



**HAL**  
open science

# Online Proximal Learning over Jointly Sparse Multitask Networks With $L_{\infty,1}$ Regularization

Danqi Jin, Jie Chen, Cédric Richard, Jingdong Chen

► **To cite this version:**

Danqi Jin, Jie Chen, Cédric Richard, Jingdong Chen. Online Proximal Learning over Jointly Sparse Multitask Networks With  $L_{\infty,1}$  Regularization. IEEE Transactions on Signal Processing, 2020, 68, pp.6319-6335. 10.1109/TSP.2020.3021247 . hal-03347269

**HAL Id: hal-03347269**

**<https://hal.science/hal-03347269v1>**

Submitted on 17 Sep 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Online Proximal Learning Over Jointly Sparse Multitask Networks With $\ell_{\infty,1}$ Regularization

Danqi Jin, *Student Member, IEEE*, Jie Chen, *Senior Member, IEEE*, Cédric Richard, *Senior Member, IEEE*, Jingdong Chen, *Senior Member, IEEE*

## Abstract

Modeling relations between local optimum parameter vectors to estimate in multitask networks has attracted much attention over the last years. This work considers a distributed optimization problem with jointly sparse structure among nodes, that is, the local solutions have the same sparse support set. Several mixed norm have been proposed to address the jointly sparse structure in the literature. Among several candidates, the (reweighted)  $\ell_{\infty,1}$ -norm is element-wise separable, it is more convenient to evaluate their approximate proximal operators. Thus by introducing a (reweighted)  $\ell_{\infty,1}$ -norm penalty term at each node, and using a proximal gradient method to minimize the regularized cost, we devise a proximal multitask diffusion LMS algorithm which can promote joint-sparsity. Analyses are provided to characterize the algorithm behavior in the mean and mean-square sense. Simulation results are presented to show its effectiveness, as well as the accuracy of the theoretical findings.

## Index Terms

Distributed optimization, diffusion strategy, joint sparsity, proximal algorithm, stochastic performance, (reweighted)  $\ell_{\infty,1}$ -norm.

## I. INTRODUCTION

Because of their superior performance and wider stability range [1], diffusion strategies have been widely used in multi-agent networks to address estimation problems in a distributed and online manner. Several diffusion strategies have been introduced, and their performance analyzed in various situations, such as the diffusion LMS [2], RLS [3], and APA [4], as well as several of their variants [5], [6].

By referring to estimating an optimal parameter vector at a node as a task, and according to the relations between the optimal parameter vectors over the entire network, diffusion networks are further divided into single-task and multitask networks. In single-task networks, all nodes estimate the same parameter vector. Typical works related to single-task networks include [7]–[11]. With multitask networks, multiple but related parameter vectors are inferred simultaneously in a cooperative manner, so as to improve the estimation accuracy by exploiting the similarities between tasks. These similarities can be promoted with appropriate regularization terms. Squared  $\ell_2$ -norm regularization is used in [12], and  $\ell_1$ -norm regularization is considered in [13], [14]. For the latter, a subgradient and a proximal algorithm are introduced in [13] and [14], respectively. In [15] and [16], the authors derive solutions for other classes of multitask problems where the relations between the nodes are defined by common latent representations or local linear equality constraints, respectively. In [17], the authors solve a multitask problem by estimating the combination matrix. In [18], the authors address multitask problems over asynchronous networks and carry out a detailed theoretical analysis. In [19], the performance of multitask diffusion networks is analyzed for correlated noise and regressors. In [20], the authors propose a combination framework that aggregates several diffusion strategies. All the algorithms cited above are based on the diffusion LMS. In [21], the authors extend the diffusion APA to multitask framework in order to improve the robustness against correlated regressors. In [22], the authors improve the performance of the multitask diffusion APA via controlled inter-cluster cooperation. In [23], the authors propose a clustered multitask partial diffusion APA that transmits only a subset of the entries of the intermediate estimates, to provide a trade-off between the estimation performance and communication cost.

Multitask learning considerably enriches the possibilities of diffusion networks. Beyond the few examples listed above, there are also applications where the optimal parameter vectors have a jointly sparse structure, namely, local solutions have the same sparse support. Applications include, for instance, distributed spectrum sensing and channel identification in underwater and wireless communication networks with multiple sensors [24]. In spectrum sensing of sparse wide-band spectra in distributed wireless sensor networks, the parameter vectors at different nodes share the same intrinsic sparse structure. However, due to different channel fading effects, shadowing effects and transmission losses, the values at the nonzero entries of these parameter vectors are different [24]–[26]. For underwater and wireless communication networks, it has been shown that real-world underwater channels [27] and wireless communication channels [28] are inherently sparse with large delay spread. In addition,

The work of Jie Chen was supported in part by NSFC grant 61671382. The work of Jingdong Chen was supported in part by NSFC grant 61425005. The work of C. Richard work was supported through the UCA JEDI Investments in the Future project with the reference number ANR-15-IDEX-0001. D. Jin, J. Chen and J. Chen are with Centre of Intelligent Acoustics and Immersive Communications at School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an, China (danqijin@mail.nwpu.edu.cn, dr.jie.chen@ieee.org, jingdongchen@ieee.org). C. Richard is with Université Côte d'Azur, OCA, CNRS, France (cedric.richard@unice.fr).

the channel supports for neighboring antennas or nodes are approximately the same [28]. Indeed, times of arrival for closely spaced nodes and antennas are quite close, though the tap weights are actually different [24]. For more detailed information about this application, see Section VI-C. Taking jointly sparse structure into consideration can greatly improve the network performance.

Several works have been proposed to address problems with jointly sparse structure over diffusion networks. In [29], the authors consider the mixed  $\ell_{2,0}$ -norm. In [24], the authors devise an algorithm with regularizers such as  $\ell_{2,0}$ ,  $\ell_{2,1}$  and reweighted  $\ell_{2,1}$  regularizers. They also conduct theoretical analyses of the algorithm behavior in the mean and mean-square sense. However, studies in [24], [29] are based on gradient-descent schemes, and use subgradient for non-differential regularization terms. For subgradient approaches, the subdifferential at a point may not be a singleton set, that is, it may be empty or consist of several elements. As a result, one may get stuck or have to choose one, respectively. Furthermore, even if the subdifferential is a singleton at each step, it might be highly discontinuous, so small deviations might lead to a singular behavior of the algorithm over iterations [30]. In this paper, we propose to use proximal operators to address the jointly sparse estimation problem and then avoid the weaknesses of subgradient-based approaches mentioned above. Proximal algorithms result in subproblems that often admit closed-form solutions or that can be solved efficiently with simple specialized methods [31]. In addition, they guarantee better convergence rates and stability than subgradient approaches [30], [32]. For single-task learning problems, proximal algorithms were first introduced in [33] to estimate sparse parameter vectors. They were used in [34] and [35] to optimize general stochastic costs with non-differential regularizers, and non-smooth regularizers, respectively. For multitask learning problems, related works only include [14]. In that paper, the authors consider the situation where the network is divided into several clusters. All nodes in a cluster are interested in estimating the same parameter vector, while nodes in adjacent clusters estimate parameter vectors that have a large number of similar entries. As a result the authors derive a closed-form proximal solution for an  $\ell_1$ -norm regularizer that promotes similarities among clusters. A theoretical analysis of the steady-state behavior is provided. The algorithm in [14] is based on an extra exchange step of observations between neighboring nodes within a same cluster as well as a fusion step that averages local estimates. The proximal algorithm derived in the current work does not include these two steps.

In this paper, we introduce a proximal diffusion LMS strategy for multitask networks with jointly sparse structure. The main contributions of this work are summarized as follows:

- We derive approximate closed-form expressions for the  $\ell_{\infty,1}$ -norm and reweighted  $\ell_{\infty,1}$ -norm proximal operators.
- We derive a proximal multitask diffusion LMS algorithm to solve problems with jointly sparse structure.
- We conduct a theoretical analysis of the algorithm performance, including a condition for stability and a study of its transient behavior in the mean sense.

**Notation.** Normal font  $x$ , boldface small letters  $\mathbf{x}$  and capital letters  $\mathbf{X}$  denote scalars, column vectors and matrices, respectively. Symbol  $[\cdot]_m$  denotes the  $m$ -th entry of its vector argument. The superscript  $(\cdot)^\top$  denotes the transpose operator. The mathematical expectation is denoted by  $\mathbb{E}\{\cdot\}$ . The Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$  is denoted by  $\mathcal{N}(\mu, \sigma^2)$ . Operator  $\|\cdot\|$  takes the absolute value of its scalar or vector argument. Operator  $\max\{\cdot, \cdot\}$  extracts the maximum value of its two arguments. Operator  $\text{diag}\{\cdot\}$  generates a diagonal matrix from its argument. Symbol  $\preceq$  denotes a component-wise inequality. The symbol  $\otimes$  denotes the Kronecker product. The set  $\mathcal{N}_k$  denotes the neighbors of node  $k$ , including  $k$  itself, and  $|\mathcal{N}_k|$  denotes its cardinality. The  $\mathcal{N}_k^-$  denotes the neighbors of node  $k$ , excluding node  $k$ . Vector  $\mathbb{1}_L$  is the all-one vector of dimension  $L \times 1$ .

