# Distributed Sparse Identification for Stochastic Dynamic Systems under Cooperative Non-Persistent Excitation Condition [★]

Die Gan [a,b], Zhixin Liu [a,b]

[a]*Key Laboratory of Systems and Control, Institute of Systems Science, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, P. R. China.*

[b]*School of Mathematical Sciences, University of Chinese Academy of Sciences, Beijing 100049, P. R. China.*

**Abstract**

This paper considers the distributed sparse identification problem over wireless sensor networks such that all sensors cooperatively estimate the unknown sparse parameter vector of stochastic dynamic systems by using the local information from neighbors. A distributed sparse least squares algorithm is proposed by minimizing a local information criterion formulated as a linear combination of accumulative local estimation error and $L_1$-regularization term. The upper bounds of the estimation error and the regret of the adaptive predictor of the proposed algorithm are presented. Furthermore, by designing a suitable adaptive weighting coefficient based on the local observation data, the set convergence of zero elements with a finite number of observations is obtained under a cooperative non-persistent excitation condition. It is shown that the proposed distributed algorithm can work well in a cooperative way even though none of the individual sensors can fulfill the estimation task. Our theoretical results are obtained without relying on the independency assumptions of regression signals that have been commonly used in the existing literature. Thus, our results are expected to be applied to stochastic feedback systems. Finally, the numerical simulations are provided to demonstrate the effectiveness of our theoretical results.

*Key words:* Distributed sparse least squares; Stochastic dynamic system; $L_1$-regularization; Regret; Cooperative non-persistent excitation.

## 1 Introduction

In recent years, wireless sensor networks (WSNs) have attracted increasing research attention because of their wide application in engineering systems including smart grids, biomedical health monitoring, target tracking and surveillance (Sayed et al., 2013; Yick et al., 2008). Distributed observation and data analysis are ubiquitous in WSNs, where sensors are interconnected to acquire and process the local information from neighbors to finish a common task. Due to various uncertainties in practical systems, the distributed identification problem over WSNs becomes one of the important topics where all the sensors collaboratively estimate an unknown parameter vector of interest by using local noisy measurements. Unlike the centralized method with a fusion center, the distributed scheme has the advantages of flexibility, robustness to node or link failures as well as reducing communication load and calculation pressure. Consequently, the theoretical analysis of distributed estimation or filtering algorithms based on several typical distributed strategies such as the incremental, the diffusion and the consensus strategies have been provided (Abdolee and Champagne, 2016; Lou et al., 2017; Battilotti et al., 2020; Liu et al., 2020).

In practical scenarios, there exist a large number of sparse systems (Bazerque and Giannakis, 2010; Vinga, 2021) where many elements in the parameter vector do not contribute or contribute marginally to the systems ( i.e., these elements are zero or near-zero). How to infer the zero elements and identify the nonzero elements

in the unknown parameter vector is an important issue in the investigation of sparse systems. Considerable progress has been made on the identification of zero and nonzero elements in an unknown sparse parameter vector (Zhao and Yu, 2006; Chiuso and Pillonetto, 2014; Eksioglu, 2013), which allows us to obtain a more reliable prediction model. One direction for the estimation of sparse signals is based on the compressed sensing (CS) theory (Candès and Tao, 2005; Baraniuk, 2007), and some estimation algorithms using CS are proposed (cf., Xu et al. (2015); Xie and Guo (2020)) in which *a priori* knowledge about the sparsity of the unknown parameter and the regression vectors are required. Another direction is the sparse optimization based on the regularization framework where the objective function is formulated as a combination of the prediction error with a penalty term. The well-known LASSO (the least absolute shrinkage and selection operator) is one of the classical algorithms to obtain the sparse signals (Tibshirani, 1996), and its variants and adaptive LASSO (Zou, 2006) are also studied. For the stochastic dynamic systems with a single sensor, the adaptive sparse estimation or filtering algorithms are studied by combing the recursive least squares (LS) and least mean squares (LMS) with regularization term (Zhao et al., 2020; Chen et al., 2009).

With the development of sensor networks, some distributed adaptive sparse estimation algorithms have been proposed, and the corresponding stability and convergence analysis are also investigated under some signal conditions. For example, Di Lorenzo and Sayed (2013) provided the convergence and mean-square performance analysis for the distributed LMS algorithm regularized by convex penalties where the assumption of independent regressors is required. Huang and Li (2015) presented theoretical analysis on the mean and mean-square performance of the distributed sparse total LS algorithm under the condition that the input signals are independent and identically distributed (i.i.d.). Shiri et al. (2018) analyzed the mean stability of distributed quasi-sparse affine projection algorithm with independent regression vectors. Huang et al. (2020) analyzed the mean stability of the sparse diffusion LMS algorithm for two regularization terms with independent regression vectors. However, for the typical models such as ARMAX (autoregressive moving-average with exogenous input) model and Hammerstein system, the regressors are often generated by the past input and output signals, so it is hard for them to satisfy the aforementioned independency assumptions.

In order to relax the independency assumption of the regressors, some attempts are made for the distributed adaptive estimation or filtering algorithms. For the unknown time-invariant parameter vector, Gan and Liu (2019) proposed a distributed stochastic gradient algorithm, and established the strong consistency of the proposed algorithm under a cooperative excitation condition. Xie et al. (2021) studied the convergence of the diffusion LS algorithm. For the time-varying parameter vector, Xie and Guo (2018) provided a cooperative information condition to guarantee the stability of the consensus-based LMS adaptive filters. Moreover, Gan et al. (2021) introduced the collective random observability condition and provided the stability analysis of the distributed Kalman filter algorithm. Nevertheless, these asymptotical results are established as the number of the observation data obtained by sensors tends to infinity, which may not be suitable for the sparse identification problem with limited observation data.

Inspired by Zhao et al. (2020) where a sparse identification algorithm for a single sensor case is put forward to infer the set of zero elements with finite observations, we develop a distributed adaptive sparse LS algorithm over sensor networks such that all sensors can cooperatively identify the unknown parameter vector and infer the zero elements with a finite number of observations. The main contributions can be summarized as follows:

- We first introduce a local information criterion for each sensor which is formulated as a linear combination of local estimation errors with $L_1$-regularization term. By minimizing this criterion, a distributed adaptive sparse identification algorithm is proposed. The upper bounds of the estimation error and the accumulative regret of the adaptive predictor are established, which can be degenerated to the results of the classical distributed LS algorithm (Xie et al., 2021) when the weighting coefficients are equal to zero.
- Then, we introduce a cooperative non-persistent excitation condition on the regressors, under which the distributed sparse LS algorithm can cooperatively identify the set of zero elements with finite observations by properly choosing the weighting coefficients. We remark that the key difference between the proposed algorithm and those in distributed sparse optimization framework (e.g., Di Lorenzo and Sayed (2013)) lies in that the weighting coefficients are generated from the local observation sequences. The cooperative excitation condition is much weaker than the widely used persistent excitations (cf., Chen et al. (2014); Zhang et al. (2021); Chen et al. (2015)) and the regularity condition (Zou, 2006).
- Different from most existing results on the distributed sparse algorithms, our theoretical results are obtained without relying on the independency assumptions of regression signals, which makes it possible for applications to the stochastic feedback systems. We also reveal that the whole sensor network can cooperatively accomplish the estimation task, even if any individual sensor can not due to lack of necessary information (Zhao et al., 2020).

The remainder of this paper is organized as follows. In Section 2, we give the problem formulation of this paper; Section 3 presents the main results of the paper includ-

ing the parameter convergence of the algorithm, the regret analysis, and the set convergence of the algorithm; the proofs of the main results are given in Section 4. A simulation example is provided in Section 5. Finally, we conclude the paper with some remarks in Section 6.

## 2 Problem formulation

### 2.1 Basic notations

In this paper, for an $m$-dimensional vector $\boldsymbol{x}$, its $L_p$-norm is defined as $\|\boldsymbol{x}\|_p = (\sum_{j=1}^m |\boldsymbol{x}(j)|^p)^{1/p}$ $(1 \leq p < \infty)$, where $\boldsymbol{x}(j)$ denotes the $j$-th element of $\boldsymbol{x}$. For $p = 1$, $\|\boldsymbol{x}\|_1$ is the sum of absolute values of all the elements in $\boldsymbol{x}$; and for $p = 2$, $\|\boldsymbol{x}\|_2$ is the Euclidean norm, we simply write $\|\cdot\|_2$ as $\|\cdot\|$. For an $m \times m$-dimensional real matrix $\boldsymbol{A}$, we use $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ to denote the largest and smallest eigenvalues of the matrix. $\|\boldsymbol{A}\|$ denotes the Euclidean norm, i.e., $\|\boldsymbol{A}\| = (\lambda_{max}(\boldsymbol{A}\boldsymbol{A}^T))^{\frac{1}{2}}$ where the notation $T$ denotes the transpose operator; $\|\boldsymbol{A}\|_F$ denotes the Frobenius norm, i.e., $\|\boldsymbol{A}\|_F = (tr(\boldsymbol{A}^T\boldsymbol{A}))^{\frac{1}{2}}$, where the notation $tr(\cdot)$ denotes the trace of the corresponding matrix. We use $col(\cdot, \cdots, \cdot)$ to denote a vector stacked by the specified vectors, and $diag(\cdot, \cdots, \cdot)$ to denote a block matrix formed in a diagonal manner of the corresponding vectors or matrices. For a symmetric matrix $\boldsymbol{A}$, if all eigenvalues of $\boldsymbol{A}$ are positive (or nonnegative), then it is a positive definite (semipositive) matrix, and we denote it as $\boldsymbol{A} > 0$ ($\geq 0$). If all elements of a matrix $\boldsymbol{A} = \{a_{ij}\} \in \mathbb{R}^{n \times n}$ are nonnegative, then it is a nonnegative matrix, and furthermore if $\sum_{j=1}^n a_{ij} = 1$ holds for all $i \in \{1, \cdots, n\}$, then it is called a stochastic matrix.

For any two positive scalar sequences $\{a_k\}$ and $\{b_k\}$, by $a_k = O(b_k)$ we mean that there exists a constant $C > 0$ independent of $k$ such that $a_k \leq Cb_k$ holds for all $k \geq 0$, and by $a_k = o(b_k)$ we mean that $\lim_{k \to \infty} a_k / b_k = 0$. For a convex function $f(x)$, we use $\partial f : x \to \partial f(x)$ to denote the subdifferential of $f$, which is a convex set. For example,

$$\partial |x| = \begin{cases} 1, & \text{if } x > 0; \\ -1, & \text{if } x < 0; \\ [\text{-1,1}], & \text{if } x = 0, \end{cases}$$

A necessary and sufficient condition that a given point $x$ belongs to the minimum set of $f$ is $0 \in \partial f(x)$ (see Rockafellar (1972)). We also need to introduce the sign function $sgn(x)$ defined as $sgn(x) = 1$ if $x \geq 0$ and $sgn(x) = -1$ if $x < 0$.

