon \mathcal{A}^\varepsilon \leq (1 + \varepsilon) \bar{\mathbf{P}}_{k+1}$. By taking $\varepsilon \rightarrow 0$ on both sides of the above equation, we can obtain that

$$\lim_{\varepsilon \rightarrow 0} \mathcal{A}^\varepsilon \mathbf{P}_{k+1}^\varepsilon \mathcal{A}^\varepsilon = \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \leq \lim_{\varepsilon \rightarrow 0} (1 + \varepsilon) \bar{\mathbf{P}}_{k+1} = \bar{\mathbf{P}}_{k+1}.$$

This completes the proof. \blacksquare

To accomplish the proof of *Theorem 3.1*, we also need the following inequality:

Lemma 4.3: For any adjacency matrix $\mathcal{A} = \{a_{ij}\} \in \mathbb{R}^{n \times n}$, and for any $k \geq 1$,

$$|\bar{\mathbf{P}}_{k+1}^{-1}| \leq |\mathbf{P}_{k+1}^{-1}|, \quad (19)$$

holds, where $\bar{\mathbf{P}}_{k+1}$ and \mathbf{P}_{k+1} are defined in (8).

Proof: By Ky Fan convex theorem [43] and noticing the definitions of $\bar{\mathbf{P}}_{k+1}$, \mathbf{P}_{k+1} , and $\mathcal{A} = \{a_{ij}\}$, we can see that

$$\begin{aligned} |\mathbf{P}_{k+1}^{-1}| &= \prod_{i=1}^n \left| \sum_{j=1}^n a_{ji} \bar{P}_{k+1,j}^{-1} \right| \\ &\geq \prod_{i=1}^n |\bar{P}_{k+1,1}^{-1}|^{a_{1i}} |\bar{P}_{k+1,2}^{-1}|^{a_{2i}} \dots |\bar{P}_{k+1,n}^{-1}|^{a_{ni}} \\ &= |\bar{P}_{k+1,1}^{-1}|^{\sum_{i=1}^n a_{1i}} |\bar{P}_{k+1,2}^{-1}|^{\sum_{i=1}^n a_{2i}} \dots |\bar{P}_{k+1,n}^{-1}|^{\sum_{i=1}^n a_{ni}} \\ &= |\bar{P}_{k+1,1}^{-1}| \cdot |\bar{P}_{k+1,2}^{-1}| \dots |\bar{P}_{k+1,n}^{-1}| = |\bar{\mathbf{P}}_{k+1}^{-1}|, \end{aligned}$$

which completes the proof. \blacksquare

To prove *Theorem 3.1*, we also need the following critical lemma:

Lemma 4.4: Let *Assumptions 3.1* and *3.2* be satisfied. Then the distributed LS defined by (7) and (8) satisfies the following relationship as $t \rightarrow \infty$:

$$\begin{aligned} &\tilde{\Theta}_{t+1}^\top \mathbf{P}_{t+1}^{-1} \tilde{\Theta}_{t+1} + [1 + o(1)] \sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k \mathbf{b}_k \Phi_k^\top \tilde{\Theta}_k \\ &+ [1 + o(1)] \sum_{k=0}^t \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &\leq \sigma_w \log(|\mathbf{P}_{t+1}^{-1}|) + o(\log(|\mathbf{P}_{t+1}^{-1}|)) + O(1), \quad \text{a.s.}, \quad (20) \end{aligned}$$

where $\Delta_{k+1} \triangleq \bar{\mathbf{P}}_{k+1} - \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \geq 0$ by *Lemma 4.2*, and $\sigma_w = \sum_{i=1}^n \sigma_i^2$, $\sigma_i^2 = \sup_{k \geq 0} \mathbb{E}[w_{k+1,i}^2 | \mathcal{F}_k]$.

Proof: Since $\mathbf{b}_k = (I_n + \Phi_k^\top \mathbf{P}_k \Phi_k)^{-1}$ and $\mathbf{c}_k = \mathbf{b}_k \otimes I_m$, then by (9), we know that

$$\tilde{\Theta}_{k+1} = \mathbf{P}_{k+1} \mathcal{A} \mathbf{P}_k^{-1} \tilde{\Theta}_k - \mathbf{P}_{k+1} \mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1}.$$