## II. PROBLEM FORMULATION

Consider a connected network consisting of  $N$  nodes. Each node has access to streaming data  $\{d_{k,n}, \mathbf{x}_{k,n}\}$ , where  $\mathbf{x}_{k,n}$  is the  $L \times 1$  real-valued regression vector at node  $k$  and time instant  $n$ , and  $d_{k,n}$  denotes the observed real-valued signal. We assume that the data at each agent  $k$  and time instant  $n$  are related via the linear model:

$$d_{k,n} = \mathbf{x}_{k,n}^\top \mathbf{w}_k^* + z_{k,n}, \quad (1)$$

where  $\mathbf{w}_k^* \in \mathbb{R}^{L \times 1}$  is the unknown system vector to estimate, and  $z_{k,n}$  is a zero-mean additive noise. We assume that  $z_{k,n}$  is independent of any other signal. Further, we assume that vectors  $\{\mathbf{w}_k^*\}_{k=1}^N$  are jointly sparse, namely, not only each  $\mathbf{w}_k^*$  is a sparse vector but, in addition, they all have the same support. The support of  $\mathbf{w}_k^*$  is defined as [36]:

$$\text{supp}(\mathbf{w}_k^*) \triangleq \{j : [\mathbf{w}_k^*]_j \neq 0\}, \quad (2)$$

where  $[\mathbf{w}_k^*]_j$  is the  $j$ -th entry of  $\mathbf{w}_k^*$ . By definition (2), jointly sparse structure means that:

$$\text{supp}(\mathbf{w}_1^*) = \dots = \text{supp}(\mathbf{w}_k^*) = \dots = \text{supp}(\mathbf{w}_N^*). \quad (3)$$

Since it is not trivial to solve problems with jointly sparse structure directly, several mixed-norms have been introduced in the literature to alleviate this problem. They include the mixed  $\ell_{2,1}$ -norm [37],  $\ell_{\infty,1}$ -norm [38] and their reweighted versions [39]. We collect  $\mathbf{w}_\ell^*$  over the entire network into an  $L \times N$  matrix, and we replace the  $k$ -th column  $\mathbf{w}_k^*$  by the optimization

variable  $\mathbf{w}_k$ . This leads us to consider:

$$\mathbf{W}_k^{\text{glob}} \triangleq [\mathbf{w}_1^* \quad \cdots \quad \mathbf{w}_{k-1}^* \quad \mathbf{w}_k \quad \mathbf{w}_{k+1}^* \quad \cdots \quad \mathbf{w}_N^*]. \quad (4)$$

Since matrices defined in (4) differ from node to node, we use the subscript  $k$  in notation  $\mathbf{W}_k^{\text{glob}}$  to distinguish them. Evaluating the mixed  $\ell_{p,1}$ -norm of matrix  $\mathbf{W}_k^{\text{glob}}$  with  $p = 2$  or  $\infty$ , so as to promote the jointly sparse structure, results in the following two steps:

**Step 1:** Evaluate the  $\ell_p$ -norm of each row of  $\mathbf{W}_k^{\text{glob}}$ , and stack the results into an intermediate vector of dimension  $L \times 1$ ;

**Step 2:** Evaluate the  $\ell_1$ -norm of the obtained intermediate vector to promote sparsity.

Though  $\ell_{2,1}$ -norm can be more efficient in some cases [40], we shall consider the  $\ell_{\infty,1}$ -norm and its reweighted form to promote the joint-sparsity. It is shown in [41] that the  $\ell_{\infty,1}$  relaxation is exact in the case of normalized nonrepeating data. [38] introduced a recursive adaptive group lasso algorithm for real-time penalized least squares prediction by using the  $\ell_{1,\infty}$  regularization (The authors defined  $\ell_{\infty,1}$ -norm as  $\ell_{1,\infty}$ -norm in [38] for a vector with group sparsity structure). Each update minimizes a convex but nondifferentiable function optimization problem. The authors developed an online homotopy method to reduce the computational complexity. [41] proposed a collaborative convex framework for factoring a data matrix into a nonnegative matrices product. The authors used  $\ell_{\infty,1}$  regularization to select the dictionary from the data and shown that this leads to an exact convex relaxation of  $\ell_0$  regularization in the case of distinct noise-free data. [42] proposed a class of group sparse RLS algorithms by using different penalty term to promote the group sparsity, where the  $\ell_{\infty,1}$ -norm is used as one of the penalty term and the subgradient of  $\ell_{\infty,1}$ -norm is evaluated. Both  $\ell_{\infty,1}$ -norm and its reweighted form are element-wise separable, which facilitates the derivation of the proximal operator in the current work. As we focus on distributed processing in the current work, only local information exchange is authorized. Thus, we restrict  $\mathbf{W}_k^{\text{glob}}$  to the local quantity:

$$\mathbf{W}_k \triangleq [\mathbf{w}_k, \mathbf{w}_\ell^* \text{ with } \ell \in \mathcal{N}_k^-] \in \mathbb{R}^{L \times |\mathcal{N}_k|}. \quad (5)$$

We consider that the columns of  $\mathbf{W}_k$  are sorted in increasing order according to the values of  $k$  and  $\ell$ . In the sequel, we shall show how to evaluate the (reweighted)  $\ell_{\infty,1}$ -norm of  $\mathbf{W}_k$ .

### III. PROXIMAL MULTITASK DIFFUSION LMS

Before proceeding, we rewrite  $\mathbf{W}_k$  in another way in order to facilitate the presentation:

$$\mathbf{W}_k = [\bar{\mathbf{w}}_{k,1}^\top \quad \cdots \quad \bar{\mathbf{w}}_{k,m}^\top \quad \cdots \quad \bar{\mathbf{w}}_{k,L}^\top]^\top, \quad (6)$$

where  $\bar{\mathbf{w}}_{k,m}$  is the  $m$ -th row of matrix  $\mathbf{W}_k$  and defined by

$$\bar{\mathbf{w}}_{k,m} \triangleq [[\mathbf{w}_k]_m, [\mathbf{w}_\ell^*]_m \text{ with } \ell \in \mathcal{N}_k^-] \in \mathbb{R}^{1 \times |\mathcal{N}_k|} \quad (7)$$

with  $[\mathbf{w}_k]_m$  being the  $m$ -th entry of  $\mathbf{w}_k$ . Definition (7) means that  $\bar{\mathbf{w}}_{k,m}$  is a *row* vector. Note that (5) and (6) are actually the same matrix, but represented in different ways. Also,  $\|\bar{\mathbf{w}}_{k,m}\|_\infty$  is defined as

$$\|\bar{\mathbf{w}}_{k,m}\|_\infty \triangleq \max\{[|\mathbf{w}_k]_m|, [|\mathbf{w}_\ell^*]_m| \text{ with } \ell \in \mathcal{N}_k^-\} \quad (8)$$

which will be used in  $g_1(\mathbf{w}_k)$  and  $g_2(\mathbf{w}_k)$  of (11) and (12), respectively.

To determine the unknown vectors  $\mathbf{w}_k^*$  with jointly sparse structure, we consider the regularized cost at node  $k$ :

$$J_k(\mathbf{w}_k) = J'_k(\mathbf{w}_k) + \lambda_k g_i(\mathbf{w}_k), \quad (9)$$

where  $J'_k(\mathbf{w}_k)$  is the mean-square error (MSE) defined as:

$$J'_k(\mathbf{w}_k) \triangleq \frac{1}{2} \mathbb{E} \{ |d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_k|^2 \}. \quad (10)$$

The nonnegative parameter  $\lambda_k$  is used to control the regularization strength,  $g_1(\mathbf{w}_k)$  and  $g_2(\mathbf{w}_k)$  evaluate the  $\ell_{\infty,1}$ -norm and reweighted  $\ell_{\infty,1}$ -norm of  $\mathbf{W}_k$ , respectively, with:

$$g_1(\mathbf{w}_k) \triangleq \sum_{m=1}^L \|\bar{\mathbf{w}}_{k,m}\|_\infty, \quad (11)$$

$$g_2(\mathbf{w}_k) \triangleq \sum_{m=1}^L \log \left[ 1 + \frac{\|\bar{\mathbf{w}}_{k,m}\|_\infty}{\varepsilon} \right], \quad (12)$$

where  $\varepsilon > 0$  is a parameter set by the user. Interpretations about the  $\ell_{\infty,1}$ -norm  $g_1(\mathbf{w}_k)$  and reweighted  $\ell_{\infty,1}$ -norm  $g_2(\mathbf{w}_k)$  are provided in Appendix A. At each node  $k$ , we then consider the convex optimization problem [43]:

$$\mathbf{w}_k^\dagger = \arg \min_{\mathbf{w}_k} J_k(\mathbf{w}_k). \quad (13)$$

Within the context of online learning, such optimization problem is usually solved via subgradient-based methods. In this paper, we propose to devise a proximal algorithm since these algorithms are usually more stable and have a better convergence rate than subgradient iterations [30], [32], [33]. Proximal gradient iteration consists of [31]:

$$\mathbf{w}_{k,n+1} = \text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)} \left( \mathbf{w}_{k,n} - \mu_k \nabla J'_k(\mathbf{w}_{k,n}) \right), \quad (14)$$

where  $\mu_k$  is a positive small step-size, and:

$$\text{prox}_{\lambda, g_i(\mathbf{w}_k)}(\mathbf{v}) \triangleq \arg \min_{\mathbf{w}_k} \left( g_i(\mathbf{w}_k) + \frac{1}{2\lambda} \|\mathbf{w}_k - \mathbf{v}\|_2^2 \right). \quad (15)$$

is the proximal operator. By introducing the intermediate quantity  $\boldsymbol{\psi}_{k,n+1}$ , we further decompose (14) into two steps:

$$\boldsymbol{\psi}_{k,n+1} = \mathbf{w}_{k,n} - \mu_k \nabla J'_k(\mathbf{w}_{k,n}), \quad (16)$$

$$\mathbf{w}_{k,n+1} = \text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)}(\boldsymbol{\psi}_{k,n+1}). \quad (17)$$

Equation (16) is the local update step, and (17) is the proximal step. Calculating the gradient of  $J'_k(\mathbf{w}_k)$  at  $\mathbf{w}_{k,n}$  and approximating the unknown statistical quantities with instantaneous quantities, (16) becomes:

$$\boldsymbol{\psi}_{k,n+1} = \mathbf{w}_{k,n} + \mu_k \mathbf{x}_{k,n} (d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n}). \quad (18)$$

Combining (17) and (18) yields the proximal multitask diffusion LMS algorithm for networks with jointly sparse structure reported in **Algorithm 1**. This algorithm differs from the standard diffusion procedure devised in [8], [9] in the sense that it does not combine, in an explicit manner, the intermediate estimates  $\boldsymbol{\psi}_{\ell,n+1}$  in the neighborhood of each node  $k$ . In our algorithm, the information on the support of the local estimates is shared by neighboring nodes via  $g_i(\mathbf{w}_k)$ .