### 2.2 Graph theory

We consider a sensor network with $n$ sensors. The communication between sensors are usually modeled as an undirected weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$, where $\mathcal{V} = \{1, 2, 3, \cdots, n\}$ is the set of sensors (or nodes), $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the edge set, and $\mathcal{A} = \{a_{ij}\} \in \mathbb{R}^{n \times n}$ is the weighted adjacency matrix. The elements of the adjacency matrix $\mathcal{A}$ satisfy $a_{ij} > 0$ if $(i, j) \in \mathcal{E}$ and $a_{ij} = 0$ otherwise. Here we assume that the matrix $\mathcal{A}$ is a symmetric and stochastic matrix. For the sensor $i$, the set of its neighbors is denoted as $N_i = \{j \in \mathcal{V} | (i, j) \in \mathcal{E}\}$, and the sensor $i$ belongs to $N_i$. The sensor $i$ can communicate information with its neighboring sensors. A path of length $\ell$ is a sequence of nodes $\{i_1, ..., i_\ell, i_{\ell+1}\}$ such that $(i_h, i_{h+1}) \in \mathcal{E}$ with $1 \leq h \leq \ell$. The graph $\mathcal{G}$ is called connected if there is a path between any two sensors. The diameter $D_{\mathcal{G}}$ of the graph $\mathcal{G}$ is defined as the maximum shortest path length between any two sensors.

### 2.3 Observation model

In this paper, we consider the parameter identification problem in a network consisting of $n$ sensors labeled $1, \cdots, n$. Assume that the data $\{y_{t,i}, \boldsymbol{\varphi}_{t,i}, t = 1, 2, \cdots\}$ collected by the sensor $i$ obeys the following discrete-time stochastic regression model,

$$y_{t+1,i} = \boldsymbol{\varphi}_{t,i}^T \boldsymbol{\theta} + w_{t+1,i}, \ t = 0, 1, 2, \cdots, \tag{1}$$

where $y_{t,i}$ is the scalar observation or output of the sensor $i$ at time $t$, $\boldsymbol{\varphi}_{t,i}$ is the $m$-dimensional stochastic regression vector which may be the function of current and past inputs and outputs, $\boldsymbol{\theta} \in \mathbb{R}^m$ is an unknown $m$-dimensional parameter to be estimated, and $\{w_{t,i}\}$ is the noise sequence. The above model (1) includes many parameterized systems, such as ARX system and Hammerstein system. We further denote the parameter vector $\boldsymbol{\theta}$ and the index set of its zero elements by

$$\begin{aligned} \boldsymbol{\theta} &\triangleq (\boldsymbol{\theta}(1), \cdots, \boldsymbol{\theta}(m))^T, \\ H^* &\triangleq \{l \in \{1, \cdots, m\} | \boldsymbol{\theta}(l) = 0\}. \end{aligned} \tag{2}$$

Our problem is to design a distributed adaptive estimation algorithm such that all sensors cooperatively infer the set $H^*$ in a finite number of steps and identify the unknown parameter $\boldsymbol{\theta}$ by using stochastic regression vectors and the observation signals from its neighbors, i.e., $\{\boldsymbol{\varphi}_{k,j}, y_{k+1,j}\}_{k=1}^t$ $(j \in N_i)$.

## 3 The main results

### 3.1 Parameter convergence

Before designing the algorithm to cooperatively estimate the unknown parameter vector and infer the set $H^*$, we first introduce the following classical distributed least squares algorithm to estimate the unknown parameter

$\boldsymbol{\theta}$ in (2), i.e.,

$$\boldsymbol{\theta}_{t+1,i} = \boldsymbol{P}_{t+1,i}\left(\sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)}\boldsymbol{\varphi}_{k,j}y_{k+1,j}\right), \qquad (3)$$

where $\boldsymbol{P}_{t+1,i} = \left(\sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)}\boldsymbol{\varphi}_{k,j}\boldsymbol{\varphi}_{k,j}^{T}\right)^{-1}$ and $a_{ij}^{(t+1-k)}$ is the $i$-th row, $j$-th column entry of the matrix $\mathcal{A}^{t+1-k}$. It is clear that the matrix $\boldsymbol{P}_{t+1,i}$ can be equivalently written as the following recursive form,

$$\boldsymbol{P}_{t+1,i}^{-1} = \sum_{j\in N_i} a_{ij}(\boldsymbol{P}_{t,j}^{-1} + \boldsymbol{\varphi}_{t,j}\boldsymbol{\varphi}_{t,j}^{T}). \qquad (4)$$

Thus, the algorithm (3) can also have the following recursive expression,

$$\boldsymbol{\theta}_{t+1,i} = \boldsymbol{P}_{t+1,i}\sum_{j\in N_i} a_{ij}(\boldsymbol{P}_{t,j}^{-1}\boldsymbol{\theta}_{t,j} + \boldsymbol{\varphi}_{t,j}y_{t+1,j}). \qquad (5)$$

Note that in the above derivation, we assume that the matrix $\sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)}\boldsymbol{\varphi}_{k,j}\boldsymbol{\varphi}_{k,j}^{T}$ is invertible which is usually not satisfied for small $t$. To solve this problem, we take the initial matrix $\boldsymbol{P}_{0,i}$ to be positive definite. By (4), we have

$$\boldsymbol{P}_{t+1,i}^{-1} = \sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)}\boldsymbol{\varphi}_{k,j}\boldsymbol{\varphi}_{k,j}^{T} + \sum_{j=1}^{n} a_{ij}^{(t+1)}\boldsymbol{P}_{0,j}^{-1}. \quad (6)$$

This modification will not affect the analysis of the asymptotic properties of the estimate of the distributed least squares algorithm.

In fact, the algorithm (5) can be obtained by minimizing the following linear combination of the estimation error $\sigma_{t+1,i}(\boldsymbol{\beta})$ between the observation signals and the prediction of the local neighbors,

$$\sigma_{t+1,i}(\boldsymbol{\beta}) = \sum_{j\in N_i} a_{ij}\left(\sigma_{t,j}(\boldsymbol{\beta}) + [y_{t+1,j} - \boldsymbol{\beta}^{T}\boldsymbol{\varphi}_{t,j}]^2\right), \quad (7)$$

with $\sigma_{0,i}(\boldsymbol{\beta}) = 0$. That is, $\boldsymbol{\theta}_{t+1,i} \triangleq \arg\min_{\boldsymbol{\beta}} \sigma_{t+1,i}(\boldsymbol{\beta})$.

Set

$$\boldsymbol{e}_{t+1}(\boldsymbol{\beta}) = col\{(y_{t+1,1} - \boldsymbol{\beta}^{T}\boldsymbol{\varphi}_{t,1})^2, ..., (y_{t+1,n} - \boldsymbol{\beta}^{T}\boldsymbol{\varphi}_{t,n})^2\},$$
$$\boldsymbol{\sigma}_t(\boldsymbol{\beta}) = col\{\sigma_{t,1}(\boldsymbol{\beta}), ..., \sigma_{t,n}(\boldsymbol{\beta})\}.$$

Hence by (7), we have

$$\begin{aligned}
\boldsymbol{\sigma}_{t+1}(\boldsymbol{\beta}) &= \mathcal{A}\boldsymbol{\sigma}_t(\boldsymbol{\beta}) + \mathcal{A}\boldsymbol{e}_{t+1}(\boldsymbol{\beta}) \\
&= \mathcal{A}^2\boldsymbol{\sigma}_{t-1}(\boldsymbol{\beta}) + \mathcal{A}^2\boldsymbol{e}_t(\boldsymbol{\beta}) + \mathcal{A}\boldsymbol{e}_{t+1}(\boldsymbol{\beta}) \\
&= \sum_{k=0}^{t} \mathcal{A}^{t+1-k}\boldsymbol{e}_{k+1}(\boldsymbol{\beta}),
\end{aligned}$$

which implies that

$$\sigma_{t+1,i}(\boldsymbol{\beta}) = \sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)}[y_{k+1,j} - \boldsymbol{\beta}^{T}\boldsymbol{\varphi}_{k,j}]^2. \qquad (8)$$

It is shown by Xie et al. (2021) that the distributed least squares algorithm (5) can generate a consistent estimate for the unknown parameter when the number of data tends to infinity. However, for the sparse unknown parameter vectors (i.e., there are many zero elements in $\boldsymbol{\theta}$), it is hard to infer the zero elements in a finite step due to the limitation of observations in practice. In order to solve this issue, we introduce the following local information criterion with $L_1$-regularization to identify the unknown sparse parameters and infer the set $H^*$,

$$J_{t+1,i}(\boldsymbol{\beta}) = \sigma_{t+1,i}(\boldsymbol{\beta}) + \alpha_{t+1,i}\|\boldsymbol{\beta}\|_1, \qquad (9)$$

where $\|\cdot\|_1$ is the $L_1$-norm, $\alpha_{t+1,i}$ is the weighting coefficient chosen to satisfy $\alpha_{t+1,i} = o(\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1}))$, and $\sigma_{t+1,i}(\boldsymbol{\beta})$ is recursively defined by (7). For the sensor $i$, we can obtain the following distributed sparse LS algorithm to estimate the unknown parameter $\boldsymbol{\theta}$ by minimizing $J_{t+1,i}(\boldsymbol{\beta})$, i.e.,

$$\boldsymbol{\beta}_{t+1,i} = \arg\min_{\boldsymbol{\beta}} J_{t+1,i}(\boldsymbol{\beta}). \qquad (10)$$

**Remark 1** *For the sensor $i$, the coefficients $\alpha_{t+1,i}$ in (9) can be dynamically adjusted by using the local observation sequence $\{\boldsymbol{\varphi}_{k,j}, y_{k+1,j}, j \in N_i\}_{k=1}^{t}$, which makes (9) be the adaptive LASSO (cf., Zou (2006)). We show that by properly choosing the coefficient $\alpha_{t+1,i}$, we can identify the set of the zero elements in the unknown sparse parameter vector $\boldsymbol{\theta}$ with a finite number of observations (see Theorem 3).*

In the following, we will first investigate the upper bound of the estimation error generated by (10), which provides the basis for the set convergence of zero elements. For this purpose, we need to introduce the following assumptions on the network topology and the observation noise.

**Assumption 1** *The communication graph $\mathcal{G}$ is connected.*

**Remark 2** *For the weighted adjacency matrix $\mathcal{A}$ of the graph $\mathcal{G}$, we denote $\mathcal{A}^l \triangleq (a_{ij}^{(l)})$ with $l \geq 1$. By the theory of product of stochastic matrices, we see that under Assumption 1, $\mathcal{A}^l$ is a positive matrix for $l \geq D_{\mathcal{G}}$, i.e., for any $i$ and $j$, $a_{ij}^{(l)} > 0$.*

**Assumption 2** *For any $i \in \{1, \cdots, n\}$, the noise sequence $\{w_{k,i}, \mathscr{F}_k\}$ is a martingale difference, and there*

*exists a constant $\delta > 2$ such that*

$$\sup_{k \geq 0} E[|w_{k+1,i}|^\delta | \mathscr{F}_k] < \infty, \quad \text{a.s.},$$

*where $\mathscr{F}_t = \sigma\{\varphi_{k,i}, w_{k,i}, k \leq t, i = 1, \cdots, n\}$ is a sequence of nondecreasing $\sigma$-algebras and $E[\cdot | \cdot]$ denotes the conditional expectation operator.*

We can verify that the i.i.d. zero-mean bounded or Gaussian noise $\{w_{k,i}\}$ which are independent of the regressors can satisfy Assumption 2.