Hence, we have the following expansion for the stochastic Lyapunov function $V_k = \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \tilde{\Theta}_k$:

$$\begin{aligned} V_{k+1} &= \tilde{\Theta}_{k+1}^\top \mathbf{P}_{k+1}^{-1} \tilde{\Theta}_{k+1} \\ &= (\tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} - \mathbf{W}_{k+1}^\top \Phi_k^\top \mathbf{P}_k \mathbf{c}_k \bar{\mathbf{P}}_{k+1}^{-1} \mathcal{A} \mathbf{P}_{k+1}) \\ &\quad \cdot (\mathcal{A} \mathbf{P}_k^{-1} \tilde{\Theta}_k - \mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1}) \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &\quad - 2 \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1} \\ &\quad + \mathbf{W}_{k+1}^\top \Phi_k^\top \mathbf{P}_k \mathbf{c}_k \bar{\mathbf{P}}_{k+1}^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \\ &\quad \cdot \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1}. \quad (21) \end{aligned}$$

Now, we proceed to estimate the right-hand-side (RHS) of (21) term by term. Firstly, we know that

$$\begin{aligned} &\tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \bar{\mathbf{P}}_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} (\mathbf{P}_k - \mathbf{P}_k \Phi_k \mathbf{b}_k \Phi_k^\top \mathbf{P}_k) \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &\quad - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \tilde{\Theta}_k - \tilde{\Theta}_k^\top \Phi_k \mathbf{b}_k \Phi_k^\top \tilde{\Theta}_k - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k \\ &= V_k - \tilde{\Theta}_k^\top \Phi_k \mathbf{b}_k \Phi_k^\top \tilde{\Theta}_k - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \mathbf{P}_k^{-1} \tilde{\Theta}_k. \quad (22) \end{aligned}$$

Moreover, by the (block) diagonal property of \mathbf{b}_k , \mathbf{c}_k , \mathbf{P}_k and Φ_k , we have $\mathbf{c}_k \mathbf{P}_k = \mathbf{P}_k \mathbf{c}_k$, $\Phi_k^\top \mathbf{c}_k = \mathbf{b}_k \Phi_k^\top$, and $\mathbf{c}_k \Phi_k = \Phi_k \mathbf{b}_k$. By the matrix inversion lemma [36], we have

$$\bar{\mathbf{P}}_{k+1}^{-1} = \mathbf{P}_k^{-1} + \Phi_k \Phi_k^\top.$$

Thus, we can estimate the second term on the RHS of (21) as follows:

$$\begin{aligned} &\tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \bar{\mathbf{P}}_{k+1}^{-1} \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1} \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \mathbf{c}_k \Phi_k \mathbf{W}_{k+1} \\ &\quad + \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \Phi_k \Phi_k^\top \mathbf{c}_k \mathbf{P}_k \Phi_k \mathbf{W}_{k+1} \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \mathbf{c}_k \Phi_k \mathbf{W}_{k+1} \\ &\quad + \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \Phi_k \mathbf{W}_{k+1} \\ &\quad - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \Phi_k \mathbf{b}_k \mathbf{W}_{k+1} \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \mathcal{A} \mathbf{P}_{k+1} \mathcal{A} \Phi_k \mathbf{W}_{k+1} \\ &= \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \bar{\mathbf{P}}_{k+1} \Phi_k \mathbf{W}_{k+1} - \tilde{\Theta}_k^\top \mathbf{P}_k^{-1} \Delta_{k+1} \Phi_k \mathbf{W}_{k+1}. \quad (23) \end{aligned}$$

As for the last term on the RHS of (21), by $\mathcal{A}P_{k+1}\mathcal{A} \leq \bar{P}_{k+1}$, we can estimate it as follows:

$$\begin{aligned}
 & \mathbf{W}_{k+1}^\top \Phi_k^\top P_k c_k \bar{P}_{k+1}^{-1} \mathcal{A} P_{k+1} \mathcal{A} \bar{P}_{k+1}^{-1} c_k P_k \Phi_k \mathbf{W}_{k+1} \\
 & \leq \mathbf{W}_{k+1}^\top \Phi_k^\top P_k c_k (P_k^{-1} + \Phi_k \Phi_k^\top) c_k P_k \Phi_k \mathbf{W}_{k+1} \\
 & = \mathbf{W}_{k+1}^\top \Phi_k^\top P_k c_k^2 \Phi_k \mathbf{W}_{k+1} \\
 & \quad + \mathbf{W}_{k+1}^\top \Phi_k^\top P_k c_k \Phi_k \Phi_k^\top c_k P_k \Phi_k \mathbf{W}_{k+1} \\
 & = \mathbf{W}_{k+1}^\top b_k^2 \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1} \\
 & \quad + \mathbf{W}_{k+1}^\top (I_n + \Phi_k^\top P_k \Phi_k) b_k^2 \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1} \\
 & \quad - \mathbf{W}_{k+1}^\top b_k^2 \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1} \\
 & = \mathbf{W}_{k+1}^\top b_k \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1}. \tag{24}
 \end{aligned}$$