---

**Algorithm 1** Proximal multitask diffusion LMS

---

Initialize  $\mathbf{w}_{k,0}$  for all  $k = 1, 2, \dots, N$ , and repeat:

$$\begin{cases} \boldsymbol{\psi}_{k,n+1} = \mathbf{w}_{k,n} + \mu_k \mathbf{x}_{k,n} (d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n}) \\ \mathbf{w}_{k,n+1} = \text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)}(\boldsymbol{\psi}_{k,n+1}) \end{cases} \quad (19)$$


---

#### IV. PROXIMAL OPERATORS EVALUATION

To apply the algorithm, we need to derive closed-form expressions for the proximal operators. Equation (17) becomes:

$$\begin{aligned} \mathbf{w}_{k,n+1} &= \text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)}(\boldsymbol{\psi}_{k,n+1}) \\ &= \arg \min_{\mathbf{w}_k} \left( g_i(\mathbf{w}_k) + \frac{1}{2\mu_k \lambda_k} \|\mathbf{w}_k - \boldsymbol{\psi}_{k,n+1}\|_2^2 \right). \end{aligned} \quad (20)$$

As  $g_i(\mathbf{w}_k)$  is separable over its all entries, the proximal operator can be evaluated in an element-wise manner as [31]:

$$\begin{aligned} [\text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)}(\boldsymbol{\psi}_{k,n+1})]_m & \\ &= \text{prox}_{\mu_k \lambda_k, g_{i,m}([\mathbf{w}_k]_m)}([\boldsymbol{\psi}_{k,n+1}]_m) \end{aligned} \quad (21)$$

with:

$$g_{1,m}([\mathbf{w}_k]_m) \triangleq \|\bar{\mathbf{w}}_{k,m}\|_\infty \quad (22)$$

$$g_{2,m}([\mathbf{w}_k]_m) \triangleq \log \left[ 1 + \frac{\|\bar{\mathbf{w}}_{k,m}\|_\infty}{\varepsilon} \right], \quad (23)$$

where  $\bar{\mathbf{w}}_{k,m}$  is the  $m$ -th row of matrix  $\mathbf{W}_k$  in (5).

##### A. Approximate proximal operator of $\ell_{\infty,1}$ -norm

Substituting  $g_{1,m}$  into (21), we obtain:

$$\begin{aligned} [\mathbf{w}_{k,n+1}]_m &= \arg \min_{[\mathbf{w}_k]_m} \left( \max\{ |[\mathbf{w}_k]_m|, |[\mathbf{w}_\ell^*]_m| \text{ with } \ell \in \mathcal{N}_k^- \} \right. \\ &\quad \left. + \frac{1}{2\mu_k \lambda_k} ([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2 \right). \end{aligned} \quad (24)$$

For the sake of simpler notations, we shall denote  $[\mathbf{w}_{k,n+1}]_m$  by  $\hat{w}$  as long as there is no ambiguity. We shall also denote the maximal value of  $|[\mathbf{w}_\ell^*]_m|$  for  $\ell \in \mathcal{N}_k^-$  as  $[\mathbf{w}_k^o]_m$ . According to the relation between  $|[\mathbf{w}_k]_m|$  and  $[\mathbf{w}_k^o]_m$ , we further split the problem into the following two cases:

- Case 1:  $|\mathbf{w}_k]_m| < [\mathbf{w}_k^o]_m$ . In this case, (24) becomes:

$$\hat{w} = \underset{\substack{[\mathbf{w}_k]_m \\ |[\mathbf{w}_k]_m| < [\mathbf{w}_k^o]_m}}{\operatorname{argmin}} [\mathbf{w}_k^o]_m + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2. \quad (25)$$

Since the first term on the right-hand side (RHS) of (25) does not depend on  $[\mathbf{w}_k]_m$ , we conclude that:

$$\hat{w} = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\mathbf{w}_k^o]_m, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m \\ -[\mathbf{w}_k^o]_m, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m. \end{cases} \quad (26)$$

- Case 2:  $|\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m$ . Equation (24) becomes:

$$\hat{w} = \underset{\substack{[\mathbf{w}_k]_m \\ |[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m}}{\operatorname{argmin}} |[\mathbf{w}_k]_m| + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2. \quad (27)$$

We shall first discard the constraint  $|\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m$  for reasons of simplicity, and denote by  $\hat{w}^o$  the solution of the unconstrained problem. As the cost function in (27) is convex on  $\mathbb{R}$ , this constraint will be taken into account in the course of the calculation. Consider first:

$$\begin{aligned} \hat{w}^o &= \underset{[\mathbf{w}_k]_m}{\operatorname{argmin}} |[\mathbf{w}_k]_m| + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2 \\ &= \operatorname{prox}_{\mu_k\lambda_k, s([\mathbf{w}_k]_m)}([\boldsymbol{\psi}_{k,n+1}]_m) \end{aligned} \quad (28)$$

where function  $s([\mathbf{w}_k]_m)$  is defined as:

$$s([\mathbf{w}_k]_m) \triangleq |[\mathbf{w}_k]_m|. \quad (29)$$

The optimality condition of (28) says that zero belongs to the subgradient set at the minimizer  $\hat{w}^o$ , that is, [31]

$$0 \in \partial|\hat{w}^o| + \frac{1}{\mu_k\lambda_k} (\hat{w}^o - [\boldsymbol{\psi}_{k,n+1}]_m) \quad (30)$$

This means that:

$$([\boldsymbol{\psi}_{k,n+1}]_m - \hat{w}^o) \in \mu_k\lambda_k \partial|\hat{w}^o| \quad (31)$$

with

$$\partial|\hat{w}^o| = \begin{cases} -1, & \text{if } \hat{w}^o < 0 \\ +1, & \text{if } \hat{w}^o > 0 \\ [-1, 1], & \text{if } \hat{w}^o = 0 \end{cases} \quad (32)$$

the subdifferential of the non-differentiable function  $|\hat{w}^o|$  [31]. The third case in (32) means that, at  $\hat{w}^o = 0$ ,  $\partial|\hat{w}^o|$  can be any value within  $[-1, 1]$ .

Condition (31) leads to:

$$\hat{w}^o = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m + \mu_k\lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m < -\mu_k\lambda_k \\ [\boldsymbol{\psi}_{k,n+1}]_m - \mu_k\lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m > \mu_k\lambda_k \\ 0, & \text{otherwise.} \end{cases} \quad (33)$$

If  $[\mathbf{w}_k^o]_m = 0$ , problem (27) becomes unconstrained and we have:

$$\hat{w} = \hat{w}^o \quad (34)$$

Otherwise, since problem (27) is convex on  $\mathbb{R}$ , considering constraint  $|\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m > 0$  with (33) yields:

$$\hat{w} = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m + \mu_k\lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m - \mu_k\lambda_k \\ -[\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{w}_k^o]_m - \mu_k\lambda_k < [\boldsymbol{\psi}_{k,n+1}]_m < 0 \\ -[\mathbf{w}_k^o]_m \text{ or } [\mathbf{w}_k^o]_m, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m = 0 \\ [\mathbf{w}_k^o]_m, & \text{if } 0 < [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{w}_k^o]_m + \mu_k\lambda_k \\ [\boldsymbol{\psi}_{k,n+1}]_m - \mu_k\lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m + \mu_k\lambda_k \end{cases} \quad (35)$$

To evaluate the proximal operator (24), several issues have to be addressed.

1. One of the main issues is that we first need to know which of (26), (34) or (35) has to be applied as the proximal operator of (24). We shall now consider the following two cases:  $[\mathbf{w}_k^o]_m = 0$  and  $[\mathbf{w}_k^o]_m > 0$ .

- Case A:  $[\mathbf{w}_k^o]_m = 0$ . Since condition  $|\mathbf{w}_k]_m| < [\mathbf{w}_k^o]_m$  of Case 1 cannot hold, we only consider Case 2. The proximal operator is given by  $\hat{w}^o$  in (33). Interestingly, note that (33) is the proximal operator of the  $\ell_1$ -norm regularizer for a scalar.
- Case B:  $[\mathbf{w}_k^o]_m > 0$ . Proximal operators (26) and (35) hold simultaneously. We shall choose the solution that minimizes the cost (24). As shown in Appendix B, we finally arrive at the following expression:

$$\hat{w} = \begin{cases} [\psi_{k,n+1}]_m + \mu_k \lambda_k, & \text{if } [\psi_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m - \mu_k \lambda_k \\ -[\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{w}_k^o]_m - \mu_k \lambda_k < [\psi_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m \\ [\psi_{k,n+1}]_m, & \text{if } |[\psi_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\mathbf{w}_k^o]_m, & \text{if } [\mathbf{w}_k^o]_m \leq [\psi_{k,n+1}]_m < [\mathbf{w}_k^o]_m + \mu_k \lambda_k \\ [\psi_{k,n+1}]_m - \mu_k \lambda_k, & \text{if } [\mathbf{w}_k^o]_m + \mu_k \lambda_k \leq [\psi_{k,n+1}]_m \end{cases} \quad (36)$$

2. Another issue is that  $\hat{w}$  cannot be evaluated with (33) and (36) since  $[\mathbf{w}_k^o]_m$  is unknown. To fix this problem, we follow a local one-step approximation strategy that has already proven its effectiveness in the literature [17], and which consists of using  $\psi_{\ell,n+1}$  as an approximation of  $\mathbf{w}_\ell^*$ . One of the benefits of this approximation is that  $\psi_{\ell,n+1}$  can be transmitted by node  $\ell$  to node  $k$  if the latter is in its neighborhood. An approximation of  $[\mathbf{w}_k^o]_m$  is then given by  $\max_{\ell \in \mathcal{N}_k^-} \{|\psi_{\ell,n+1}]_m|\}$ , which allows node  $k$  to evaluate its proximal operator. This approximation will be taken into consideration in the theoretical analysis in Section V-A.

3. Before addressing the last issue, we need to point out the prominent role of Case A compared to Case B: unlike the proximal operator (36) in Case B, only the proximal operator (33) in Case A has the capability to drive  $[\mathbf{w}_k]_m$  to zero and promote sparsity. Observe that condition  $[\mathbf{w}_k^o]_m = 0$  has to be satisfied to trigger Case A, otherwise Case B is considered. Within the context of online learning with stochastic gradient descent algorithms, due to the existence of gradient noise, the estimates of  $[\mathbf{w}_k^o]_m$  for zero-valued entries are actually nonzero-valued but vary around zero. Condition  $[\mathbf{w}_k^o]_m = 0$  of Case A is thus seldom satisfied, and  $[\psi_{k,n+1}]_m$  is not driven to zero. To promote the sparsity of the estimates, since the true values of the non-zero entries are usually far away from zero, we introduce a small positive threshold value  $\tau_1$  instead of zero to make a distinguish between zero-valued and nonzero-valued entries. As a consequence, we arrive at conditions  $[\mathbf{w}_k^o]_m \leq \tau_1$  to trigger Case A and  $[\mathbf{w}_k^o]_m > \tau_1$  to select Case B.

We summarize our method in **Algorithm 2**.