Assume that there are $d$ nonzero elements in the unknown parameter vector $\boldsymbol{\theta}$. Without loss of generality, we assume $\boldsymbol{\theta} = (\boldsymbol{\theta}(1), \cdots, \boldsymbol{\theta}(d), \boldsymbol{\theta}(d+1), \cdots, \boldsymbol{\theta}(m))^T$ with $\boldsymbol{\theta}(l) \neq 0, l = 1, \cdots, d$, and $\boldsymbol{\theta}(j) = 0, j = d+1, \cdots, m$. For the estimate $\boldsymbol{\beta}_{t+1,i}$ obtained by the distributed sparse LS algorithm (10), we denote the estimate error as

$$\widetilde{\boldsymbol{\beta}}_{t+1,i} = \boldsymbol{\beta}_{t+1,i} - \boldsymbol{\theta}. \tag{11}$$

Then we have the following result concerning the upper bound of the estimation error $\widetilde{\boldsymbol{\beta}}_{t,i}$.

**Theorem 1** *Let $\boldsymbol{P}_{t+1,i}^{-1}$ be generated by (4) with arbitrarily initial matrix $\boldsymbol{P}_{0,i} > 0$. Then under Assumptions 1 and 2, we have for all $i \in \{1, \cdots, n\}$*

$$\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\| = O\left(\frac{\alpha_{t+1,i}}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})} + \sqrt{\frac{\log r_t}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})}}\right), \text{a.s.}$$

*where $r_t = \max\limits_{1 \leq i \leq n} \lambda_{\max}\{\boldsymbol{P}_{0,i}^{-1}\} + \sum_{i=1}^{n} \sum_{k=0}^{t} \|\varphi_{k,i}\|^2$.*

The proof of Theorem 1 is provided in Subsection 4.1.

**Remark 3** *By (6), we have for $t \geq D_{\mathcal{G}}$,*

$$\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1}) \geq a_{\min} \lambda_{\min}^{n,t}, \tag{12}$$

*where $a_{\min} \triangleq \min_{i,j \in \mathcal{V}} a_{ij}^{(D_{\mathcal{G}})} > 0$ and*

$$\lambda_{\min}^{n,t} = \lambda_{\min}\left\{\sum_{j=1}^{n} \boldsymbol{P}_{0,j}^{-1} + \sum_{j=1}^{n} \sum_{k=0}^{t-D_{\mathcal{G}}+1} \varphi_{k,j} \varphi_{k,j}^T\right\}.$$

*From Theorem 1, if the coefficient $\alpha_{t+1,i}$ is chosen to satisfy $\alpha_{t+1,i} = o(\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1}))$ and the regression vectors satisfy the weakest possible cooperative excitation condition $\log r_t = o(\lambda_{\min}^{n,t})$ (cf., Xie et al. (2021)), then the almost sure convergence of the distributed sparse LS algorithm can be obtained, i.e., $\boldsymbol{\beta}_{t+1,i} \xrightarrow[t \to \infty]{} \boldsymbol{\theta}$.*

## 3.2 Analysis of the regret

Regret is one of the key metrics for evaluating the performance of the online learning algorithms (Hosseini et al., 2016; Shahrampour and Jadbabaie, 2018). For each sensor $i \in \{1, \cdots, n\}$, we construct an adaptive predictor $\hat{y}_{t+1,i}$ by using the estimate $\boldsymbol{\beta}_{t,i}$ defined in (10) at the time instant $t$,

$$\hat{y}_{t+1,i} = \varphi_{t,i}^T \boldsymbol{\beta}_{t,i}.$$

The prediction error can be described by the following loss function $\rho_{t+1,i}(\boldsymbol{\beta}_{t,i})$, i.e.,

$$\rho_{t+1,i}(\boldsymbol{\beta}_{t,i}) = E\left[(y_{t+1,i} - \hat{y}_{t+1,i})^2 | \mathscr{F}_k\right]$$
$$= E\left[(y_{t+1,i} - \varphi_{t,i}^T \boldsymbol{\beta}_{t,i})^2 | \mathscr{F}_t\right].$$

Then the cumulative regret over the whole network is defined as

$$R_t = \sum_{i=1}^{n} \sum_{k=0}^{t} \rho_{k+1,i}(\boldsymbol{\beta}_{k,i}) - \min_{\boldsymbol{\zeta} \in \mathbb{R}^m} \sum_{i=1}^{n} \sum_{k=0}^{t} \rho_{k+1,i}(\boldsymbol{\zeta}).$$

The regret defined above reflects the difference between the cumulative loss $\rho_{k+1,i}(\boldsymbol{\beta}_{k,i})$ when the unknown parameter is estimated by (10) and the optimal static value of the cumulative loss function $\rho_{k+1,i}(\cdot)$. Due to existence of the noise, it is generally desired that the average regret $R_t/nt$ is small or even goes to zero as $t \to \infty$.

In the following, we analyze the asymptotic property of the regret $R_t$ over the sensor network. By Assumption 2 and the fact $\varphi_{k,i}^T \widetilde{\boldsymbol{\beta}}_{k,i} \in \mathscr{F}_k$, we have

$$R_t = \sum_{i=1}^{n} \sum_{k=0}^{t} E((y_{k+1,i} - \hat{y}_{k+1,i})^2 | \mathscr{F}_k)$$
$$- \min_{\boldsymbol{\zeta} \in \mathbb{R}^m} \sum_{i=1}^{n} \sum_{k=0}^{t} E((y_{k+1,i} - \varphi_{k+1,i}^T \boldsymbol{\zeta})^2 | \mathscr{F}_k)$$
$$= \sum_{i=1}^{n} \sum_{k=0}^{t} E(\varphi_{k,i}^T \widetilde{\boldsymbol{\beta}}_{k,i} + w_{k+1,i})^2 | \mathscr{F}_k)$$
$$- \min_{\boldsymbol{\zeta} \in \mathbb{R}^m} \sum_{i=1}^{n} \sum_{k=0}^{t} E((\varphi_{k,i}^T(\boldsymbol{\theta} - \boldsymbol{\zeta}) + w_{k+1,i})^2 | \mathscr{F}_k)$$
$$= \sum_{i=1}^{n} \sum_{k=0}^{t} (\varphi_{k,i}^T \widetilde{\boldsymbol{\beta}}_{k,i})^2. \tag{13}$$

**Theorem 2** *Under Assumption 2, if $\boldsymbol{\Phi}_t^T \boldsymbol{P}_t \boldsymbol{\Phi}_t = O(1)$, and $\alpha_{t,i} = O\left(\sqrt{\lambda_{\min}(\boldsymbol{P}_{t,i}^{-1})}\right)$, we have*

$$R_t = O(\log r_t), \quad \text{a.s.}$$

*where $\boldsymbol{\Phi}_t \triangleq diag\{\varphi_{t,1}, ...\varphi_{t,n}\}, \boldsymbol{P}_t \triangleq diag\{\boldsymbol{P}_{t,1}, ..., \boldsymbol{P}_{t,n}\}$, and $r_t$ is defined in Theorem 1.*

The proof of Theorem 2 is given in Subsection 4.2.

**Remark 4** *We know that for the bounded regressors $\varphi_{t,i}$, $r_t$ will be of the order $O(t)$. Consequently, by Theorem 2, the upper bound of the regret $R_t$ over the sensor network is sublinear with respect to $nt$, i.e., $R_t/nt = O(\log t/t) \to 0$ as $t \to \infty$. The analysis of the regret does not require any excitation condition on the regression signals. Theorem 1 and Theorem 2 can be degenerated to the results of the classical distributed LS algorithm in Xie et al. (2021) when $\alpha_{t+1,i}$ is equal to zero.*

### 3.3 Set convergence

In the last two subsections, we have obtained the asymptotic results concerning the parameter convergence and the regret analysis. Inspired by Zhao et al. (2020), we propose the following distributed sparse adaptive algorithm (Algorithm 1) to identify the set of zero elements with a finite number of observations by choosing $\alpha_{t,i}$ adaptively.

**Algorithm 1.**
**Step 1:** Based on $\{\varphi_{k,j}, y_{k+1,j}\}_{k=1}^{t}$ ($j \in N_i$), begin with an initial vector $\boldsymbol{\theta}_{0,i}$ and an initial matrix $\boldsymbol{P}_{0,i} > 0$, compute the matrix $\boldsymbol{P}_{t+1,i}^{-1}$ defined by (4) and the local estimate $\boldsymbol{\theta}_{t+1,i}$ of $\boldsymbol{\theta}$ by (5), and further define

$$\hat{\boldsymbol{\theta}}_{t+1,i}(l)$$
$$\triangleq \boldsymbol{\theta}_{t+1,i}(l) + \text{sgn}(\boldsymbol{\theta}_{t+1,i}(l))\sqrt{\frac{\log(\lambda_{\max}(\boldsymbol{P}_{t+1,i}^{-1}))}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})}}, \quad (14)$$

**Step 2:** Choose a positive sequence $\{\alpha_{k,i}\}_{k=1}^{t+1}$ satisfying

$$\alpha_{k,i} = o(\lambda_{\min}(\boldsymbol{P}_{k,i}^{-1})),$$
$$\lambda_{\max}(\boldsymbol{P}_{k,i}^{-1})\sqrt{\frac{\log(\lambda_{\max}(\boldsymbol{P}_{k,i}^{-1}))}{\lambda_{\min}(\boldsymbol{P}_{k,i}^{-1})}} = o(\alpha_{k,i}). \quad (15)$$

**Step 3:** Optimize the convex objective local function,

$$\bar{J}_{t+1,i}(\boldsymbol{\xi}) = \sigma_{t+1,i}(\boldsymbol{\xi}) + \alpha_{t+1,i}\sum_{l=1}^{m}\frac{1}{|\hat{\boldsymbol{\theta}}_{t+1,i}(l)|}|\boldsymbol{\xi}(l)| \quad (16)$$

with $\sigma_{t+1,i}(\boldsymbol{\xi})$ defined in (7), and obtain

$$\boldsymbol{\xi}_{t+1,i} = (\boldsymbol{\xi}_{t+1,i}(1), \cdots, \boldsymbol{\xi}_{t+1,i}(m))^T$$
$$\triangleq \arg\min_{\boldsymbol{\xi}} \bar{J}_{t+1,i}(\boldsymbol{\xi}), \quad (17)$$
$$H_{t+1,i} \triangleq \{l = 1, \cdots, m | \boldsymbol{\xi}_{t+1,i}(l) = 0\}. \quad (18)$$

In the convex objective function (16), different components in $\boldsymbol{\xi}$ are assigned different weights, which is an adaptive LASSO estimator since the weights $\alpha_{t+1,i}/\hat{\boldsymbol{\theta}}_{t+1,i}(l)$ are generated from the local observation sequence $\{\varphi_{k,j}, y_{k+1,j}, j \in N_i\}_{k=1}^{t}$. The $\hat{\boldsymbol{\theta}}_{t+1,i}(l)$ appearing in the denominator satisfies that $|\hat{\boldsymbol{\theta}}_{t+1,i}(l)| \geq \sqrt{\frac{\log(\lambda_{\max}(\boldsymbol{P}_{t+1,i}^{-1}))}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})}} > 0$, which makes (16) well defined. Moreover, if $\hat{\boldsymbol{\theta}}_{t+1,i}(l) \to 0$ for some $l = 1, \cdots, m$ and hence $1/\hat{\boldsymbol{\theta}}_{t+1,i}(l) \to \infty$, then the corresponding minimizer $\boldsymbol{\xi}_{t+1,i}(l)$ should be exactly zero. This provides an intuitive explanation for the sparse solution of Algorithm 1 with a finite number of observations. The set $H_{t+1,i}$ generated from the convex optimization problem (17) serves as the estimate for the set $H^*$ defined in (2). There exist some typical algorithms such as basic pursuit and interior-point algorithms to solve the convex optimization problem (17) in the literature (see e.g., Kim et al. (2007); Gill et al. (2011)).