By (21)-(24), we have

$$\begin{aligned}
 V_{k+1} & \leq V_k - \tilde{\Theta}_k^\top \Phi_k b_k \Phi_k^\top \tilde{\Theta}_k - \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} P_k^{-1} \tilde{\Theta}_k \\
 & \quad - 2 \tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \mathbf{W}_{k+1} \\
 & \quad + 2 \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} \Phi_k \mathbf{W}_{k+1} \\
 & \quad + \mathbf{W}_{k+1}^\top b_k \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1}. \tag{25}
 \end{aligned}$$

Summing from $k = 0$ to t yields

$$\begin{aligned}
 & V_{t+1} + \sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k b_k \Phi_k^\top \tilde{\Theta}_k + \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} P_k^{-1} \tilde{\Theta}_k \\
 & \leq V_0 - 2 \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \mathbf{W}_{k+1} \\
 & \quad - 2 \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} (-\Delta_{k+1}) \Phi_k \mathbf{W}_{k+1} \\
 & \quad + \sum_{k=0}^t \mathbf{W}_{k+1}^\top b_k \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1}. \tag{26}
 \end{aligned}$$

Next, we estimate the last three terms on the RHS of (26) separately. By *Assumptions 3.1* and *3.2*, and $\tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \in \mathcal{F}_k$, $\tilde{\Theta}_k^\top P_k^{-1} (-\Delta_{k+1}) \Phi_k \in \mathcal{F}_k$, we can use the martingale estimation theorem (*Theorem 2.8* in [28]) to get the following estimation for any $\delta > 0$:

$$\begin{aligned}
 & \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \mathbf{W}_{k+1} \\
 & = O\left(\left\{\sum_{k=0}^t \|\tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k\|^2\right\}^{\frac{1}{2}+\delta}\right), \quad \text{a.s.}, \tag{27}
 \end{aligned}$$

and

$$\begin{aligned}
 & \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} (-\Delta_{k+1}) \Phi_k \mathbf{W}_{k+1} \\
 & = O\left(\left\{\sum_{k=0}^t \|\tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} \Phi_k\|^2\right\}^{\frac{1}{2}+\delta}\right), \quad \text{a.s.} \tag{28}
 \end{aligned}$$

To further analyze (27) and (28), we note from the definitions of \bar{P}_{k+1} and b_k that

$$\begin{aligned}
 & P_k^{-1} \bar{P}_{k+1} \Phi_k = \Phi_k - c_k \Phi_k \Phi_k^\top P_k \Phi_k \\
 & = \Phi_k - c_k \Phi_k (I_n + \Phi_k^\top P_k \Phi_k) + c_k \Phi_k = \Phi_k b_k.
 \end{aligned}$$

Hence, it is easy to see that

$$\begin{aligned}
 \|\tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k\|^2 & = \tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \Phi_k^\top \bar{P}_{k+1} P_k^{-1} \tilde{\Theta}_k \\
 & = \tilde{\Theta}_k^\top \Phi_k b_k^2 \Phi_k^\top \tilde{\Theta}_k \leq \tilde{\Theta}_k^\top \Phi_k b_k \Phi_k^\top \tilde{\Theta}_k,
 \end{aligned}$$

By taking $0 < \delta < \frac{1}{2}$, we have from (27) that

$$\begin{aligned}
 & \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k \mathbf{W}_{k+1} \\
 & = O(1) + o\left(\left\{\sum_{k=0}^t \|\tilde{\Theta}_k^\top P_k^{-1} \bar{P}_{k+1} \Phi_k\|^2\right\}\right) \\
 & = O(1) + o\left(\left\{\sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k b_k \Phi_k^\top \tilde{\Theta}_k\right\}\right), \quad \text{a.s.} \tag{29}
 \end{aligned}$$