---

**Algorithm 2** Approximate proximal operator of  $\ell_{\infty,1}$ -norm

---

**Initialization:** Choose threshold value  $\tau_1 > 0$ .

**Proximal operator:** At each instant  $n \geq 0$ , for each node  $k$ , utilize  $\psi_{k,n+1}$  to evaluate  $\mathbf{w}_{k,n+1}$  in an elementwise manner:

- 1) Calculate  $[\mathbf{w}_k^o]_m$  as the maximal value of  $|\psi_{\ell,n+1}]_m|$  for all  $\ell \in \mathcal{N}_k^-$ ;
  - 2) If  $[\mathbf{w}_k^o]_m \leq \tau_1$ , then calculate  $[\mathbf{w}_{k,n+1}]_m$  as  $\hat{w}^o$  via (33);
  - 3) If  $[\mathbf{w}_k^o]_m > \tau_1$ , then calculate  $[\mathbf{w}_{k,n+1}]_m$  as  $\hat{w}$  via (36).
- 

*B. Approximate proximal operator of reweighted  $\ell_{\infty,1}$ -norm*

Substituting the expression of  $g_{2,m}$  into (21), we arrive at:

$$[\mathbf{w}_{k,n+1}]_m = \underset{[\mathbf{w}_k]_m}{\operatorname{argmin}} \left[ \log \left( 1 + \frac{\max\{|\mathbf{w}_k]_m|, |\mathbf{w}_\ell^*]_m| : \forall \ell \in \mathcal{N}_k^-\}}{\varepsilon} \right) + \frac{1}{2\mu_k \lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2 \right]. \quad (37)$$

Considering the same definition for  $[\mathbf{w}_k^o]_m$  as before, and denoting  $[\mathbf{w}_{k,n+1}]_m$  by  $\hat{w}$  for the sake of conciseness, the problem can be split into two cases depending on the order relation between  $[\mathbf{w}_k^o]_m$  and  $|\mathbf{w}_k]_m|$ :

- Case 1:  $|\mathbf{w}_k]_m| < [\mathbf{w}_k^o]_m$ . Equation (37) becomes:

$$\hat{w} = \underset{\substack{[\mathbf{w}_k]_m \\ |\mathbf{w}_k]_m| < [\mathbf{w}_k^o]_m}}{\operatorname{argmin}} \log \left( 1 + \frac{[\mathbf{w}_k^o]_m}{\varepsilon} \right) + \frac{1}{2\mu_k \lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \quad (38)$$

This problem has a closed-form solution given by:

$$\hat{w} = \begin{cases} [\psi_{k,n+1}]_m, & \text{if } |[\psi_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\mathbf{w}_k^o]_m, & \text{if } [\psi_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m \\ -[\mathbf{w}_k^o]_m, & \text{if } [\psi_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m. \end{cases} \quad (39)$$

- Case 2:  $|[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m$ . In this case, (37) becomes:

$$\begin{aligned} \hat{w} = \underset{\substack{[\mathbf{w}_k]_m \\ |[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m}}{\text{argmin}} \quad & \log\left(1 + \frac{|[\mathbf{w}_k]_m|}{\varepsilon}\right) \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \end{aligned} \quad (40)$$

By first discarding the constraint  $|[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m$  as we did it before, we obtain the following unconstrained optimization problem:

$$\begin{aligned} \hat{w}^o = \underset{[\mathbf{w}_k]_m}{\text{argmin}} \quad & \log\left(1 + \frac{|[\mathbf{w}_k]_m|}{\varepsilon}\right) \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \end{aligned} \quad (41)$$

Deriving the solution  $\hat{w}^o$  directly from the optimality condition is a tough problem. We shall now alleviate this issue with a Majorization-Minimization (MM) approach [44], [45]. Conceptually, MM algorithms work by iteratively minimizing a simple surrogate function majorizing a given objective function. By introducing the auxiliary variable  $[\mathbf{u}_k]_m$ , problem (41) can be written as:

$$\begin{aligned} \hat{w}^o = \underset{\substack{[\mathbf{w}_k]_m, [\mathbf{u}_k]_m \\ |[\mathbf{w}_k]_m| \leq [\mathbf{u}_k]_m}}{\text{argmin}} \quad & \log\left(1 + \frac{[\mathbf{u}_k]_m}{\varepsilon}\right) \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \end{aligned} \quad (42)$$

Define:

$$f([\mathbf{u}_k]_m) \triangleq \log\left(1 + \frac{[\mathbf{u}_k]_m}{\varepsilon}\right). \quad (43)$$

It can be observed that function  $f([\mathbf{u}_k]_m)$  is concave and below its tangent. So we have:

$$\begin{aligned} f([\mathbf{u}_k]_m) \leq & f([\mathbf{u}_k^{(t)}]_m) \\ & + f'([\mathbf{u}_k^{(t)}]_m) \cdot ([\mathbf{u}_k]_m - [\mathbf{u}_k^{(t)}]_m) \end{aligned} \quad (44)$$

with  $f'(\cdot)$  the first-order derivative of  $f(\cdot)$ , and where the superscript  $(t)$  denotes the iteration index. Using (44), we can construct a surrogate function for (42) as follows:

$$\begin{aligned} h([\mathbf{u}_k]_m, [\mathbf{w}_k]_m) \triangleq & \log\left(1 + \frac{[\mathbf{u}_k^{(t)}]_m}{\varepsilon}\right) + \frac{[\mathbf{u}_k]_m - [\mathbf{u}_k^{(t)}]_m}{\varepsilon + [\mathbf{u}_k^{(t)}]_m} \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2 \end{aligned} \quad (45)$$

Using the driving principle of MM algorithms, this yields the following optimization problem:

$$\begin{aligned} ([\mathbf{w}_k^{(t+1)}]_m, [\mathbf{u}_k^{(t+1)}]_m) = \underset{\substack{[\mathbf{w}_k]_m, [\mathbf{u}_k]_m \\ |[\mathbf{w}_k]_m| \leq [\mathbf{u}_k]_m}}{\text{argmin}} \quad & \frac{[\mathbf{u}_k]_m}{\varepsilon + [\mathbf{u}_k^{(t)}]_m} \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2 \end{aligned} \quad (46)$$

which is equivalent to:

$$\begin{aligned} [\mathbf{w}_k^{(t+1)}]_m = \underset{[\mathbf{w}_k]_m}{\text{argmin}} \quad & \frac{|[\mathbf{w}_k]_m|}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \\ & + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \end{aligned} \quad (47)$$



Note that, in the construction of (46), we have dropped several constant terms unrelated to  $[\mathbf{u}_k]_m$  and  $[\mathbf{w}_k]_m$ . Using the optimality condition of (47) at the minimizer  $[\mathbf{w}_k^{(t+1)}]_m$ , we have:

$$([\boldsymbol{\psi}_{k,n+1}]_m - [\mathbf{w}_k^{(t+1)}]_m) \in \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \cdot \partial |[\mathbf{w}_k^{(t+1)}]_m|. \quad (48)$$

Considering (32) with (48) yields:

$$[\mathbf{w}_k^{(t+1)}]_m = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m + \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m < \frac{-\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \\ [\boldsymbol{\psi}_{k,n+1}]_m - \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m > \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \\ 0 & \text{otherwise.} \end{cases} \quad (49)$$

Performing one iteration of (49) with  $[\mathbf{w}_k^{(t)}]_m = [\mathbf{w}_{k,n}]_m$  at each instant  $n$  can be sufficient for approximating  $\hat{w}^o$  with enough precision, as illustrated in the sequel. This value is then used to obtain an approximate solution  $\hat{w}$  of problem (40) by taking the constraint  $|[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m$  into account.

If  $[\mathbf{w}_k^o]_m = 0$ , problem (40) becomes unconstrained and we have:

$$\hat{w} = \hat{w}^o \quad (50)$$

where  $\hat{w}^o$  is obtained by (49).

Otherwise, since problem (40) is convex on  $\mathbb{R}$ , considering constraint  $|[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m > 0$  with (49) yields:

$$\hat{w} = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m + \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{c}_k]_m \\ -[\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{c}_k]_m < [\boldsymbol{\psi}_{k,n+1}]_m < 0 \\ -[\mathbf{w}_k^o]_m \text{ or } [\mathbf{w}_k^o]_m, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m = 0 \\ [\mathbf{w}_k^o]_m, & \text{if } 0 < [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{c}_k]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m - \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{c}_k]_m \end{cases} \quad (51)$$

with  $[\mathbf{c}_k]_m$  a constant defined as:

$$[\mathbf{c}_k]_m \triangleq [\mathbf{w}_k^o]_m + \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}. \quad (52)$$

As with the  $\ell_{\infty,1}$ -norm, we have to fix several issues.

**1.** We first need to check which of (39), (49) and (51) has to be applied to calculate the proximal operator. We consider the following two cases:  $[\mathbf{w}_k^o]_m = 0$  and  $[\mathbf{w}_k^o]_m > 0$ .

- Case A:  $[\mathbf{w}_k^o]_m = 0$ . We focus on Case 2 since condition for Case 1 cannot hold. In Case 2, observe that any  $[\mathbf{w}_k]_m$  calculated as  $[\mathbf{w}_k^{(t+1)}]_m$  in (49) satisfies  $|[\mathbf{w}_k]_m| \geq [\mathbf{w}_k^o]_m = 0$ . Thus the proximal operator is given by (49).
- Case B:  $[\mathbf{w}_k^o]_m > 0$ . Case 1 and Case 2 can hold simultaneously. We must choose, among the resulting proximal operators, the one that minimizes cost (37). The derivation is provided in Appendix C. We obtain the following solution:

$$\hat{w} = \begin{cases} [\boldsymbol{\psi}_{k,n+1}]_m + \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{c}_k]_m \\ -[\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{c}_k]_m < [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\mathbf{w}_k^o]_m, & \text{if } [\mathbf{w}_k^o]_m \leq [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{c}_k]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m - \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{c}_k]_m \end{cases} \quad (53)$$

**2.** Unlike the proximal operator (53) in Case B, only the proximal operator (49) in Case A has the capability to drive  $[\mathbf{w}_k]_m$  to zero and promote sparsity. Thus, as in Section IV-A, we relax the condition  $[\mathbf{w}_k^o]_m = 0$  by introducing a small

positive tolerance  $\tau_2$  to distinguish between zero and nonzero entries. This leads to the condition  $[\mathbf{w}_k^o]_m \leq \tau_2$  for Case A, and  $[\mathbf{w}_k^o]_m > \tau_2$  for Case B.

We summarize our method in **Algorithm 3**.