We introduce the following cooperative non-persistent excitation condition to study the convergence of the sets of zero elements in the unknown sparse parameter vector with a finite number of observations, which is different from the asymptotic analysis given in the last two subsections.

**Assumption 3** *(Cooperative Non-Persistent Excitation Condition) The following condition is satisfied,*

$$\frac{r_t}{\lambda_{\min}^{n,t}}\sqrt{\frac{\log(r_t)}{\lambda_{\min}^{n,t}}} \xrightarrow[t\to\infty]{} 0, \quad \text{a.s.} \quad (19)$$

*where $r_t$ and $\lambda_{\min}^{n,t}$ are respectively defined in Theorem 1 and Remark 3.*

**Remark 5** *For the single sensor case with $n = 1$ and $D_{\mathcal{G}} = 1$, the condition (19) reduces to the excitation condition given by Zhao et al. (2020). Assumption 3 reveals the cooperative effect of multiple sensors in the sense that the condition (19) can make it possible for Algorithm 1 to estimate the unknown parameter $\boldsymbol{\theta}$ and the sets of zero elements by the cooperation of multiple sensors even if any individual sensor cannot due to lack of adequate excitation, which is also shown in the simulation example given in Section 5.*

For the set $H_{t,i}$ obtained by (18), we get the following finite time convergence result, which shows that the set of zero elements in $\boldsymbol{\theta}$ can be correctly identified with a finite number of observations.

**Theorem 3** *(Set convergence) Under Assumptions 1-3, if $\log r_t = O(\log r_{t-D_{\mathcal{G}}+1})$, then there exists a positive integer $T_0$ (which may depend on the sample $\omega$) such that for all $i \in \{1, \cdots, n\}$*

$$\boldsymbol{\xi}_{t+1,i}(d+1) = \cdots = \boldsymbol{\xi}_{t+1,i}(m) = 0, \quad t \geq T_0.$$

That is, $H_{t+1,i} = H^*$ for $t \geq T_0$, where $H^*$ and $H_{t+1,i}$ are defined in (2) and (18).

The detailed proof of Theorem 3 is given in Subsection 4.3.

**Remark 6** *From Theorem 3 (also Theorem 1 and Theorem 2 ), we see that the parameter convergence, regret analysis, and set convergence results in this paper are derived without using the independency assumption on the regression vectors, which makes it possible to apply our algorithm to practical feedback systems.*

## 4   Proofs of the main results

In order to prove the main theorems of the paper, we first give two preliminary lemmas.

Denote the estimation error of the classical distributed LS algorithm (5) as $\widetilde{\boldsymbol{\theta}}_{t+1,i} \triangleq \boldsymbol{\theta}_{t+1,i} - \boldsymbol{\theta}$, and $\widetilde{\boldsymbol{\Theta}}_t = col\{\widetilde{\boldsymbol{\theta}}_{t,1}, ..., \widetilde{\boldsymbol{\theta}}_{t,n}\}$.

**Lemma 1** (Xie et al., 2021) *Under Assumptions 1 and 2, we have the following results for the classical distributed LS algorithm (5),*

1) $\sum_{i=1}^{n} \|\widetilde{\boldsymbol{\theta}}_{t,i}\|^2 = O\left(\dfrac{\log r_t}{\lambda_{\min}^{n,t}}\right),$

2) $\sum_{k=0}^{t} \lambda_{\max}(\boldsymbol{d}_k \boldsymbol{\Phi}_k^T \boldsymbol{P}_k \boldsymbol{\Phi}_k) = O(\log r_t),$

3) $\sum_{k=0}^{t} \widetilde{\boldsymbol{\Theta}}_k^T \boldsymbol{\Phi}_k \boldsymbol{d}_k \boldsymbol{\Phi}_k^T \widetilde{\boldsymbol{\Theta}}_k = O(\log r_t),$

*where $\boldsymbol{P}_k$ and $\boldsymbol{\Phi}_k$ are defined in Theorem 2, $r_t \triangleq \max\limits_{1 \leq i \leq n} \lambda_{\max}\{\boldsymbol{P}_{0,i}^{-1}\} + \sum_{i=1}^{n} \sum_{k=0}^{t} \|\boldsymbol{\varphi}_{k,i}\|^2$ and $\boldsymbol{d}_t \triangleq diag\left\{\dfrac{1}{1+\boldsymbol{\varphi}_{t,1}^T \boldsymbol{P}_{t,1} \boldsymbol{\varphi}_{t,1}}, ..., \dfrac{1}{1+\boldsymbol{\varphi}_{t,n}^T \boldsymbol{P}_{t,n} \boldsymbol{\varphi}_{t,n}}\right\}.$*

The following lemma provides an upper bound for the cumulative summation of the noises.

**Lemma 2** (Gan and Liu, 2022) *Under Assumptions 1 and 2, for any $i \in \{1, ..., n\}$, we have*

$$\left\| \boldsymbol{P}_{t,i}^{\frac{1}{2}} \left( \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} w_{k+1,j} \right) \right\| = O(\sqrt{\log(r_t)}).$$

### 4.1   Proof of Theorem 1

**Proof.** By noting that $\boldsymbol{\beta}_{t+1,i}$ is the minimizer of

$J_{t+1,i}(\boldsymbol{\beta})$, it follows that

$$\begin{aligned} 0 &\geq J_{t+1,i}(\boldsymbol{\beta}_{t+1,i}) - J_{t+1,i}(\boldsymbol{\theta}) \\ &= J_{t+1,i}(\widetilde{\boldsymbol{\beta}}_{t+1,i} + \boldsymbol{\theta}) - J_{t+1,i}(\boldsymbol{\theta}). \end{aligned} \tag{20}$$

Since $\boldsymbol{\theta}(j) = 0$, $j = d+1, \cdots, m$, by (1), (8) and (9), we have

$$\begin{aligned} &J_{t+1,i}(\widetilde{\boldsymbol{\beta}}_{t+1,i} + \boldsymbol{\theta}) \\ &= \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} [w_{k+1,j} - \widetilde{\boldsymbol{\beta}}_{t+1,i}^T \boldsymbol{\varphi}_{k,j}]^2 \\ &\quad + \alpha_{t+1,i} \sum_{l=1}^{d} |\widetilde{\boldsymbol{\beta}}_{t+1,i}(l) + \boldsymbol{\theta}(l)| + \alpha_{t+1,i} \sum_{l=d+1}^{m} |\widetilde{\boldsymbol{\beta}}_{t+1,i}(l)| \\ &= \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} w_{k+1,j}^2 \\ &\quad - 2\widetilde{\boldsymbol{\beta}}_{t+1,i}^T \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} w_{k+1,j} \\ &\quad + \widetilde{\boldsymbol{\beta}}_{t+1,i}^T \sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} \boldsymbol{\varphi}_{k,j}^T \widetilde{\boldsymbol{\beta}}_{t+1,i} \\ &\quad + \alpha_{t+1,i} \sum_{l=1}^{d} |\widetilde{\boldsymbol{\beta}}_{t+1,i}(l) + \boldsymbol{\theta}(l)| + \alpha_{t+1,i} \sum_{l=d+1}^{m} |\widetilde{\boldsymbol{\beta}}_{t+1,i}(l)|. \end{aligned} \tag{21}$$

Similarly, we have

$$\begin{aligned} &J_{t+1,i}(\boldsymbol{\theta}) \\ &= \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} [y_{k+1,j} - \boldsymbol{\theta}^T \boldsymbol{\varphi}_{k,j}]^2 + \alpha_{t+1,i} \sum_{l=1}^{d} |\boldsymbol{\theta}(l)| \\ &= \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} w_{k+1,j}^2 + \alpha_{t+1,i} \sum_{l=1}^{d} |\boldsymbol{\theta}(l)|. \end{aligned} \tag{22}$$

Hence by (21) and (22), we have

$$\begin{aligned} &J_{t+1,i}(\widetilde{\boldsymbol{\beta}}_{t+1,i} + \boldsymbol{\theta}) - J_{t+1,i}(\boldsymbol{\theta}) \\ &\geq \widetilde{\boldsymbol{\beta}}_{t+1,i}^T \sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} \boldsymbol{\varphi}_{k,j}^T \widetilde{\boldsymbol{\beta}}_{t+1,i} \\ &\quad - 2\widetilde{\boldsymbol{\beta}}_{t+1,i}^T \sum_{j=1}^{n} \sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} w_{k+1,j} \\ &\quad + \alpha_{t+1,i} \sum_{l=1}^{d} (|\widetilde{\boldsymbol{\beta}}_{t+1,i}(l) + \boldsymbol{\theta}(l)| - |\boldsymbol{\theta}(l)|) \\ &\triangleq M_{t+1,i}^{(1)} - 2M_{t+1,i}^{(2)} + M_{t+1,i}^{(3)}. \end{aligned} \tag{23}$$

In the following, we estimate $M_{t+1,i}^{(1)}$, $M_{t+1,i}^{(2)}$ and $M_{t+1,i}^{(3)}$ separately. Denote $\boldsymbol{V}_{t+1,i} = \boldsymbol{P}_{t+1,i}^{-\frac{1}{2}} \widetilde{\boldsymbol{\beta}}_{t+1,i}$. By Lemma 2,

we have

$$|M^{(2)}_{t+1,i}|$$

$$= \left| \widetilde{\boldsymbol{\beta}}^T_{t+1,i} \boldsymbol{P}^{-\frac{1}{2}}_{t+1,i} \boldsymbol{P}^{\frac{1}{2}}_{t+1,i} \sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j} w_{k+1,j} \right|$$

$$= O\left(\sqrt{\log(r_t)}\right) \|\boldsymbol{V}_{t+1,i}\|.$$

Hence, there exists a positive constant $c_1$ such that for large $t$,

$$M^{(1)}_{t+1,i} - 2M^{(2)}_{t+1,i}$$
$$\geq \frac{1}{2}\|\boldsymbol{V}_{t+1,i}\|^2 - c_1\sqrt{\log(r_t)}\|\boldsymbol{V}_{t+1,i}\|. \tag{24}$$

By $C_r$-inequality, we have

$$|M^{(3)}_{t+1,i}| \leq \alpha_{t+1,i} \sum_{l=1}^d |\widetilde{\boldsymbol{\beta}}_{t+1,i}(l)| \leq \alpha_{t+1,i}\sqrt{d}\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\|. \tag{25}$$