Moreover, since $\Delta_{k+1} = \bar{P}_{k+1} - \mathcal{A}P_{k+1}\mathcal{A} \leq \bar{P}_{k+1}$, then

$$\begin{aligned}
 & \Delta_{k+1}^{\frac{1}{2}} \Phi_k \Phi_k^\top \Delta_{k+1}^{\frac{1}{2}} \leq \lambda_{\max}\{\Phi_k^\top \Delta_{k+1} \Phi_k\} \cdot I_{mn} \\
 & \leq \lambda_{\max}\{\Phi_k^\top \bar{P}_{k+1} \Phi_k\} \cdot I_{mn} \\
 & = \lambda_{\max}\{\Phi_k^\top (P_k - c_k P_k \Phi_k \Phi_k^\top P_k) \Phi_k\} \cdot I_{mn} \\
 & = \lambda_{\max}\{b_k \Phi_k^\top P_k \Phi_k\} \cdot I_{mn} < I_{mn}.
 \end{aligned}$$

Hence, we have $\Delta_{k+1} \Phi_k \Phi_k^\top \Delta_{k+1} \leq \Delta_{k+1}$, and so we have

$$\begin{aligned}
 \|\tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} \Phi_k\|^2 & = \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} \Phi_k \Phi_k^\top \Delta_{k+1} P_k^{-1} \tilde{\Theta}_k \\
 & \leq \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} P_k^{-1} \tilde{\Theta}_k. \tag{30}
 \end{aligned}$$

By taking $0 < \delta < \frac{1}{2}$, we know from (28) that

$$\begin{aligned}
 & \sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} (-\Delta_{k+1}) \Phi_k \mathbf{W}_{k+1} \\
 & = O(1) + o\left(\left\{\sum_{k=0}^t \|\tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} \Phi_k\|^2\right\}\right) \\
 & = O(1) + o\left(\left\{\sum_{k=0}^t \tilde{\Theta}_k^\top P_k^{-1} \Delta_{k+1} P_k^{-1} \tilde{\Theta}_k\right\}\right), \quad \text{a.s.} \tag{31}
 \end{aligned}$$

We now proceed to estimate the last term in (26). Firstly, we know that

$$\begin{aligned}
 & \mathbf{W}_{k+1}^\top b_k \Phi_k^\top P_k \Phi_k \mathbf{W}_{k+1} \leq \|b_k \Phi_k^\top P_k \Phi_k\| \cdot \|\mathbf{W}_{k+1}\|^2 \\
 & = \lambda_{\max}\{b_k \Phi_k^\top P_k \Phi_k\} \cdot \left\{\sum_{i=1}^n w_{k+1,i}^2\right\}. \tag{32}
 \end{aligned}$$

Following a similar proof idea as in the traditional single sensor case ([25], see also [28]), from $\bar{P}_{k+1} = P_k - c_k P_k \Phi_k \Phi_k^\top P_k$, we have $P_k^{-1} = \bar{P}_{k+1}^{-1} (I_{mn} - c_k P_k \Phi_k \Phi_k^\top)$. By taking determinants on both sides of the above identity, and noticing $0 \leq b_k \Phi_k^\top P_k \Phi_k \leq I_n$, we have

$$\begin{aligned}
 |P_k^{-1}| & = |\bar{P}_{k+1}^{-1}| \cdot |I_n - b_k \Phi_k^\top P_k \Phi_k| \\
 & \leq |\bar{P}_{k+1}^{-1}| \cdot (1 - \lambda_{\max}\{b_k \Phi_k^\top P_k \Phi_k\}).
 \end{aligned}$$

Moreover, we know from *Lemma 4.3* that

$$\lambda_{\max}\{b_k \Phi_k^\top P_k \Phi_k\} \leq \frac{|\bar{P}_{k+1}^{-1}| - |P_k^{-1}|}{|\bar{P}_{k+1}^{-1}|} \leq \frac{|P_{k+1}^{-1}| - |P_k^{-1}|}{|P_{k+1}^{-1}|}.$$

Therefore

$$\begin{aligned} \sum_{k=0}^t \lambda_{\max}\{\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k\} &\leq \sum_{k=0}^t \frac{|\mathbf{P}_{k+1}^{-1}| - |\mathbf{P}_k^{-1}|}{|\mathbf{P}_{k+1}^{-1}|} \\ &\leq \sum_{k=0}^t \int_{|\mathbf{P}_k^{-1}|}^{|\mathbf{P}_{k+1}^{-1}|} \frac{dx}{x} = \log(|\mathbf{P}_{t+1}^{-1}|) - \log(|\mathbf{P}_0^{-1}|). \end{aligned} \quad (33)$$