---

**Algorithm 3** Approximate proximal operator of reweighted  $\ell_{\infty,1}$ -norm

---

**Initialization:** Choose threshold value  $\tau_2 > 0$  and  $\varepsilon > 0$ .

**Proximal operator:** At each instant  $n \geq 0$ , for each node  $k$ , utilize  $\boldsymbol{\psi}_{k,n+1}$  to evaluate  $\mathbf{w}_{k,n+1}$  in an element-wise manner:

- 1) Calculate  $[\mathbf{w}_k^o]_m$  as the maximal value of  $|\boldsymbol{\psi}_{\ell,n+1}|_m$  for all  $\ell \in \mathcal{N}_k^-$ ;
  - 2) If  $[\mathbf{w}_k^o]_m \leq \tau_2$ , then calculate  $[\mathbf{w}_{k,n+1}]_m$  as  $[\mathbf{w}_k^{(t+1)}]_m$  via (49) with  $[\mathbf{w}_k^{(t)}]_m = [\mathbf{w}_{k,n}]_m$ ;
  - 3) If  $[\mathbf{w}_k^o]_m > \tau_2$ , then calculate  $[\mathbf{w}_{k,n+1}]_m$  as  $\hat{w}$  via (53).
- 

## V. PERFORMANCE AND CONVERGENCE ANALYSES

In this section, we analyse the performance and convergence property of **Algorithm 2** and **Algorithm 3** in an unified framework. Quantities specifically related to  $\ell_{\infty,1}$ -norm or reweighted  $\ell_{\infty,1}$ -norm are distinguished by the superscripts <sup>(1)</sup> and <sup>(2)</sup>, respectively. Observe that the closed-forms expressions of both approximate proximal operators can be written compactly as:

$$\text{prox}_{\mu_k \lambda_k, g_i(\mathbf{w}_k)}(\boldsymbol{\psi}_{k,n+1}) = \boldsymbol{\psi}_{k,n+1} - \boldsymbol{\gamma}_{k,n+1}^{(i)}, \quad (54)$$

with  $\boldsymbol{\gamma}_{k,n+1}^{(i)}$  a  $(L \times 1)$ -dimensional vector.

For the  $\ell_{\infty,1}$ -norm, the  $m$ -th entry of  $\boldsymbol{\gamma}_{k,n+1}^{(1)}$  is given by:

$$[\boldsymbol{\gamma}_{k,n+1}^{(1)}]_m = \begin{cases} -\mu_k \lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m < -\mu_k \lambda_k \\ [\boldsymbol{\psi}_{k,n+1}]_m, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| \leq \mu_k \lambda_k \\ \mu_k \lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m > \mu_k \lambda_k \end{cases} \quad (55)$$

if  $[\mathbf{w}_k^o]_m \leq \tau_1$ , and:

$$[\boldsymbol{\gamma}_{k,n+1}^{(1)}]_m = \begin{cases} -\mu_k \lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m - \mu_k \lambda_k \\ [\boldsymbol{\psi}_{k,n+1}]_m + [\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{w}_k^o]_m - \mu_k \lambda_k < [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m \\ 0, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m - [\mathbf{w}_k^o]_m, & \text{if } [\mathbf{w}_k^o]_m \leq [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{w}_k^o]_m + \mu_k \lambda_k \\ \mu_k \lambda_k, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m + \mu_k \lambda_k \end{cases} \quad (56)$$

if  $[\mathbf{w}_k^o]_m > \tau_1$ . For the reweighted  $\ell_{\infty,1}$ -norm, the  $m$ -th entry of  $\boldsymbol{\gamma}_{k,n+1}^{(2)}$  is given by:

$$[\boldsymbol{\gamma}_{k,n+1}^{(2)}]_m = \begin{cases} \frac{-\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m < \frac{-\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \\ [\boldsymbol{\psi}_{k,n+1}]_m, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| \leq \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \\ \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m > \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \end{cases} \quad (57)$$

if  $[\mathbf{w}_k^o]_m \leq \tau_2$ , and:

$$[\boldsymbol{\gamma}_{k,n+1}^{(2)}]_m = \begin{cases} \frac{-\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{c}_k]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m + [\mathbf{w}_k^o]_m, & \text{if } -[\mathbf{c}_k]_m < [\boldsymbol{\psi}_{k,n+1}]_m \leq -[\mathbf{w}_k^o]_m \\ 0, & \text{if } |[\boldsymbol{\psi}_{k,n+1}]_m| < [\mathbf{w}_k^o]_m \\ [\boldsymbol{\psi}_{k,n+1}]_m - [\mathbf{w}_k^o]_m, & \text{if } [\mathbf{w}_k^o]_m \leq [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{c}_k]_m \\ \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}, & \text{if } [\boldsymbol{\psi}_{k,n+1}]_m \geq [\mathbf{c}_k]_m \end{cases} \quad (58)$$

if  $[\mathbf{w}_k^2]_m > \tau_2$ . We observe in (55)–(56) that:

$$|[\boldsymbol{\gamma}_{k,n+1}^{(1)}]_m| \leq \mu_k \lambda_k \quad (59)$$

Thus,  $\boldsymbol{\gamma}_{k,n+1}^{(1)}$  is absolutely bounded:  $|\boldsymbol{\gamma}_{k,n+1}^{(1)}| \preceq \mu_k \lambda_k \mathbb{1}_L$ .

From (57)–(58), we conclude that:

$$|\boldsymbol{\gamma}_{k,n+1}^{(2)}| \preceq b_k \mathbb{1}_L \quad (60)$$

where  $b_k$  is a constant defined as:

$$b_k \triangleq \frac{\mu_k \lambda_k}{\varepsilon}. \quad (61)$$

Substituting (54) into (19), we obtain:

$$\begin{cases} \boldsymbol{\psi}_{k,n+1} = \mathbf{w}_{k,n} + \mu_k \mathbf{x}_{k,n} (d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n}) \\ \mathbf{w}_{k,n+1} = \boldsymbol{\psi}_{k,n+1} - \boldsymbol{\gamma}_{k,n+1}^{(i)}. \end{cases} \quad (62)$$

We shall now analyse the performance of proximal multitask diffusion LMS algorithm based on expression (62).

Define

$$\tilde{\mathbf{w}}_{k,n+1} \triangleq \mathbf{w}_{k,n+1} - \mathbf{w}_k^*, \quad (63)$$

$$\tilde{\boldsymbol{\psi}}_{k,n+1} \triangleq \boldsymbol{\psi}_{k,n+1} - \mathbf{w}_k^*. \quad (64)$$

By collecting  $\mathbf{w}_k^*$ ,  $\mathbf{w}_{k,n+1}$ ,  $\boldsymbol{\psi}_{k,n+1}$ ,  $\tilde{\mathbf{w}}_{k,n+1}$ ,  $\tilde{\boldsymbol{\psi}}_{k,n+1}$ ,  $\boldsymbol{\gamma}_{k,n+1}^{(i)}$  over the entire network into block column vectors, we obtain quantities  $\mathbf{w}^*$ ,  $\mathbf{w}_{n+1}$ ,  $\boldsymbol{\psi}_{n+1}$ ,  $\tilde{\mathbf{w}}_{n+1}$ ,  $\tilde{\boldsymbol{\psi}}_{n+1}$ ,  $\boldsymbol{\gamma}_{n+1}^{(i)}$ , respectively. To facilitate the theoretical analysis, we introduce the following assumptions on the regression data and step-size. These assumptions are widely used in the analysis of adaptive filters [46]–[48] and diffusion networks [7], [9], [14].

**A1 (Independent Regressors):** The regression vector  $\mathbf{x}_{k,n}$ , generated from a zero-mean random process, is temporally stationary, white (over  $n$ ) and spatially independent (over  $k$ ) with covariance matrix  $\mathbf{R}_{x,k} = \mathbb{E}\{\mathbf{x}_{k,n} \mathbf{x}_{k,n}^\top\} > 0$ .

**A2 (Small step-sizes):** The step-sizes  $\mu_k$  of the network are small enough, so that terms on the higher-order powers of the step-sizes can be ignored.

### A. Mean behavior analysis

Subtracting  $\mathbf{w}_k^*$  from both sides of the first equation of (62), and using signal model (1) and block notations, we obtain:

$$\tilde{\boldsymbol{\psi}}_{n+1} = \mathbf{B}_n \tilde{\mathbf{w}}_n + \mathbf{U} \mathbf{h}_n, \quad (65)$$

where

$$\mathbf{B}_n \triangleq \mathbf{I} - \mathbf{U} \mathbf{M}_n \quad (66)$$

with

$$\mathbf{U} \triangleq \text{diag}\{\mu_1, \mu_2, \dots, \mu_N\} \otimes \mathbf{I}_L, \quad (67)$$

$$\mathbf{M}_n \triangleq \text{diag}\{\mathbf{x}_{1,n} \mathbf{x}_{1,n}^\top, \mathbf{x}_{2,n} \mathbf{x}_{2,n}^\top, \dots, \mathbf{x}_{N,n} \mathbf{x}_{N,n}^\top\}, \quad (68)$$

$$\mathbf{h}_n \triangleq \text{col}\{\mathbf{x}_{1,n} z_{1,n}, \mathbf{x}_{2,n} z_{2,n}, \dots, \mathbf{x}_{N,n} z_{N,n}\}. \quad (69)$$

Subtracting  $\mathbf{w}_k^*$  from both sides of the second equation of (62) and using block notations, we obtain:

$$\tilde{\mathbf{w}}_{n+1} = \tilde{\boldsymbol{\psi}}_{n+1} - \boldsymbol{\gamma}_{n+1}^{(i)}. \quad (70)$$

Combining (65) and (70) leads to:

$$\tilde{\mathbf{w}}_{n+1} = \mathbf{B}_n \tilde{\mathbf{w}}_n + \mathbf{U} \mathbf{h}_n - \boldsymbol{\gamma}_{n+1}^{(i)}. \quad (71)$$

Using **A1** and taking the expectation of (71), we obtain:

$$\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\} = \mathbf{B} \mathbb{E}\{\tilde{\mathbf{w}}_n\} - \mathbb{E}\{\boldsymbol{\gamma}_{n+1}^{(i)}\}, \quad (72)$$

where

$$\mathbf{B} \triangleq \mathbb{E}\{\mathbf{B}_n\} = \mathbf{I} - \mathbf{U} \mathbf{M}, \quad (73)$$

$$\mathbb{E}\{\gamma_{k,n+1}^{(1)}\}_m = \begin{cases} -\mu_k \lambda_k, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \leq -\mathbb{E}\{\mathbf{w}_k^o\}_m - \mu_k \lambda_k \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m + \mathbb{E}\{\mathbf{w}_k^o\}_m, & \text{if } -\mathbb{E}\{\mathbf{w}_k^o\}_m - \mu_k \lambda_k < [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \leq -\mathbb{E}\{\mathbf{w}_k^o\}_m \\ 0, & \text{if } |[\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m| < \mathbb{E}\{\mathbf{w}_k^o\}_m \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m - \mathbb{E}\{\mathbf{w}_k^o\}_m, & \text{if } \mathbb{E}\{\mathbf{w}_k^o\}_m \leq [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m < \mathbb{E}\{\mathbf{w}_k^o\}_m + \mu_k \lambda_k \\ \mu_k \lambda_k, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \geq \mathbb{E}\{\mathbf{w}_k^o\}_m + \mu_k \lambda_k \end{cases} \quad (77)$$

with

$$\mathbf{M} \triangleq \mathbb{E}\{\mathbf{M}_n\} = \text{diag}\{\mathbf{R}_{x,1}, \mathbf{R}_{x,2}, \dots, \mathbf{R}_{x,N}\}, \quad (74)$$

$$\mathbb{E}\{\boldsymbol{\gamma}_{n+1}^{(i)}\} = \text{col}\{\mathbb{E}\{\boldsymbol{\gamma}_{1,n+1}^{(i)}\}, \dots, \mathbb{E}\{\boldsymbol{\gamma}_{N,n+1}^{(i)}\}\}. \quad (75)$$

To analyze iteration (72), we need to derive the explicit expression of  $\mathbb{E}\{\boldsymbol{\gamma}_{n+1}^{(i)}\}$  firstly. Since it has different expressions for  $\ell_{\infty,1}$ -norm and reweighted  $\ell_{\infty,1}$ -norm, we shall now evaluate it separately.