Hence by (20) and (23)-(25), we have for large $t$

$$0 \geq \frac{\|\boldsymbol{V}_{t+1,i}\|^2}{2} - c_1\sqrt{\log(r_t)}\|\boldsymbol{V}_{t+1,i}\| - \sqrt{d}\alpha_{t+1,i}\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\|,$$

which implies that

$$\|\boldsymbol{V}_{t+1,i}\| \leq \sqrt{c_1^2 \log r_t + 2\sqrt{d}\alpha_{t+1,i}\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\|} + \sqrt{c_1 \log r_t}. \tag{26}$$

Note that by the definition of $\boldsymbol{V}_{t+1,i}$, we have

$$\|\boldsymbol{V}_{t+1,i}\|^2 \geq \lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\|^2.$$

Combining this with (26), we have

$$\left( \|\widetilde{\boldsymbol{\beta}}_{t+1,i}\| - \frac{2\sqrt{d}\alpha_{t+1,i}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})} \right)^2$$

$$\leq \left( \frac{2\sqrt{d}\alpha_{t+1,i}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})} \right)^2 + \frac{(2c_1^2 + 2c_1)\log r_t}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})}.$$

Thus, we have

$$\|\widetilde{\boldsymbol{\beta}}_{t+1,i}\| = O\left( \frac{\alpha_{t+1,i}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})} + \sqrt{\frac{\log r_t}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t+1,i})}} \right), \tag{27}$$

which completes the proof of the theorem. □

### 4.2 Proof of Theorem 2

**Proof.** By (8), we obtain the subdifferential of (9),

$$\partial J_{t+1,i}(\boldsymbol{\beta}) = -2\sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(y_{k+1,j} - \boldsymbol{\varphi}_{k,j}^T\boldsymbol{\beta})$$
$$+ \alpha_{t+1,i}\partial\|\boldsymbol{\beta}\|_1,$$

where $\partial\|\boldsymbol{\beta}\|_1$ is the subdifferential of $\|\boldsymbol{\beta}\|_1$. Since $\boldsymbol{\beta}_{t+1,i}$ is the minimizer of $J_{t+1,i}(\boldsymbol{\beta})$, we have $\boldsymbol{0} \in \partial J_{t+1,i}(\boldsymbol{\beta}_{t+1,i})$ with $\boldsymbol{0} \triangleq (\underbrace{0,0,\cdots,0}_m)^T$, i.e.,

$$\boldsymbol{0} \in -2\sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(y_{k+1,j} - \boldsymbol{\varphi}_{k,j}^T\boldsymbol{\beta}_{t+1,i})$$
$$+ \alpha_{t+1,i}\partial\|\boldsymbol{\beta}_{t+1,i}\|_1. \tag{28}$$

Let us write (28) in a component form, i.e., for all $l \in \{1,\cdots,m\}$,

$$0 \in -\sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(l)y_{k+1,j}$$
$$+ \sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(l)\left( \sum_{s\neq l} \boldsymbol{\varphi}_{k,j}(s)\boldsymbol{\beta}_{t+1,i}(s) \right)$$
$$+ \sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}^2(l)\boldsymbol{\beta}_{t+1,i}(l) + \frac{\alpha_{t+1,i}}{2}\partial|\boldsymbol{\beta}_{t+1,i}(l)|$$
$$\triangleq \boldsymbol{D}_{t+1,i}(l) + \sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}^2(l)\boldsymbol{\beta}_{t+1,i}(l)$$
$$+ \frac{\alpha_{t+1,i}}{2}\partial|\boldsymbol{\beta}_{t+1,i}(l)|. \tag{29}$$

Note that

$$\partial|\boldsymbol{\beta}_{t+1,i}(l)| = \begin{cases} 1, & \text{if } \boldsymbol{\beta}_{t+1,i}(l) > 0 \\ -1, & \text{if } \boldsymbol{\beta}_{t+1,i}(l) < 0 \\ \in [-1,1], & \text{if } \boldsymbol{\beta}_{t+1,i}(l) = 0 \end{cases}.$$

Set $\sum_{j=1}^n \sum_{k=0}^t a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}\boldsymbol{\varphi}_{k,j}^T \triangleq \boldsymbol{\Psi}_{t+1,i}$. Combining the above equation with (29) yields for large $t$

$$\boldsymbol{\beta}_{t+1,i}(l)$$
$$= \begin{cases} \dfrac{-\boldsymbol{D}_{t+1,i}(l) + \frac{\alpha_{t+1,i}}{2}}{\boldsymbol{\Psi}_{t+1,i}(l,l)}, & \text{if } \boldsymbol{D}_{t+1,i}(l) > \dfrac{\alpha_{t+1,i}}{2} \\[3ex] \dfrac{-\boldsymbol{D}_{t+1,i}(l) - \frac{\alpha_{t+1,i}}{2}}{\boldsymbol{\Psi}_{t+1,i}(l,l)}, & \text{if } \boldsymbol{D}_{t+1,i}(l) < -\dfrac{\alpha_{t+1,i}}{2} \\[3ex] 0, & \text{if } |\boldsymbol{D}_{t+1,i}(l)| \leq \dfrac{\alpha_{t+1,i}}{2} \end{cases},$$

with $\boldsymbol{\Psi}_{t+1,i}(l,l)$ being the $l$-th diagonal element of the matrix $\boldsymbol{\Psi}_{t+1,i}$. This implies that

$$\boldsymbol{\Psi}_{t+1,i}(l,l)\boldsymbol{\beta}_{t+1,i}(l) = -\boldsymbol{D}_{t+1,i}(l) + \boldsymbol{\gamma}_{t+1,i}(l), \qquad (30)$$

where

$$\boldsymbol{\gamma}_{t+1,i}(l) = \begin{cases} \dfrac{\alpha_{t+1,i}}{2}, & \text{if } \boldsymbol{D}_{t+1,i}(l) > \dfrac{\alpha_{t+1,i}}{2} \\[2mm] -\dfrac{\alpha_{t+1,i}}{2}, & \text{if } \boldsymbol{D}_{t+1,i}(l) < -\dfrac{\alpha_{t+1,i}}{2}. \\[2mm] \boldsymbol{D}_{t+1,i}(l), & \text{if } |\boldsymbol{D}_{t+1,i}(l)| \leq \dfrac{\alpha_{t+1,i}}{2} \end{cases}$$

Then by (30) and the definition of $\boldsymbol{D}_{t+1,i}(l)$, we have for all $l \in \{1, \cdots, m\}$

$$\sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(l)\boldsymbol{\varphi}_{k,j}^{T}\boldsymbol{\beta}_{t+1,i}$$
$$= \sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}(l)y_{k+1,j} + \boldsymbol{\gamma}_{t+1,i}(l).$$

We rewrite the above equation into the matrix form, and obtain the following equation by (3) for large $t$

$$\boldsymbol{\beta}_{t+1,i} = \boldsymbol{P}_{t+1,i} \left( \sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}y_{k+1,j} + \boldsymbol{\gamma}_{t+1,i} \right)$$
$$= \boldsymbol{\theta}_{t+1,i} + \boldsymbol{P}_{t+1,i}\boldsymbol{\gamma}_{t+1,i}, \qquad (31)$$

where $\boldsymbol{\gamma}_{t+1,i} = (\boldsymbol{\gamma}_{t+1,i}(1), \cdots, \boldsymbol{\gamma}_{t+1,i}(m))^{T}$ and $\boldsymbol{\theta}_{t+1,i}$ is defined in (3). Note that for all $i \in \{1, \cdots, n\}$ and $l \in \{1, \cdots, m\}$, $\|\boldsymbol{\gamma}_{t+1,i}(l)\| \leq \frac{\alpha_{t+1,i}}{2}$, hence by Lemma 1, we obtain

$$\sum_{k=0}^{t} \boldsymbol{\gamma}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k}\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\gamma}_{k}$$
$$\leq \sum_{k=0}^{t} \lambda_{\max}(\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k})\boldsymbol{\gamma}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\gamma}_{k}$$
$$\leq \sum_{k=0}^{t} \left[ \lambda_{\max}(\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k}) \left( \sum_{i=1}^{n} \lambda_{\max}(\boldsymbol{P}_{k,i})\|\boldsymbol{\gamma}_{k,i}\|^{2} \right) \right]$$
$$= O \left( \sum_{k=0}^{t} \left[ \lambda_{\max}(\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k}) \left( \sum_{i=1}^{n} \frac{\alpha_{k,i}^{2}}{\lambda_{\min}(\boldsymbol{P}_{k,i}^{-1})} \right) \right] \right)$$
$$= O(\log r_{t}), \qquad (32)$$

where $\boldsymbol{\gamma}_{t+1} = col\{\boldsymbol{\gamma}_{t+1,1}, \cdots, \boldsymbol{\gamma}_{t+1,n}\}$. By the definition of $\boldsymbol{d}_{t}$ in Lemma 1, we have $\boldsymbol{I}_{n} = \boldsymbol{d}_{k} + \boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k}$. By (31), we have $\widetilde{\boldsymbol{\beta}}_{t+1} = \widetilde{\boldsymbol{\Theta}}_{t+1} + \boldsymbol{P}_{t+1}\boldsymbol{\gamma}_{t+1}$, where $\widetilde{\boldsymbol{\beta}}_{t} = col\{\widetilde{\boldsymbol{\beta}}_{t,1}, ..., \widetilde{\boldsymbol{\beta}}_{t,n}\}$. Hence by (32), Lemma 1 and the con-

dition $\boldsymbol{\Phi}_{t}^{T}\boldsymbol{P}_{t}\boldsymbol{\Phi}_{t} = O(1)$, we have

$$R_{t} = \sum_{k=0}^{t} \widetilde{\boldsymbol{\beta}}_{k}^{T}\boldsymbol{\Phi}_{k}\boldsymbol{\Phi}_{k}^{T}\widetilde{\boldsymbol{\beta}}_{k}$$
$$= \sum_{k=0}^{t} \widetilde{\boldsymbol{\beta}}_{k}^{T}\boldsymbol{\Phi}_{k}\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\widetilde{\boldsymbol{\beta}}_{k} + \sum_{k=0}^{t} \widetilde{\boldsymbol{\beta}}_{k}^{T}\boldsymbol{\Phi}_{k}(\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k})\boldsymbol{\Phi}_{k}^{T}\widetilde{\boldsymbol{\beta}}_{k}$$
$$= O \left( \sum_{k=0}^{t} \widetilde{\boldsymbol{\beta}}_{k}^{T}\boldsymbol{\Phi}_{k}\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\widetilde{\boldsymbol{\beta}}_{k} \right)$$
$$= O \left( \sum_{k=0}^{t} \widetilde{\boldsymbol{\Theta}}_{k}^{T}\boldsymbol{\Phi}_{k}\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\widetilde{\boldsymbol{\Theta}}_{k} + \sum_{k=0}^{t} \boldsymbol{\gamma}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\Phi}_{k}\boldsymbol{d}_{k}\boldsymbol{\Phi}_{k}^{T}\boldsymbol{P}_{k}\boldsymbol{\gamma}_{k} \right)$$
$$= O(\log r_{t}).$$