Consequently, by using the martingale estimation theorem (Theorem 2.8 in [28]), we have for any $\forall \eta > 0$,

$$\begin{aligned} \sum_{k=0}^t \lambda_{\max}\{\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k\} \left\{ \sum_{i=1}^n w_{k+1,i}^2 - \mathbb{E} \left[\sum_{i=1}^n w_{k+1,i}^2 | \mathcal{F}_k \right] \right\} \\ = O \left(S_t \left(\frac{\beta}{2} \right) \left\{ \log \left(S_t \left(\frac{\beta}{2} \right) + e \right) \right\}^{\frac{2}{\beta} + \eta} \right), \quad \text{a.s.}, \end{aligned} \quad (34)$$

where

$$S_t \left(\frac{\beta}{2} \right) \triangleq \left[\sum_{k=0}^t (\lambda_{\max}\{\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k\})^{\frac{\beta}{2}} \right]^{\frac{2}{\beta}}.$$

Since $\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k \leq I_n$ and $\frac{\beta}{2} > 1$, we have from (33) that

$$S_t \left(\frac{\beta}{2} \right) = O(1) + O((\log |\mathbf{P}_{t+1}^{-1}|)^{\frac{2}{\beta}}).$$

From this, we can get from (32)-(34) that

$$\begin{aligned} \sum_{k=0}^t \mathbf{W}_{k+1}^\top \mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k \mathbf{W}_{k+1} \\ \leq \sum_{i=1}^n \sigma_i^2 \sum_{k=0}^t \lambda_{\max}\{\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k\} + o(\log |\mathbf{P}_{t+1}^{-1}|) + O(1) \\ \leq \sigma_w \log |\mathbf{P}_{t+1}^{-1}| + o(\log |\mathbf{P}_{t+1}^{-1}|) + O(1). \end{aligned}$$

Finally, substituting this into (26), we know that the desired result (20) is true. This completes the proof. \blacksquare

Proof of Theorem 3.1:

By the definitions of $\bar{P}_{t,i}^{-1}$ and $P_{t,i}^{-1}$, it is easy to know that for any $t \geq 0$,

$$P_{t+1,i}^{-1} = \sum_{j=1}^n a_{ji} \bar{P}_{t+1,j}^{-1} = \sum_{j=1}^n a_{ji} (P_{t,j}^{-1} + \varphi_{t,j} \varphi_{t,j}^\top).$$

Consequently, we have

$$\begin{aligned} \max_{1 \leq i \leq n} \lambda_{\max}\{P_{t+1,i}^{-1}\} \\ = \max_{i=1, \dots, n} \lambda_{\max} \left\{ \sum_{j=1}^n a_{ji} (P_{t,j}^{-1} + \varphi_{t,j} \varphi_{t,j}^\top) \right\} \\ \leq \max_{1 \leq i \leq n} \sum_{j=1}^n a_{ji} (\lambda_{\max}\{P_{t,j}^{-1}\} + \lambda_{\max}\{\varphi_{t,j} \varphi_{t,j}^\top\}) \\ \leq \max_{1 \leq i \leq n} \lambda_{\max}\{P_{t,i}^{-1}\} + \sum_{j=1}^n \|\varphi_{t,j}\|^2 \\ \leq \max_{1 \leq i \leq n} \lambda_{\max}\{P_{0,i}^{-1}\} + \sum_{j=1}^n \sum_{k=0}^t \|\varphi_{k,j}\|^2 \\ \leq \lambda_{\max}\{\mathbf{P}_0^{-1}\} + \sum_{j=1}^n \sum_{k=0}^t \|\varphi_{k,j}\|^2. \end{aligned} \quad (35)$$

From (35) and the connection between determinant and eigenvalues of a matrix, it is easy to conclude that

$$\log(|\mathbf{P}_{t+1}^{-1}|) \leq mn \log \left(\max_{1 \leq i \leq n} \lambda_{\max}\{P_{t+1,i}^{-1}\} \right) \leq mn \log(r_t).$$

Consequently, Theorem 3.1 follows from this and Lemma 4.4 immediately.