1)  $\ell_{\infty,1}$ -norm: Taking the expectation of (55), (56), and approximating the inequality conditions by conditions on the expectation terms, the  $m$ -th entry of  $\mathbb{E}\{\boldsymbol{\gamma}_{k,n+1}^{(1)}\}$  is given by:

$$\mathbb{E}\{\boldsymbol{\gamma}_{k,n+1}^{(1)}\}_m = \begin{cases} -\mu_k \lambda_k, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m < -\mu_k \lambda_k \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m, & \text{if } |[\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m| \leq \mu_k \lambda_k \\ \mu_k \lambda_k, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m > \mu_k \lambda_k \end{cases} \quad (76)$$

if  $\mathbb{E}\{\mathbf{w}_k^o\}_m \leq \tau_1$ , and it is defined by (77) if  $\mathbb{E}\{\mathbf{w}_k^o\}_m > \tau_1$ . To facilitate the derivation of the transient mean behavior in (76) and (77), we approximate the quantities appeared in ‘‘if’’ conditions by their corresponding expectations. It is difficult, if not impossible, to conduct theoretical analysis without this approximation. Note that this approximation is only used in the transient mean behavior analysis and have not been used in any other places. It is observed that the theoretical transient mean behavior obtained with this approximation match well with the results obtained via the Monte-Carlo simulation in Fig. 2 of Section VI-A. Vector  $\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}$  is the  $k$ -th block vector of  $\mathbb{E}\{\boldsymbol{\psi}_{n+1}\}$ , which is evaluated as follows:

$$\mathbb{E}\{\boldsymbol{\psi}_{n+1}\} = \mathbb{E}\{\tilde{\boldsymbol{\psi}}_{n+1}\} + \mathbf{w}^*, \quad (78)$$

with:

$$\mathbb{E}\{\tilde{\boldsymbol{\psi}}_{n+1}\} = \mathbf{B}\mathbb{E}\{\tilde{\mathbf{w}}_n\}. \quad (79)$$

On the other hand,  $\mathbb{E}\{\mathbf{w}_k^o\}_m$  is approximated by the maximal value of  $|[\mathbb{E}\{\boldsymbol{\psi}_{\ell,n+1}\}]_m|$  for all  $\ell \in \mathcal{N}_k^-$ .

2) *Reweighted  $\ell_{\infty,1}$ -norm*: As the distribution of  $[\mathbf{w}_k^{(t)}]_m$  is unknown, we cannot evaluate exactly the expectations involving  $[\mathbf{w}_k^{(t)}]_m$  in (57) and (58). We therefore consider the first-order Taylor series expansion as in [49]–[51]; the performance results obtained with the first-order Taylor series expansion match well with the Monte-Carlo results in Fig. 3 of Section VI-A. By expanding  $f([\mathbf{w}_k^{(t)}]_m) \triangleq \mu_k \lambda_k / [\varepsilon + |[\mathbf{w}_k^{(t)}]_m|]$  around  $\mathbb{E}\{[\mathbf{w}_k^{(t)}]_m\}$  and taking expectation, we obtain the following approximation:

$$\mathbb{E}\left\{\frac{\mu_k \lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|}\right\} \approx \frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{[\mathbf{w}_k^{(t)}]_m\}|}. \quad (80)$$

Taking the expectation of (57), (58) and using (80), as well as approximating the inequality conditions by conditions on the expectation terms, we obtain the explicit expression of  $\mathbb{E}\{\boldsymbol{\gamma}_{k,n+1}^{(2)}\}$ , with its  $m$ -th entry given by (83) or (84) when  $\mathbb{E}\{\mathbf{w}_k^o\}_m \leq \tau_2$  and  $\mathbb{E}\{\mathbf{w}_k^o\}_m > \tau_2$ , respectively, with quantities:

$$\mathbb{E}\{[\mathbf{w}_k^{(t)}]_m\} = [\mathbb{E}\{\mathbf{w}_{k,n}\}]_m \quad (81)$$

$$\mathbb{E}\{[\mathbf{c}_k]_m\} = \mathbb{E}\{[\mathbf{w}_k^o]_m\} + \frac{\mu_k \lambda_k}{\varepsilon + |[\mathbb{E}\{\mathbf{w}_{k,n}\}]_m|}. \quad (82)$$

Note that several approximations, such as the first-order Taylor series expansion [49]–[51] and approximated the inequality conditions, are only used in the transient mean behavior of the reweighted  $\ell_{\infty,1}$ -norm, not in any other places. It is observed that the theoretical transient mean behavior obtained with these approximations match well with the results obtained via the Monte-Carlo simulation in Fig. 3 of Section VI-A.

From iteration (72), we obtain the following **Theorem 1** for the mean stability of proximal multitask diffusion LMS algorithm (19).

$$[\mathbb{E}\{\gamma_{k,n+1}^{(2)}\}]_m = \begin{cases} \frac{-\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|}, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m < \frac{-\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|} \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m, & \text{if } |[\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m| \leq \frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|} \\ \frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|}, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m > \frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|} \end{cases} \quad (83)$$

$$[\mathbb{E}\{\gamma_{k,n+1}^{(2)}\}]_m = \begin{cases} -\frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|}, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \leq -\mathbb{E}\{\mathbf{c}_k\}_m \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m + \mathbb{E}\{\mathbf{w}_k^o\}_m, & \text{if } -\mathbb{E}\{\mathbf{c}_k\}_m < [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \leq -\mathbb{E}\{\mathbf{w}_k^o\}_m \\ 0, & \text{if } |[\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m| < \mathbb{E}\{\mathbf{w}_k^o\}_m \\ [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m - \mathbb{E}\{\mathbf{w}_k^o\}_m, & \text{if } \mathbb{E}\{\mathbf{w}_k^o\}_m \leq [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m < \mathbb{E}\{\mathbf{c}_k\}_m \\ \frac{\mu_k \lambda_k}{\varepsilon + |\mathbb{E}\{\mathbf{w}_k^{(t)}\}_m|}, & \text{if } [\mathbb{E}\{\boldsymbol{\psi}_{k,n+1}\}]_m \geq \mathbb{E}\{\mathbf{c}_k\}_m \end{cases} \quad (84)$$

**Theorem 1.** (*Mean stability*) Assume data model (1) and assumption **A1** hold. Then for any initial conditions, distributed networks endowed with the proximal multitask diffusion LMS algorithm (19) are stable in the mean, if step-sizes  $\mu_k$  satisfy:

$$0 < \mu_k < \frac{2}{\lambda_{\max}\{\mathbf{R}_{x,k}\}}, \quad k = 1, \dots, N, \quad (85)$$

where  $\lambda_{\max}\{\cdot\}$  denotes the maximal eigenvalue of its matrix argument. The block maximum norm of the bias can be upper bounded as:

$$\lim_{n \rightarrow \infty} \|\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}\|_{b,\infty} \leq \frac{\sqrt{L} \cdot \max_k \{\mu_k \lambda_k\}}{1 - \|\mathbf{B}\|_{b,\infty}} \quad (86)$$

$$\lim_{n \rightarrow \infty} \|\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}\|_{b,\infty} \leq \frac{1}{\varepsilon} \frac{\sqrt{L} \cdot \max_k \{\mu_k \lambda_k\}}{1 - \|\mathbf{B}\|_{b,\infty}} \quad (87)$$

for the  $\ell_{\infty,1}$ -norm and reweighted  $\ell_{\infty,1}$ -norm, respectively.

*Proof:* By iterating the RHS of (72) from time instant  $n = 0$ , and proving the convergence of the obtained series, we arrive at condition (85) for step-size to ensure the mean stability of (19). For more details, see Appendix D. ■

**Remark 1:** Equation (85) provides an upper bound for step-size  $\mu_k$  to ensure the mean stability of the distributed networks with the proximal multitask diffusion LMS algorithm (19). The upper bound is closely related to the second-order statistics  $\mathbf{R}_{x,k}$  of the input signals. Equations (86) and (87) indicate that the proximal multitask diffusion LMS algorithm (19) is biased. The upper bound of the biases are proportional to the length of the system vector  $L$ , the step-size  $\mu_k$  and the regularization parameter  $\lambda_k$ . The bias can be reduced by using a sufficiently small step-size  $\mu_k$  or regularization parameter  $\lambda_k$ . Besides, for the reweighted  $\ell_{\infty,1}$ -norm, the bias is also inversely proportional to parameter  $\varepsilon$ . In addition, the upper bound of bias for the reweighted  $\ell_{\infty,1}$ -norm has an improvement with a factor  $\frac{1}{\varepsilon}$  than that of the  $\ell_{\infty,1}$ -norm.