This completes the proof of the theorem. $\qquad\square$

### 4.3 Proof of Theorem 3

**Proof.** Denote the estimation error between $\boldsymbol{\xi}_{t+1,i}$ obtained by Algorithm 1 and $\boldsymbol{\theta}$ as

$$\widetilde{\boldsymbol{\xi}}_{t+1,i} = \boldsymbol{\xi}_{t+1,i} - \boldsymbol{\theta}. \qquad (33)$$

By Assumption 3 and Lemma 1, we see that the limits of $\boldsymbol{\theta}_{t+1,i}(l)$ and $\hat{\boldsymbol{\theta}}_{t+1,i}(l)$, $l = 1, \cdots, d$ are nonzero. Similar to the proof of Theorem 1, we also have the following result,

$$\|\widetilde{\boldsymbol{\xi}}_{t+1,i}\| = O \left( \frac{\alpha_{t+1,i}}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})} + \sqrt{\frac{\log r_{t}}{\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1})}} \right). \quad (34)$$

By the definition of $\widetilde{\boldsymbol{\xi}}_{t+1,i}$ in (33), it suffices to prove that there exists a positive integer $T_{0}$ such that for all $i \in \{1, \cdots, n\}$

$$\widetilde{\boldsymbol{\xi}}_{t+1,i}(d+1) = \cdots = \widetilde{\boldsymbol{\xi}}_{t+1,i}(m) = 0, \quad t \geq T_{0}.$$

Otherwise, if for some $s_{l} \in \{d+1, \cdots, m\}$, some sensor $i_{0}$, and some subsequence $\{t_{p}\}_{p\geq1}$ such that $\widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}(s_{l}) \neq 0$, $p \geq 1$. Thus for $p \geq 1$, we have $\|\widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}\| > 0$.

Denote

$$\widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}} = \begin{pmatrix} \widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}^{(1)} \\ \widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}^{(2)} \end{pmatrix} \text{ and } \bar{\boldsymbol{\xi}}_{t_{p}+1,i_{0}} = \begin{pmatrix} \widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}^{(1)} \\ \boldsymbol{0} \end{pmatrix}, \qquad (35)$$

where $\widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}^{(1)} \in \mathbb{R}^{d}$ and $\widetilde{\boldsymbol{\xi}}_{t_{p}+1,i_{0}}^{(2)} \in \mathbb{R}^{m-d}$. By noting that $\boldsymbol{\xi}_{t_{p}+1,i_{0}}$ is the minimizer of $\bar{J}_{t_{p}+1,i_{0}}(\boldsymbol{\xi})$ defined by

9

(16) , it follows that

$$0 \geq \bar{J}_{t_p+1,i_0}(\boldsymbol{\xi}_{t_p+1,i_0}) - \bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \bar{\boldsymbol{\xi}}_{t_p+1,i_0})$$
$$= \bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}) - \bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \bar{\boldsymbol{\xi}}_{t_p+1,i_0}). \quad (36)$$

Denote

$$\boldsymbol{\Psi}_{t+1,i} = \sum_{j=1}^{n}\sum_{k=0}^{t} a_{ij}^{(t+1-k)} \boldsymbol{\varphi}_{k,j}\boldsymbol{\varphi}_{k,j}^{T} \triangleq \begin{pmatrix} \boldsymbol{\Psi}_{t+1,i}^{(11)} & \boldsymbol{\Psi}_{t+1,i}^{(12)} \\ \boldsymbol{\Psi}_{t+1,i}^{(21)} & \boldsymbol{\Psi}_{t+1,i}^{(22)} \end{pmatrix},$$
$$\text{and } \boldsymbol{\varphi}_{k,j} \triangleq \begin{pmatrix} \boldsymbol{\varphi}_{k,j}^{(1)} \\ \boldsymbol{\varphi}_{k,j}^{(2)} \end{pmatrix}. \quad (37)$$

Similar to (21), we have for $\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}$

$$\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}) - \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} w_{k+1,j}^2$$

$$= -2\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} \boldsymbol{\varphi}_{k,j}^{(1)} w_{k+1,j}$$
$$-2\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} \boldsymbol{\varphi}_{k,j}^{(2)} w_{k+1,j}$$
$$+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(11)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(21)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)}$$
$$+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(12)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(22)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}$$
$$+\alpha_{t_p+1,i_0} \sum_{l=1}^{d} \frac{1}{\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)} |\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}(l) + \boldsymbol{\theta}(l)|$$
$$+\alpha_{t_p+1,i_0} \sum_{l=d+1}^{m} \frac{1}{|\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)|} |\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}(l)|. \quad (38)$$

For $\bar{\boldsymbol{\xi}}_{t_p+1,i_0}$ defined in (35), we have

$$\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \bar{\boldsymbol{\xi}}_{t_p+1,i_0}) - \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} w_{k+1,j}^2$$

$$= -2\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} \boldsymbol{\varphi}_{k,j}^{(1)} w_{k+1,j}$$
$$+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(11)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)}$$
$$+\alpha_{t_p+1,i_0} \sum_{l=1}^{d} \frac{1}{|\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)|} |\bar{\boldsymbol{\xi}}_{t_p+1,i_0}(l) + \boldsymbol{\theta}(l)|. \quad (39)$$

By (38) and (39), we have

$$\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}) - \bar{J}_{t_p+1,i_0}(\boldsymbol{\theta} + \bar{\boldsymbol{\xi}}_{t_p+1,i_0})$$

$$= -2\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} \boldsymbol{\varphi}_{k,j}^{(2)} w_{k+1,j}$$
$$+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(22)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)} + \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(12)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}$$
$$+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} \boldsymbol{\Psi}_{t_p+1,i_0}^{(21)} \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(1)}$$
$$+\alpha_{t_p+1,i_0} \sum_{l=d+1}^{m} \frac{1}{|\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)|} |\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}(l)|$$
$$\triangleq -2I_{t_p+1,i_0}^{(1)} + I_{t_p+1,i_0}^{(2)} + I_{t_p+1,i_0}^{(3)} + I_{t_p+1,i_0}^{(4)} + I_{t_p+1,i_0}^{(5)}. \quad (40)$$

In the following, we estimate $I_{t_p+1,i_0}^{(1)}$, $I_{t_p+1,i_0}^{(2)}$, $I_{t_p+1,i_0}^{(3)}$, $I_{t_p+1,i_0}^{(4)}$, $I_{t_p+1,i_0}^{(5)}$ separately. By (6) and (37), we have

$$\boldsymbol{P}_{t+1,i}^{-1} = \boldsymbol{\Psi}_{t+1,i} + \sum_{j=1}^{n} a_{ij}^{(t+1)} \boldsymbol{P}_{0,j}^{-1} \triangleq \begin{pmatrix} \boldsymbol{Q}_{t+1,i}^{(11)} & \boldsymbol{Q}_{t+1,i}^{(12)} \\ \boldsymbol{Q}_{t+1,i}^{(21)} & \boldsymbol{Q}_{t+1,i}^{(22)} \end{pmatrix}.$$

By (37) and Lemma 2, we have

$$|I_{t_p+1,i_0}^{(1)}| = \left| \widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)T} (\boldsymbol{Q}_{t_p+1,i_0}^{(22)})^{\frac{1}{2}} (\boldsymbol{Q}_{t_p+1,i_0}^{(22)})^{-\frac{1}{2}} \right.$$
$$\left. \sum_{j=1}^{n}\sum_{k=0}^{t_p} a_{i_0j}^{(t_p+1-k)} \boldsymbol{\varphi}_{k,j}^{(2)} w_{k+1,j} \right|$$
$$= \|(\boldsymbol{Q}_{t_p+1,i_0}^{(22)})\|^{\frac{1}{2}} \|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\| O\left(\sqrt{\log r_{t_p}^{(2)}}\right),$$

where $r_t^{(2)} \triangleq \max_{1\leq i\leq n} \lambda_{\max}\{\boldsymbol{Q}_{0,i}^{(22)}\} + \sum_{i=1}^{n}\sum_{k=0}^{t} \|\boldsymbol{\varphi}_{k,i}^{(2)}\|^2$.

Note that $\lambda_{\max}(\boldsymbol{Q}_{t_p+1,i_0}^{(22)}) \leq \lambda_{\max}(\boldsymbol{P}_{t_p+1,i_0}^{-1})$ and $\lambda_{\min}(\boldsymbol{Q}_{t_p+1,i_0}^{(22)}) \geq \lambda_{\min}(\boldsymbol{P}_{t_p+1,i_0}^{-1})$. Hence, we have $r_{t_p}^{(2)} \leq r_{t_p}$. We obtain that for large $p$ and some positive constant $c_2$

$$-2I_{t_p+1,i_0}^{(1)} + I_{t_p+1,i_0}^{(2)}$$
$$\geq \lambda_{\min}(\boldsymbol{\Psi}_{t_p+1,i_0}^{(22)})\|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\|^2$$
$$-c_2\|(\boldsymbol{Q}_{t_p+1,i_0}^{(22)})\|^{\frac{1}{2}} \|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\| \sqrt{\log r_{t_p}^{(2)}}$$
$$\geq \frac{1}{2}\lambda_{\min}(\boldsymbol{Q}_{t_p+1,i_0}^{(22)})\|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\|^2$$
$$-c_2\|(\boldsymbol{Q}_{t_p+1,i_0}^{(22)})\|^{\frac{1}{2}} \|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\| \sqrt{\log r_{t_p}^{(2)}}$$
$$\geq \frac{1}{2}\lambda_{\min}(\boldsymbol{P}_{t_p+1,i_0}^{-1})\|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\|^2$$
$$-c_2\sqrt{\lambda_{\max}(\boldsymbol{P}_{t_p+1,i_0}^{-1})}\|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}^{(2)}\| \sqrt{\log r_{t_p}}. \quad (41)$$

By (34) and Lemma 1, and based on the equivalence of norms in a finite dimensional space, we have

$$
\begin{aligned}
|I^{(3)}_{t_p+1,i_0}| &= |\widetilde{\boldsymbol{\xi}}^{(1)T}_{t_p+1,i_0}\boldsymbol{\Psi}^{(12)}_{t_p+1,i_0}\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}| \\
&\leq \|\widetilde{\boldsymbol{\xi}}^{(1)}_{t_p+1,i_0}\|\|\boldsymbol{\Psi}^{(12)}_{t_p+1,i_0}\|\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\| \\
&\leq c_3\|\widetilde{\boldsymbol{\xi}}^{(1)}_{t_p+1,i_0}\|\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\|\boldsymbol{\Psi}^{(12)}_{t_p+1,i_0}\|_F \\
&\leq c_3\|\widetilde{\boldsymbol{\xi}}^{(1)}_{t_p+1,i_0}\|\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\|\boldsymbol{\Psi}_{t_p+1,i_0}\|_F \\
&\leq c_4\|\widetilde{\boldsymbol{\xi}}^{(1)}_{t_p+1,i_0}\|\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\|\boldsymbol{\Psi}_{t_p+1,i_0}\| \\
&\leq c_4\|\widetilde{\boldsymbol{\xi}}^{(1)}_{t_p+1,i_0}\|\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0}) \\
&= O\left(\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})\left[\frac{\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\right.\right. \\
&\left.\left. +\sqrt{\frac{\log(r_{t_p})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}\right]\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\right),
\end{aligned} \tag{42}
$$

where $c_3$ and $c_4$ are two positive constants.