B. Proof of Theorem 3.2

By the definition of \mathbf{b}_k in (8), we know that

$$\Phi_k \Phi_k^\top = \Phi_k \mathbf{b}_k \Phi_k^\top + \Phi_k (\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k) \Phi_k^\top.$$

Then by noticing that \mathbf{b}_k , Φ_k and \mathbf{P}_k are (block) diagonal matrices, and $\Phi_k^\top \mathbf{P}_k \Phi_k = O(1)$, a.s., we know that

$$\begin{aligned} \sum_{i=1}^n \sum_{k=0}^t R_{k,i} &= \sum_{i=1}^n \sum_{k=0}^t (\varphi_{k,i}^\top \tilde{\theta}_{k,i})^2 = \sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k \Phi_k^\top \tilde{\Theta}_k \\ &= \sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k \mathbf{b}_k \Phi_k^\top \tilde{\Theta}_k + \sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k (\mathbf{b}_k \Phi_k^\top \mathbf{P}_k \Phi_k) \Phi_k^\top \tilde{\Theta}_k \\ &= O \left(\sum_{k=0}^t \tilde{\Theta}_k^\top \Phi_k \mathbf{b}_k \Phi_k^\top \tilde{\Theta}_k \right). \end{aligned} \quad (36)$$

Substituting this into Theorem 3.1 1), we conclude that (12) holds.

C. Proof of Theorem 3.3

For ease of representation, let $a_{ij}^{(s)}$ be the (i, j) -th entry of the matrix \mathcal{A}^s , $s \geq 1$. Note that $a_{ij}^{(1)} = a_{ij}$. By Assumption 3.3 and Remark 3.1, we know that $a_{ji}^{(D_G)} \geq a_{\min} > 0$, where $a_{\min} = \min_{i,j \in \mathcal{V}} a_{ij}^{(D_G)} > 0$, and D_G is the diameter of the graph \mathcal{G} . Consequently, it is not difficult to see that for any $k > D_G$, $a_{ji}^{(k)} \geq a_{\min}$ holds.

By (8), it is easy to see that for any $t \geq 0$,

$$\begin{aligned} \text{vec}\{\mathbf{P}_{t+1}^{-1}\} &= \mathcal{A} \text{vec}\{\bar{\mathbf{P}}_{t+1}^{-1}\} = \mathcal{A} \text{vec}\{\mathbf{P}_t^{-1}\} + \mathcal{A} \text{vec}\{\Phi_t \Phi_t^\top\} \\ &= \mathcal{A}^{t+1} \text{vec}\{\mathbf{P}_0^{-1}\} + \sum_{k=0}^t \mathcal{A}^{t-k+1} \text{vec}\{\Phi_k \Phi_k^\top\}, \end{aligned}$$

which implies that for any $t \geq D_G$,

$$\begin{aligned} P_{t+1,i}^{-1} &= \sum_{j=1}^n a_{ji}^{(t+1)} P_{0,j}^{-1} + \sum_{j=1}^n \sum_{k=0}^t a_{ji}^{(t-k+1)} \varphi_{k,j} \varphi_{k,j}^\top \\ &\geq \sum_{j=1}^n a_{ji}^{(t+1)} P_{0,j}^{-1} + \sum_{j=1}^n \sum_{k=0}^{t-D_G+1} a_{ji}^{(t-k+1)} \varphi_{k,j} \varphi_{k,j}^\top \\ &\geq a_{\min} \sum_{j=1}^n P_{0,j}^{-1} + a_{\min} \sum_{j=1}^n \sum_{k=0}^{t-D_G+1} \varphi_{k,j} \varphi_{k,j}^\top. \end{aligned} \quad (37)$$

Then we have

$$\lambda_{\min}\{\mathbf{P}_{t+1}^{-1}\} \geq a_{\min} \lambda_{\min} \left\{ \sum_{j=1}^n P_{0,j}^{-1} + \sum_{j=1}^n \sum_{k=0}^{t-D_G+1} \varphi_{k,j} \varphi_{k,j}^\top \right\}.$$

Note also that

$$\|\tilde{\Theta}_{t+1}\|^2 \leq \tilde{\Theta}_{t+1}^\top \left[\frac{\mathbf{P}_{t+1}^{-1}}{\lambda_{\min}\{\mathbf{P}_{t+1}^{-1}\}} \right] \tilde{\Theta}_{t+1}.$$

Hence, by 2) in Theorem 3.1, we know that Theorem 3.3 holds.

V. CONCLUDING REMARKS

In this paper, we have established a convergence theory for a basic class of distributed LS algorithms, under quite general conditions on the measured information or data used in the estimation. The accumulated regret of adaptive predictors has been shown to have a celebrated logarithm increase without any excitation condition imposed on the system data, and the convergence rate of the distributed LS estimates has also been established under a cooperative excitation condition, which can be regarded as an extension of the weakest possible excitation condition known for the convergence of the classical LS. Neither independence and stationarity, nor Gaussian property, are required in our results. Moreover, the cooperative excitation condition introduced and used in the paper indicates that the distributed LS can fulfill the estimation task cooperatively, even if any individual sensor cannot due to lack of necessary excitation.