### B. Mean-square behavior analysis

Under **A1** and using (71), then for any semi-positive definite matrix  $\boldsymbol{\Sigma}$  of compatible dimension, the weighted mean-square behavior of  $\tilde{\mathbf{w}}_{n+1}$  evaluates as:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\Sigma}}^2\} &= \mathbb{E}\{\|\tilde{\mathbf{w}}_n\|_{\boldsymbol{\Sigma}'}^2\} + \mathbb{E}\{\|\mathbf{U}\mathbf{h}_n\|_{\boldsymbol{\Sigma}}^2\} + \mathbb{E}\{\|\gamma_{n+1}^{(i)}\|_{\boldsymbol{\Sigma}}^2\} \\ &\quad - 2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}_n^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\} - 2\mathbb{E}\{\mathbf{h}_n^\top \mathbf{U}^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\}, \end{aligned} \quad (88)$$

where  $\|\mathbf{x}\|_{\boldsymbol{\Sigma}}^2 \triangleq \mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x}$ , and

$$\boldsymbol{\Sigma}' \triangleq \mathbb{E}\{\mathbf{B}_n^\top \boldsymbol{\Sigma} \mathbf{B}_n\}. \quad (89)$$

Let  $\boldsymbol{\sigma} \triangleq \text{vec}\{\boldsymbol{\Sigma}\}$  and  $\boldsymbol{\sigma}' \triangleq \text{vec}\{\boldsymbol{\Sigma}'\}$ , where  $\text{vec}\{\cdot\}$  operator stacks the columns of its matrix argument on top of each other. Using the property of  $\text{vec}\{\cdot\}$  operator, (89) becomes:

$$\boldsymbol{\sigma}' = \mathbb{E}\{\mathbf{B}_n^\top \otimes \mathbf{B}_n^\top\} \boldsymbol{\sigma}. \quad (90)$$

Under **A2** and ignoring terms on the second-order of the maximal step-size, we have the approximation for (90):

$$\boldsymbol{\sigma}' \approx \mathbf{F}\boldsymbol{\sigma} \quad (91)$$

with  $\mathbf{F} \triangleq \mathbf{B}^\top \otimes \mathbf{B}^\top$ . Define:

$$\mathbf{H} \triangleq \mathbf{U} \text{diag}\{\sigma_{z,1}^2 \mathbf{R}_{x,1}, \dots, \sigma_{z,N}^2 \mathbf{R}_{x,N}\} \mathbf{U}^\top. \quad (92)$$

We then have:

$$\mathbb{E}\{\|\mathbf{U}\mathbf{h}_n\|_{\boldsymbol{\Sigma}}^2\} = [\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma}. \quad (93)$$

To make the analysis tractable, we adopt approximation:

$$\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}_n^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\} \approx \mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\}. \quad (94)$$

Since  $|\gamma_{k,n+1}^{(1)}| \preceq \mu_k \lambda_k \mathbb{1}_L$  and  $|\gamma_{k,n+1}^{(2)}| \preceq \frac{\mu_k \lambda_k}{\varepsilon} \mathbb{1}_L$ , we conclude that  $\gamma_{n+1}^{(i)}$  is at most of the same order as the step-size. This implies that the last term on the RHS of (88) contains higher-order powers of the step-size, and can be ignored according to assumption **A2**. Finally, by using (91), (93) and (94), expression (88) becomes:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\} &= \mathbb{E}\{\|\tilde{\mathbf{w}}_n\|_{\mathbf{F}\boldsymbol{\sigma}}^2\} + [\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma} + \mathbb{E}\{\|\gamma_{n+1}^{(i)}\|_{\boldsymbol{\sigma}}^2\} \\ &\quad - 2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\}, \end{aligned} \quad (95)$$

where we use the notations  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\Sigma}}^2\}$  and  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\}$  interchangeably.

From iteration (95), we obtain the following **Theorem 2** to ensure the mean-square stability of the proximal multitask diffusion LMS algorithm (19).

**Theorem 2.** (*Mean-square stability*) Assume data model (1) and assumptions **A1**, **A2** hold. Further assume that approximation (91) is reasonable for sufficiently small step-sizes. Then for any initial conditions, distributed networks endowed with proximal multitask diffusion LMS algorithm (19) is stable in the mean-square sense, if the step-sizes  $\mu_k$  are sufficiently small and satisfy (85).

*Proof:* By proving that the last two terms on the RHS of (95) is bounded, and iterating (95) from time instant  $n = 0$ , we obtain the condition for step-size to ensure the mean-square stability of (19). For more details, see Appendix E. ■

**Remark 2:** The weighted mean-square behavior of  $\tilde{\mathbf{w}}_{n+1}$ , that is  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\Sigma}}^2\}$ , can be decomposed as:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\Sigma}}^2\} &= \mathbb{E}\left\{\|\mathbf{w}_{n+1} - \mathbf{w}^* - \mathbb{E}\{\mathbf{w}_{n+1}\} + \mathbb{E}\{\mathbf{w}_{n+1}\}\|_{\boldsymbol{\Sigma}}^2\right\} \\ &= \underbrace{\mathbb{E}\left\{\|\tilde{\mathbf{w}}_{n+1} - \mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}\|_{\boldsymbol{\Sigma}}^2\right\}}_{\text{Variance term}} + \underbrace{\|\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}\|_{\boldsymbol{\Sigma}}^2}_{\text{Bias term}} \end{aligned} \quad (96)$$

By substituting  $\tilde{\mathbf{w}}_{n+1}$  of (71) and  $\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}$  of (72) into (96), we obtain an equivalent form of (95). Relation (96) is called the bias-variance decomposition. From (96), we observe that the stability of  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\Sigma}}^2\}$  in **Theorem 2** ensures the stabilities of both the Variance term and the Bias term  $\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}$  of **Theorem 1**.

## VI. SIMULATION RESULTS

In this section, we present simulation results to validate the effectiveness of the algorithm. With the exception of the simulation results presented in Section VI-A, used to validate theoretical results in the mean behavior analysis and obtained by averaging over 500 independent Monte-Carlo runs, all other simulated curves were obtained by averaging over 100 independent Monte-Carlo runs.

### A. Theoretical validation

We considered a connected network consisting of 16 nodes and 36 edges. The number of edges at each node was between 2 and 7, without taking into account the possible self-loops connecting each node to itself. Other characteristics of the network are listed in Table I. Each regressor  $\mathbf{x}_{k,n}$  was generated from a zero-mean Gaussian distribution with covariance matrix  $\mathbf{R}_{x,k} = \sigma_{x,k}^2 \mathbf{I}_{30}$ . Each additive noise  $z_{k,n}$  was generated from a zero-mean Gaussian distribution with variance  $\sigma_{z,k}^2$ . Variances  $\sigma_{x,k}^2$  and  $\sigma_{z,k}^2$  at each node were generated randomly from a Gaussian distribution as shown in Fig. 1. Note that these variances settings are the same as in [12], [17].

To validate the theoretical results reported in the mean behavior analysis, we considered a stationary system identification scenario. The unknown system coefficients  $\mathbf{w}_k^*$  were generated such that the entire network has a jointly sparse structure with sparsity degree of 10/30. Each nonzero element of  $\mathbf{w}_k^*$  was generated independently from a standard Gaussian distribution. The simulation results presented in this section were obtained by averaging over 500 independent runs.

TABLE I

CHARACTERISTICS OF THE NETWORK USED FOR MODEL VALIDATION.  $\mathbf{L}$  IS THE LAPLACIAN MATRIX ASSOCIATED WITH THE GRAPH,  $\lambda_2(\mathbf{L})$  IS THE ALGEBRAIC CONNECTIVITY [52] OF THE GRAPH, SIZE IS THE NUMBER OF NODES, DENSITY IS THE NUMBER OF NON-ZERO ENTRIES OF THE ADJACENCY MATRIX, AND DIAMETER IS THE MAXIMUM DISTANCE BETWEEN ANY TWO NODES [53].

Network	Size	Density	$\lambda_2(\mathbf{L})$	Diameter
Net	16	34%	0.9983	4

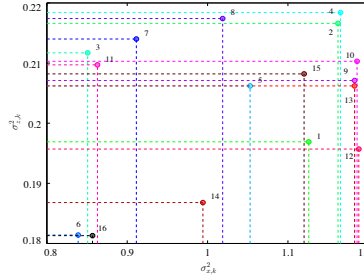
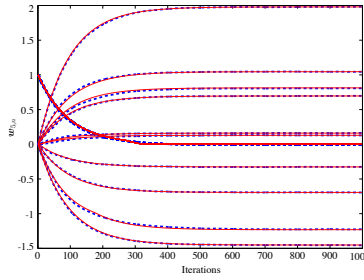
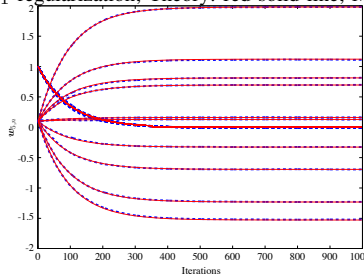
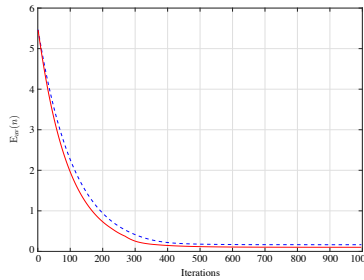


Fig. 1. Agent input and noise variances.

Fig. 2. Validation of the mean behavior model for  $\ell_{\infty,1}$ -regularization; Theory: red solid line; Monte-Carlo simulations: blue dashed line.Fig. 3. Validation of the mean behavior model in the case of the reweighted  $\ell_{\infty,1}$ -regularization; Theory: red solid line; Monte-Carlo: blue dashed line.Fig. 4. The average error  $E_{\text{av}}(n)$  of  $\ell_{\infty,1}$ -regularization; Theory: red solid line; Monte-Carlo simulations: blue dashed line.

The results are illustrated in Fig. 2 and Fig. 3 for the proximal multitask LMS with  $\ell_{\infty,1}$ -regularization and reweighted  $\ell_{\infty,1}$ -regularization, respectively. The good match between the theoretical results and the Monte-Carlo curves illustrates the accuracy of theoretical results in the mean behavior analysis.

We also considered the average error over the entire network defined by:

$$E_{\text{av}}(n) \triangleq \frac{1}{N} \sum_{k=1}^N \|\mathbf{w}_{k,n} - \mathbf{w}_k^*\| \quad (97)$$

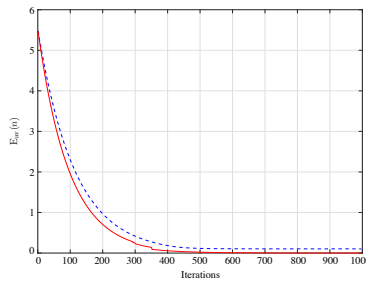


Fig. 5. The average error  $E_{\text{av}}(n)$  of the reweighted  $\ell_{\infty,1}$ -regularization; Theory: red solid line; Monte-Carlo simulations: blue dashed line.

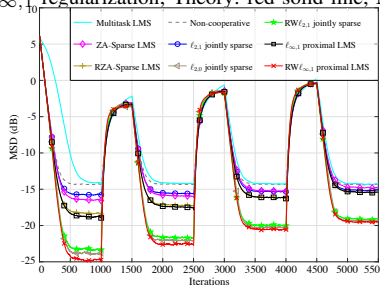


Fig. 6. Comparison of the proposed algorithms with several state-of-the-art algorithms for white inputs.