Similarly, we have

$$
\begin{aligned}
|I^{(4)}_{t_p+1,i_0}| \leq O\left(\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})\left[\frac{\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\right.\right. \\
\left.\left.+\sqrt{\frac{\log(r_{t_p})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}\right]\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\right).
\end{aligned} \tag{43}
$$

Then by the definition of $\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)$ in (14), and the condition $\log r_t = O(\log r_{t-D_{\mathcal{G}}+1})$, we have for $l = d+1,\cdots,m$,

$$
L_{t_p+1,i_0} \leq |\hat{\boldsymbol{\theta}}_{t_p+1,i_0}(l)| \leq c_5 L_{t_p+1,i_0},
$$

where $c_5 > 0$ is a positive constant, and

$$
L_{t_p+1,i_0} = \sqrt{\frac{\log(\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0}))}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}.
$$

Hence we have

$$
\begin{aligned}
I^{(5)}_{t_p+1,i_0} &\geq \alpha_{t_p+1,i_0}\frac{1}{c_5 L_{t_p+1,i_0}}\sum_{l=d+1}^{m}|\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}(l)| \\
&\geq \alpha_{t_p+1,i_0}\frac{1}{c_5 L_{t_p+1,i_0}}\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|.
\end{aligned} \tag{44}
$$

Thus, by (40)-(44), for some $c_6 > 0$, we obtain

$$
\begin{aligned}
&\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta}+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0})-\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta}+\bar{\boldsymbol{\xi}}_{t_p+1,i_0}) \\
&\geq \lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\cdot \\
&\left(\left(\frac{\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|}{2}-c_2\sqrt{\frac{\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}\sqrt{\frac{\log r_{t_p}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}\right.\right. \\
&-\frac{c_6\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\left[\frac{\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}+\right. \\
&\left.\left.\left.\sqrt{\frac{\log(r_{t_p})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}}\right]+\frac{\alpha_{t_p+1,i_0}}{c_5\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})L_{t_p+1,i_0}}\right)\right).
\end{aligned} \tag{45}
$$

By (12), (15) and Assumption 3, we have

$$
\begin{aligned}
&\frac{\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\sqrt{\frac{\log r_{t_p}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}} \\
&\leq \frac{\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\sqrt{\frac{\log r_{t_p}}{\lambda^{n,t_p}_{\min}}} \\
&= o\left(\frac{\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})L_{t_p+1,i_0}}\right).
\end{aligned} \tag{46}
$$

By (12) and Assumption 3, we have

$$
\begin{aligned}
&\frac{\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})\alpha_{t_p+1,i_0}}{\lambda^2_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}\Big/\frac{\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})L_{t_p+1,i_0}} \\
&= L_{t_p+1,i_0}\frac{\lambda_{\max}(\boldsymbol{P}^{-1}_{t_p+1,i_0})}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})} \\
&= O\left(\frac{r_{t_p}}{\lambda^{n,t_p}_{\min}}\sqrt{\frac{\log(r_{t_p})}{\lambda^{n,t_p}_{\min}}}\right) = o(1).
\end{aligned} \tag{47}
$$

From (45)-(47), we have

$$
\begin{aligned}
&\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta}+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0})-\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta}+\bar{\boldsymbol{\xi}}_{t_p+1,i_0}) \\
&\geq \lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|\cdot \\
&\left(\frac{\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\|}{2}+\frac{[\frac{1}{c_5}+o(1)]\alpha_{t_p+1,i_0}}{\lambda_{\min}(\boldsymbol{P}^{-1}_{t_p+1,i_0})L_{t_p+1,i_0}}\right).
\end{aligned} \tag{48}
$$

Note that $\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0}(s_l) \neq 0$ for some $s_l \in \{d+1,\cdots,m\}$. Hence $\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t_p+1,i_0}\| > 0$. Then by (48), we have $J_{t_p+1,i_0}(\boldsymbol{\theta}+\widetilde{\boldsymbol{\xi}}_{t_p+1,i_0})-\bar{J}_{t_p+1,i_0}(\boldsymbol{\theta}+\bar{\boldsymbol{\xi}}_{t_p+1,i_0}) > 0$, which contradicts (36). This implies that $\|\widetilde{\boldsymbol{\xi}}^{(2)}_{t+1,i}\| = 0$ for all large $t$ and all $i \in \{1,\cdots,n\}$. We complete the proof of the theorem. $\square$

## 5 A simulation example

In this section, we provide an example to illustrate the performance of the distributed sparse identification algorithm (i.e., Algorithm 1) proposed in this paper.

**Example 1** *Consider a network composed of $n = 6$ sensors whose dynamics obey the model (1) with the dimension $m = 5$. The noise sequence $\{w_{t,i}, t \geq 1, i = 1, \cdots, n\}$ in (1) is independent and identically distributed with $w_{t,i} \sim \mathcal{N}(0, 0.1)$ (Gaussian distribution with zero mean and variance $0.1$). Let the regression vectors $\boldsymbol{\varphi}_{t,i} \in \mathbb{R}^m$ $(i = 1, \cdots, n, \ t \geq 1)$ be generated by the following state space model,*

$$\begin{aligned} \boldsymbol{x}_{t,i} &= \boldsymbol{A}_i \boldsymbol{x}_{t-1,i} + \boldsymbol{B}_i \varepsilon_{t,i}, \\ \boldsymbol{\varphi}_{t,i} &= \boldsymbol{C}_i \boldsymbol{x}_{t,i}, \end{aligned} \tag{49}$$

*where $\boldsymbol{x}_{t,i} \in \mathbb{R}^m$ is the state of the above system with $\boldsymbol{x}_{0,i} = [\underbrace{1, \cdots, 1}_{m}]^T$, the matrices $\boldsymbol{A}_i$, $\boldsymbol{B}_i$ and $\boldsymbol{C}_i$ $(i = 1, 2, \cdots, n)$ are chosen according to the following way such that the regression vector $\boldsymbol{\varphi}_{t,i}$ is lack of adequate excitation for any individual sensor,*

$$\boldsymbol{A}_i = diag\{\underbrace{1.1, \cdots, 1.1}_{m}\},$$

$$\boldsymbol{B}_i = \boldsymbol{e_j} \in \mathbb{R}^m,$$
$$\boldsymbol{C}_i = col\{\underbrace{0, \cdots, 0, \boldsymbol{e_j}, 0, \cdots, 0}_{j^{th}}\} \in \mathbb{R}^{m \times m},$$

*where $j = \mod (i, m)$ and $\boldsymbol{e_j}$ $(j = 1, \cdots, m)$ is the $j$th column of the identity matrix $\boldsymbol{I}_m$ $(m = 5)$. Let the noise sequence $\{\varepsilon_{t,i}, t \geq 1, i = 1, \cdots, n\}$ in (49) be independent and identically distributed with $\varepsilon_{t,i} \sim \mathcal{N}(0, 0.2)$. All sensors will estimate an unknown parameter*

$$\boldsymbol{\theta} = [\boldsymbol{\theta}(1), \boldsymbol{\theta}(2), \boldsymbol{\theta}(3), \boldsymbol{\theta}(4), \boldsymbol{\theta}(5)]^T = [0.8, 1.6, 0, 0, 0]^T.$$

*The initial estimate is taken as $\boldsymbol{\xi}_{0,i} = [1, 1, 1, 1, 1]^T$ for $i = 1, 2, \cdots, 6$. We use the Metropolis rule (Xiao et al., 2005) to construct the weights of the network, i.e.,*

$$a_{li} = \begin{cases} 1 - \sum_{j \neq i} a_{ij} & \text{if } l = i \\ 1/(\max\{n_i, n_l\}) & \text{if } l \in N_i \setminus \{i\} \end{cases} \tag{50}$$

*where $n_i$ is the degree of the node $i$.*

It can be verified that for each sensor $i$ $(i = 1 \cdots, 6)$, the regression signals $\boldsymbol{\varphi}_{t,i}$ ( generated by (49)) have no adequate excitation to estimate the unknown parameter, but they can cooperate to satisfy Assumption 3. We repeat the simulation for $s = 100$ times with the same initial states.

1) We estimate the unknown parameter $\boldsymbol{\theta}$ by using the non-cooperative sparse identification algorithm (i.e., the adjacency matrix is the unit matrix) and the distributed sparse identification algorithm (Algorithm 1) proposed in this paper respectively. We adopt the Matlab CVX tools (http://cvxr.com/cvx/) to solve the convex optimization problem (16), and take the weight coefficient as $\alpha_{t,i} = (\lambda_{\min}(\boldsymbol{P}_{t+1,i}^{-1}))^{0.75}$. The average estimation error generated by these two algorithms is shown in Fig. 1. We see that the estimation error generated by distributed sparse identification algorithm converges to zero as $t$ increases, while the estimation error of the non-cooperative sparse identification algorithm does not. The estimate sequences $\{\boldsymbol{\xi}_{t,i}(1), \boldsymbol{\xi}_{t,i}(2), \boldsymbol{\xi}_{t,i}(3), \boldsymbol{\xi}_{t,i}(4), \boldsymbol{\xi}_{t,i}(5)\}_{t=0}^{200}$ $(i = 1, \cdots, 6)$ generated by Algorithm 1 are given in Fig. 2. We see from these figures that the estimates can converge to the true value $\boldsymbol{\theta}$. Therefore, the estimation task can be fulfilled through exchanging information between sensors even though any individual sensor can not.



Fig. 1. The estimation errors of the distributed sparse identification algorithm and non-cooperative sparse identification algorithm

2) We estimate the unknown parameter $\boldsymbol{\theta}$ by using the classical distributed LS algorithm studied by Xie et al. (2021) and Algorithm 1 proposed in this paper under the same network topology. Table 1 and Table 2 show the estimates for $\boldsymbol{\theta}(3)$, $\boldsymbol{\theta}(4)$, $\boldsymbol{\theta}(5)$ by these two algorithms at different time instants $t$. From Table 1 and Table 2, we can see that, compared with the distributed LS algorithm in Xie et al. (2021), Algorithm 1 can generate sparser and more accurate estimates for the unknown parameters and thus give us valuable information in inferring the zero and nonzero elements in the unknown parameters.