REFERENCES

[1] 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 Processing magazine*, vol. 30, no. 3, pp. 155–171, May 2013.

[2] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, April 2014.

[3] S. Y. Xie and L. Guo, "A necessary and sufficient condition for stability of LMS-based consensus adaptive filters," *Automatica*, vol. 93, pp. 12–19, July 2018.

[4] S. Y. Xie and L. Guo, "Analysis of normalized least mean squares-based consensus adaptive filters under a general information condition," *SIAM Journal on Control and Optimization*, vol. 56, no. 5, pp. 3404–3431, Sept. 2018.

[5] S. Y. Xie and L. Guo, "Analysis of distributed adaptive filters based on diffusion strategies over sensor networks," *IEEE Transactions on Automatic Control*, vol. 63, no. 11, pp. 3643–3658, Nov. 2018.

[6] I. E. K. Harrane, R. Flamary and C. Richard, "On reducing the communication cost of the diffusion LMS algorithm," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 5, no. 1, pp. 100–112, March 2019.

[7] A. H. Sayed and C. G. Lopes, "Distributed recursive least-squares strategies over adaptive networks," in *40th Asilomar Conference on Signals, Systems and Computers*, Pacific, Grove, CA, Oct. 2006, pp. 233–237.

[8] G. Mateos and G. B. Giannakis, "Distributed recursive least-squares: stability and performance analysis," *IEEE Transactions on Signal Processing*, vol. 60, no. 7, pp. 3740–3754, July 2012.

[9] L. Xiao, S. Boyd, and S. Lall, "A space-time diffusion scheme for peer-to-peer least-squares estimation," in *Proceedings of 5th International Conference on Information Processing in Sensor Networks (IPSN 2006)*, Nashville, TN, April 2006, pp. 168–176.

[10] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion recursive least squares for distributed estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 56, no. 5, pp. 1865–1877, May 2008.

[11] A. Bertrand, M. Moonen, and A. H. Sayed, "Diffusion bias-compensated RLS estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 59, no. 11, pp. 5212–5224, Nov. 2011.

[12] R. Arablouei, K. Dogancay, S. Werner, and Y-F Huang, "Adaptive distributed estimation based on recursive least-squares and partial diffusion," *IEEE Transactions on Signal Processing*, vol. 62, no. 14, pp. 3510–3522, July 2014.

[13] V. Vahidpour, A. Rastegarnia, A. Khalili and S. Sanei, "Analysis of partial diffusion recursive least squares adaptation over noisy links," *IET Signal Processing*, vol. 11, no. 6, pp. 749–757, August 2017.

[14] A. Rastegarnia, "Reduced-communication diffusion RLS for distributed estimation over multi-agent networks," *IEEE Transactions on Circuits and Systems II: Express Briefs*, doi: 10.1109/TCSII.2019.2899194.

[15] Y. Yu, H. Zhao, R. C. de Lamare, Y. Zakharov and L. Lu, "Robust distributed diffusion recursive least squares algorithms with side information for adaptive networks," *IEEE Transactions on Signal Processing*, vol. 67, no. 6, pp. 1566–1581, March 2019.

[16] R. Carli, A. Chiuso, L. Schenato, and S. Zampieri, "Distributed Kalman filtering based on consensus strategies," *IEEE Journal on Selected Areas in Communications*, vol. 26, no. 4, pp. 622–633, May 2008.

[17] G. Battistelli, L. Chisci, G. Mugnai, A. Farina, and A. Graziano, "Consensus-based linear and nonlinear filtering," *IEEE Transactions on Automatic Control*, vol. 60, no. 5, pp. 1410–1415, May 2015.

[18] Q. Liu, Z. Wang, X. He, and D. H. Zhou, "On Kalman-consensus filtering with random link failures over sensor networks," *IEEE Transactions on Automatic Control*, vol. 63, no. 8, pp. 2701–2708, Aug. 2018.

[19] S. Das and J. Moura, "Consensus+innovations distributed Kalman filter with optimized gains," *IEEE Transactions on Signal Processing*, vol. 65, no. 2, pp. 467–481, Jan. 2017.

[20] K. J. Aström and B. Wittenmark, "On self tuning regulators," *Automatica*, vol. 9, no. 2, pp. 185–199, March 1973.

[21] L. Guo, "Convergence and logarithm laws of self-tuning regulators," *Automatica*, vol. 31, no. 3, pp. 435–450, March 1995.

[22] L. Ljung, "Consistency of the least-squares identification method," *IEEE Transactions on Automatic Control*, vol. 21, no. 5, pp. 779–781, Oct. 1976.

[23] J. B. Moore, "On strong consistency of least squares identification algorithm," *Automatica*, vol. 14, no. 5, pp. 505–509, Sep. 1978.

[24] H. F. Chen, "Strong consistency and convergence rate of least squares identification," *Science in China Series A - Mathematics, Physics, Astronomy & Technological Science*, vol. 25, no. 7, pp. 771–784, 1982.

[25] T. L. Lai and C. Z. Wei, "Least squares estimates in stochastic regression models with applications to identification and control dynamic systems," *Annals of Statistics*, vol. 10, no. 1, pp. 154–166, March 1982.

[26] T. L. Lai and C. Z. Wei, "Extended least squares and their applications to adaptive control and prediction in linear systems," *IEEE Transactions on Automatic Control*, vol. 31, no. 10, pp. 898–906, Oct. 1986.

[27] H. F. Chen and L. Guo, "Convergence rate of least-squares identification and adaptive control for stochastic systems," *International Journal of Control*, vol. 44, no. 5, pp. 1459–1476, Nov. 1986.

[28] H. F. Chen and L. Guo, *Identification and Stochastic Adaptive Control*, Birkhäuser, Boston, 1991.

[29] L. Guo and H. F. Chen, "The Åström-Wittenmark self-tuning regulator revised and ELS-based adaptive tracker," *IEEE Transactions on Automatic Control*, vol. 36, no. 7, pp. 802–812, July 1991.

[30] L. Ljung, "Analysis of recursive stochastic algorithms," *IEEE Transactions on Automatic Control*, vol. 22, no. 4, pp. 551–575, Aug. 1977.

[31] G. C. Goodwin, P. J. Ramadge, and P. E. Caines, "Discrete time stochastic adaptive control," *SIAM Journal on Control and Optimization*, vol. 19, no. 6, pp. 829–853, 1981.

[32] P. R. Kumar, "Convergence of adaptive control schemes using least-squares parameter estimates," *IEEE Transactions on Automatic Control*, vol. 35, no. 4, pp. 416–424, April 1990.

[33] S. Julier and J. Uhlmann, "A non-divergent estimation algorithm in the presence of unknown correlations," in *Proceedings of the 1997 American Control Conference*, Albuquerque, NM, USA, April 1997, pp. 2369–2373.

[34] L. Chen, P. Arambel, and P. Mehra, "Estimation under unknown correlation: covariance intersection revisited," *IEEE Transactions on Automatic Control*, vol. 47, no. 11, pp. 1879–1882, May 2002.

[35] Y. S. Chow and H. Teicher, *Probability Theory*, New York: Springer, March 1978.

[36] L. Guo, *Time-Varying Stochastic Systems—Stability and Adaptive Theory (Second Edition)*, Science Press, 2020.

[37] M. Zhong and C. G. Cassandras, "Asynchronous distributed optimization with event-driven communication," *IEEE Transactions on Automatic Control*, vol. 55, no. 12, pp. 2735–2750, Dec. 2010.

[38] S. Y. Xie and L. Guo, "Analysis of compressed distributed adaptive filters," *Automatica*, vol. 112, 108707, Feb. 2020.

[39] C. Godsil and G. Royle, *Algebraic Graph Theory*, Springer-Verlag, 2014.

[40] S. Hosseini, A. Chapman, and M. Mesbahi, "Online distributed convex optimization on dynamic networks," *IEEE Transactions on Automatic Control*, vol. 61, no. 11, pp. 3545–3550, Nov. 2016.

[41] S. Shahrampour and A. Jadbabaie, "Distributed online optimization in dynamic environments using mirror descent," *IEEE Transactions on Automatic Control*, vol. 63, no. 3, pp. 714–725, March 2018.

[42] T. L. Lai, "Asymptotically efficient adaptive control in stochastic regression models," *Advances in Applied Mathematics*, vol. 7, no. 1, pp. 23–45, March 1986.

[43] K. Fan, "On a theorem of Weyl concerning eigenvalues of linear transformations," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 36, no. 1, pp. 31–35, Jan. 1950.