## 6 Concluding remarks

In this paper, we first introduced a local information criterion which is formulated as a linear combination of

Estimates for $\theta(1)$ of all sensors  Estimates for $\theta(2)$ of all sensors  Estimates for $\theta(3)$ of all sensors  Estimates for $\theta(4)$ of all sensors  Estimates for $\theta(5)$ of all sensors

sensor 1, sensor 2, sensor 3, sensor 4, sensor 5, sensor 6

Fig. 2. The estimate sequences $\{\boldsymbol{\xi}_{t,i}\}_{t=0}^{200}$ of all sensors

Table 1
Estimates by the distributed LS algorithm in Xie et al. (2021) and Algorithm 1 for $t = 50$

|  | sensor 1 | sensor 2 | sensor 3 | sensor 4 | sensor 5 | sensor 6 |
|---|---|---|---|---|---|---|
| Estimate for $\boldsymbol{\theta}(3)$ |  |  |  |  |  |  |
| By distributed LS | $2.5892 \times 10^{-4}$ | $1.4805 \times 10^{-4}$ | $3.1352 \times 10^{-4}$ | $2.7231 \times 10^{-4}$ | $2.9085 \times 10^{-4}$ | $2.9085 \times 10^{-4}$ |
| By Algorithm 1 | $-2.8518 \times 10^{-6}$ | $-6.3009 \times 10^{-12}$ | $-8.3539 \times 10^{-18}$ | $1.4030 \times 10^{-6}$ | $-4.6969 \times 10^{-7}$ | $-2.7547 \times 10^{-18}$ |
| Estimate for $\boldsymbol{\theta}(4)$ |  |  |  |  |  |  |
| By distributed LS | $2.7949 \times 10^{-4}$ | $2.7949 \times 10^{-4}$ | $2.7949 \times 10^{-4}$ | $0.0011$ | $2.7949 \times 10^{-4}$ | $2.7949 \times 10^{-4}$ |
| By Algorithm 1 | $7.2376 \times 10^{-18}$ | $1.6087 \times 10^{-8}$ | $-6.2511 \times 10^{-5}$ | $1.1212 \times 10^{-6}$ | $-3,6619 \times 10^{-10}$ | $-7.8179 \times 10^{-7}$ |
| Estimate for $\boldsymbol{\theta}(5)$ |  |  |  |  |  |  |
| By distributed LS | $2.1450 \times 10^{-4}$ | $8.1487 \times 10^{-5}$ | $1.7771 \times 10^{-4}$ | $1.7014 \times 10^{-4}$ | $4.9508 \times 10^{-5}$ | $1.7350 \times 10^{-4}$ |
| By Algorithm 1 | $-2.8248 \times 10^{-6}$ | $1.3601 \times 10^{-10}$ | $-3.8278 \times 10^{-5}$ | $7.7698 \times 10^{-18}$ | $-1.3398 \times 10^{-8}$ | $-2.6207 \times 10^{-6}$ |

Table 2
Estimates by the distributed LS algorithm in Xie et al. (2021) and Algorithm 1 for $t = 100$

|  | sensor 1 | sensor 2 | sensor 3 | sensor 4 | sensor 5 | sensor 6 |
|---|---|---|---|---|---|---|
| Estimate for $\boldsymbol{\theta}(3)$ |  |  |  |  |  |  |
| By distributed LS | $3.9929 \times 10^{-6}$ | $4.0720 \times 10^{-6}$ | $4.5125 \times 10^{-6}$ | $3.9929 \times 10^{-6}$ | $3.9929 \times 10^{-6}$ | $4.5015 \times 10^{-6}$ |
| By Algorithm 1 | $-4.1586 \times 10^{-13}$ | $2.6792 \times 10^{-12}$ | $1.7980 \times 10^{-13}$ | $-1.6160 \times 10^{-12}$ | $8.8066 \times 10^{-14}$ | $6.9114 \times 10^{-13}$ |
| Estimate for $\boldsymbol{\theta}(4)$ |  |  |  |  |  |  |
| By distributed LS | $6.4080 \times 10^{-6}$ | $3.6820 \times 10^{-6}$ | $3.2300 \times 10^{-6}$ | $2.5931 \times 10^{-6}$ | $6.4080 \times 10^{-6}$ | $3.2300 \times 10^{-6}$ |
| By Algorithm 1 | $-5.9666 \times 10^{-12}$ | $-4.0833 \times 10^{-19}$ | $-8.79535 \times 10^{-12}$ | $-1.2865 \times 10^{-11}$ | $-7.3473 \times 10^{-12}$ | $-6.2931 \times 10^{-12}$ |
| Estimate for $\boldsymbol{\theta}(5)$ |  |  |  |  |  |  |
| By distributed LS | $4.5652 \times 10^{-6}$ | $4.8154 \times 10^{-6}$ | $4.6507 \times 10^{-6}$ | $5.9311 \times 10^{-6}$ | $5.5863 \times 10^{-6}$ | $4.6507 \times 10^{-6}$ |
| By Algorithm 1 | $1.4196 \times 10^{-12}$ | $1.9062 \times 10^{-12}$ | $-1.8454 \times 10^{-12}$ | $3.0918 \times 10^{-12}$ | $-2.1729 \times 10^{-15}$ | $-1.9412 \times 10^{-12}$ |

the local estimation error with $L_1$-regularization term. By minimizing this criterion, we proposed a distributed sparse identification algorithm to estimate an unknown parameter vector of a stochastic system. The upper bounds of the estimation error and the averaged accumulated regrets of adaptive prediction are obtained without excitation conditions. Furthermore, we showed that under the cooperative non-persistent excitation conditions, the set of zero elements in the unknown parameter vector can be correctly identified with a finite number of observations by properly choosing the weighting coefficient. We remark that our theoretical results are established without using such stringent conditions as independency of the regression vectors, which makes it possible to combine the distributed adaptive estimation with the distributed control. For future research, it will be interesting to consider the combination of the distributed sparse identification algorithm with the distributed control, and design a recursive distributed sparse adaptive algorithm.

# References

Abdolee, R. and Champagne, B. (2016). Diffusion LMS strategies in sensor networks with noisy input data. *IEEE/ACM Transactions on Networking*, 24(1):3–14.

Baraniuk, R. G. (2007). Compressive sensing. *IEEE Signal Processing Magazine*, 24(4):118–121.

Battilotti, S., Cacace, F., d'Angelo, M., and Germani, A. (2020). Asymptotically optimal consensus-based distributed filtering of continuous-time linear systems. *Automatica*, 122:109189.

Bazerque, J. A. and Giannakis, G. B. (2010). Distributed spectrum sensing for cognitive radio networks by exploiting sparsity. *IEEE Transactions on Signal Processing*, 58(3):1847–1862.

Candès, E. J. and Tao, T. (2005). Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215.

Chen, W., Hua, S., and Zhang, H. (2015). Consensus-based distributed cooperative learning from closed-loop neural control systems. *IEEE Transactions on Neural Networks and Learning Systems*, 26(2):331–345.

Chen, W., Wen, C., Hua, S., and Sun, C. (2014). Distributed cooperative adaptive identification and control for a group of continuous-time systems with a cooperative PE condition via consensus. *IEEE Transactions on Automatic Control*, 59(1):91–106.

Chen, Y., Gu, Y., and Hero, A. O. (2009). Sparse LMS for system identification. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3125–3128.

Chiuso, A. and Pillonetto, G. (2014). Bayesian and nonparametric methods for system identification and model selection. In *2014 European Control Conference*, pages 2376–2381.

Di Lorenzo, P. and Sayed, A. H. (2013). Sparse distributed learning based on diffusion adaptation. *IEEE Transactions on Signal Processing*, 61(6):1419–1433.

Eksioglu, E. M. (2013). Group sparse RLS algorithms. *International Journal of Adaptive Control and Signal Processing*, 28(12):1398–1412.

Gan, D. and Liu, Z. (2019). Strong consistency of the distributed stochastic gradient algorithm. In *Proceedings of the 58th IEEE Conference on Decision and Control*, pages 5082–5087, Nice, France.

Gan, D. and Liu, Z. (2022). Distributed order estimation of ARX model under cooperative excitation condition. *SIAM Journal on Control and Optimization*, arXiv:2110.09826.

Gan, D., Xie, S., and Liu, Z. (2021). Stability of the distributed Kalman filter using general random coefficients. *Science China Information Sciences*, 64:172204.

Gill, P. R., Wang, A., and Molnar, A. (2011). The in-crowd algorithm for fast basis pursuit denoising. *IEEE Transactions on Signal Processing*, 59(10):4595–4605.

Hosseini, S., Chapman, A., and Mesbahi, M. (2016). Online distributed convex optimization on dynamic networks. *IEEE Transactions on Automatic Control*, 61(11):3545–3550.

Huang, S. and Li, C. (2015). Distributed sparse total least-squares over networks. *IEEE Transactions on Signal Processing*, 63(11):2986–2998.

Huang, W., Chen, C., Yao, X., and Li, Q. (2020). Diffusion fused sparse LMS algorithm over networks. *Signal Processing*, 171:107497.

Kim, S.-J., Koh, K., Lustig, M., Boyd, S., and Gorinevsky, D. (2007). An interior-point method for large-scale $\ell_1$-regularized least squares. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):606–617.

Liu, Y., Liu, J., Xu, C., Li, G., and He, Y. (2020). Fully distributed variational Bayesian non-linear filter with unknown measurement noise in sensor networks. *Science China Information Sciences*, 63:210202.

Lou, J., Jia, L., Tao, R., and Wang, Y. (2017). Distributed incremental bias-compensated RLS estimation over multi-agent networks. *Science China Information Sciences*, 60:032204.

Rockafellar, R. T. (1972). *Convex analysis.* Princeton University Press.

Sayed, A. H., Tu, S.-Y., Chen, J., Zhao, X., and Towfic, Z. J. (2013). Diffusion strategies for adaptation and learning over networks: an examination of distributed strategies and network behavior. *IEEE Signal Processing Magazine*, 30(3):155–171.

Shahrampour, S. and Jadbabaie, A. (2018). Distributed online optimization in dynamic environments using mirror descent. *IEEE Transactions on Automatic Control*, 63(3):714–725.

Shiri, H., Tinati, M. A., Codreanu, M., and Daneshvar, S. (2018). Distributed sparse diffusion estimation based on set membership and affine projection algorithm. *Digital Signal Processing*, 73:47–61.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological)*, 58(1):267–288.

Vinga, S. (2021). Structured sparsity regularization for analyzing high-dimensional omics data. *Brief Bioinform*, 22(1):77–87.

Xiao, L., Boyd, S., and Lall, S. (2005). A scheme for robust distributed sensor fusion based on average consensus. In *Proceedings of the 4th Fourth International Symposium on Information Processing in Sensor Networks*, pages 63–70, Boise, ID, USA.

Xie, S. and Guo, L. (2018). A necessary and sufficient condition for stability of LMS-based consensus adaptive filters. *Automatica*, 93:12–19.

Xie, S. and Guo, L. (2020). Analysis of compressed distributed adaptive filters. *Automatica*, 112:108707.

Xie, S., Zhang, Y., and Guo, L. (2021). Convergence of a distributed least squares. *IEEE Transactions on Automatic Control*, 66(10):4952–4959.

Xu, S., de Lamare, R. C., and Poor, H. V. (2015). Distributed compressed estimation based on compressive sensing. *IEEE Signal Processing Letters*, 22(9):1311–1315.

Yick, J., Mukherjee, B., and Ghosal, D. (2008). Wire-

less sensor network survey. *Computer Networks*, 52(12):2292–2330.

Zhang, H., Wang, T., and Zhao, Y. (2021). Asymptotically efficient recursive identification of fir systems with binary-valued observations. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 51(5):2687–2700.

Zhao, P. and Yu, B. (2006). On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2563.

Zhao, W., Yin, G., and Bai, E.-W. (2020). Sparse system identification for stochastic systems with general observation sequences. *Automatica*, 121:109162.

Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429.