Though there is a small bias between theoretical and Monte-Carlo curves, the results reported in Fig. 4 and Fig. 5 confirm the accuracy of the theoretical models.

## B. Numerical Simulations

1) *Comparison with existing algorithms:* We firstly considered a non-stationary jointly sparse system identification scenario with  $w_k^*$  varying over time. Each nonzero entry of  $w_k^*$  was generated independently from a standard Gaussian distribution. The evolution of  $w_k^*$  was divided into four stationary stages and three transient stages. During stationary stages, sparse vectors  $w_k^*$  were set to sparsity degree of 3/30, 5/30, 8/30 and 10/30, respectively. The transient stages were designed by using linear interpolation over 500 time instants. The regressors  $x_{k,n}$  were generated as those in Section VI-A for white inputs, while they were generated according to a zero-mean Gaussian distribution with covariance matrix  $R_{x,k} = \sigma_{x,k}^2 R^\dagger$  for colored inputs, where  $R^\dagger$  is an  $30 \times 30$  Hermite matrix with eigenvalue spread  $\lambda_{\max}\{R^\dagger\}/\lambda_{\min}\{R^\dagger\} = 21$ , and symbol  $\lambda_{\min}\{\cdot\}$  denotes the minimal eigenvalue of its matrix argument. For comparison purpose, non-cooperative diffusion LMS algorithm, non-cooperative sparse diffusion LMS [6] with zero-attracting (ZA) regularizer and reweighted zero-attracting (RZA) regularizer, multitask diffusion LMS with adaptive combiner [17] and jointly sparse multitask diffusion LMS [24] with  $\ell_{2,1}$ -regularization, reweighted  $\ell_{2,1}$ -regularization (RW $\ell_{2,1}$ ) and  $\ell_{2,0}$ -regularization were taken into consideration. We adopted a uniform step-size 0.01 for all algorithms. For these algorithms, we set their parameters so that they reach their best performance. To enable reproducible research, details about parameters used by these algorithms can be found in Table III of Appendix F.

The results are illustrated in Fig. 6 for white inputs. We observe that multitask LMS with adaptive combiner is the worst one among all competing algorithms, since it utilizes similarities between neighboring nodes to improve estimation accuracy. This does not necessarily exist in jointly sparse scenarios and may deteriorate the MSD performance. Since jointly sparse system can be regarded as a special case of general sparse systems, by using additional information about system sparsity, sparse diffusion LMS with ZA regularizer and RZA regularizer have better performance than the non-cooperative LMS. Similarly, all jointly sparse multitask algorithms considered in this comparative experiment perform better than the non-cooperative LMS. Observe that sparse diffusion LMS algorithms can perform slightly better than jointly sparse diffusion LMS algorithms as the parameter vectors to estimate become sparser, while they do not perform as well when these parameter vectors are less sparse. These findings illustrate the interest of exploiting joint sparsity as prior information. On the one hand, the proposed proximal multitask LMS with reweighted  $\ell_{\infty,1}$ -regularization performs better than all other algorithms as evidenced by its lowest steady-state MSD. On the other hand, the proposed proximal multitask LMS with  $\ell_{\infty,1}$ -regularization has similar performance to the sparse diffusion LMS with RZA regularizer when there are more zeros in the jointly sparse system to estimate, but the former performs better than the latter as the number of nonzero entries increases.

The results are illustrated in Fig. 7 for colored input signals. It can be observed that the convergence rate of all algorithms is slower than in the case of white inputs, and some algorithms have poorer performance than the non-cooperative algorithm when the parameter vectors to estimate are less sparse. Besides some conclusions that have been drawn for white inputs, we observe that the proximal multitask LMS algorithm with reweighted  $\ell_{\infty,1}$ -regularization still has the best performance.

Then we examined the performance of all algorithms as a function of the sparsity degree of the jointly sparse system to estimate. We considered the steady-state MSD as the measure of performance in this experiment. The results are illustrated in



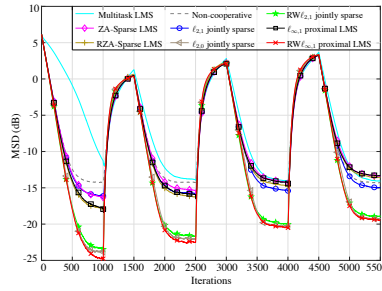


Fig. 7. Comparison of the proposed algorithms with several state-of-the-art algorithms for colored inputs.

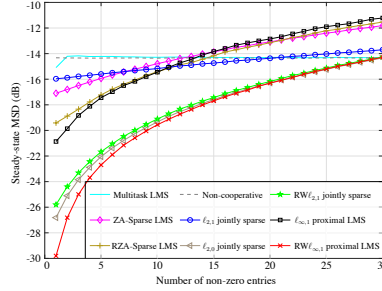


Fig. 8. Steady-state MSDs of all algorithms as a function of the sparsity degree.

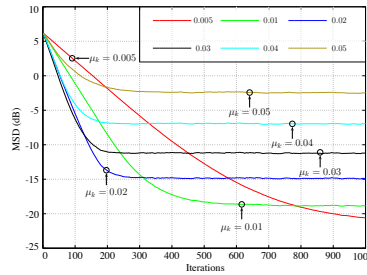


Fig. 9. Effect of the step-size  $\mu_k$  on the MSD of the proximal LMS algorithm with  $\ell_{\infty,1}$ -regularization.

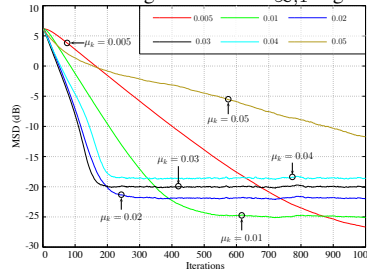


Fig. 10. Effect of the step-size  $\mu_k$  on the MSD of the proximal LMS algorithm with reweighted  $\ell_{\infty,1}$ -regularization.

Fig. 8. Besides some conclusions that have been drawn before, we observe that some of the multitask algorithms with sparse and jointly sparse regularizers may have poorer performance than the non-cooperative LMS when the number of nonzero entries increases. Three algorithms, including the proximal multitask LMS with reweighted  $\ell_{\infty,1}$ -regularization, uniformly show better performance than all other algorithms for all sparsity degrees, and have similar performance to the non-cooperative LMS for totally non-sparse systems.

2) *Effects of parameters setting:* To examine the effects of parameters setting, including the step-size  $\mu_k$ , the regularization parameter  $\lambda_k$ , the threshold values  $\tau_1, \tau_2$  and the parameter  $\varepsilon$ , we considered a stationary system identification problem with sparsity degree of  $3/30$ . Each parameter was set to a same value for all nodes in the network. We examined the influence of one selected parameter at a time, setting all other parameters to fixed values, in order to facilitate comparison.

The effects of the step-size  $\mu_k$  are illustrated in Fig. 9 and Fig. 10 for the proximal LMS with  $\ell_{\infty,1}$ -regularization and reweighted  $\ell_{\infty,1}$ -regularization, respectively. Observe that  $\mu_k$  allows to control the trade-off between convergence rate and steady-state performance. A larger step-size results in a faster convergence rate at the cost of a larger steady-state MSD. A small step-size results in a more accurate estimation at the cost of a slower convergence speed.

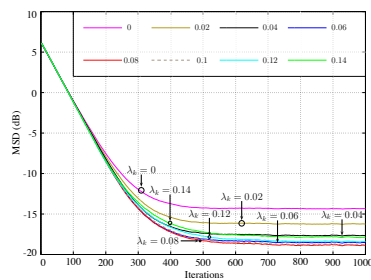


Fig. 11. Effect of the regularization parameter  $\lambda_k$  on the MSD of the proximal LMS algorithm with  $\ell_{\infty,1}$ -regularization.

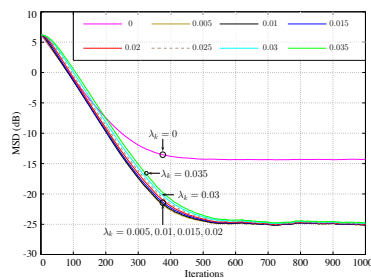


Fig. 12. Effect of the regularization parameter  $\lambda_k$  on the MSD of the proximal LMS algorithm with reweighted  $\ell_{\infty,1}$ -regularization.

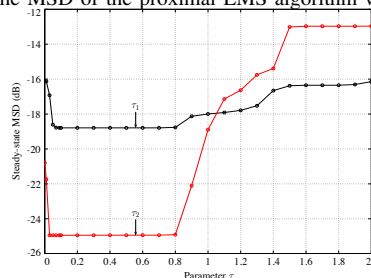


Fig. 13. Effect of parameters  $\tau_1$  and  $\tau_2$  on the steady-state MSD of the proximal LMS algorithm with  $\ell_{\infty,1}$ -regularization (black) and reweighted  $\ell_{\infty,1}$ -regularization (red), respectively.

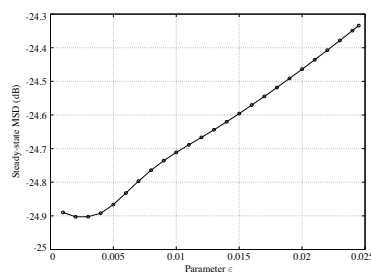


Fig. 14. Effect of parameter  $\epsilon$  on the steady-state MSD of the proximal LMS algorithm with reweighted  $\ell_{\infty,1}$ -regularization.

The effects of the regularization parameter  $\lambda_k$  are illustrated in Fig. 11 and Fig. 12. We observe that increasing  $\lambda_k$  improves the convergence speed and the MSD at steady-state at first, and then degrades them. In this experiment, the critical values were 0.08 and 0.01 for the  $\ell_{\infty,1}$  regularizer and the reweighted  $\ell_{\infty,1}$  regularizer, respectively.

The effects of parameters  $\tau_1$  and  $\tau_2$  are shown in Fig. 13. Similar behaviors can be observed on the MSD at steady-state of both algorithms with respect to  $\tau_1$  and  $\tau_2$ . The best values were obtained over interval  $[0.05, 0.8]$  for both regularizers. First, these results show the need for introducing these two parameters. Second, they show that these two parameters can be appropriately selected over a large interval. Consider the  $\ell_{\infty,1}$ -norm regularizer. When  $[\mathbf{w}_k^o]_m \leq \tau_1$ , observe that the approximate proximal operator is given by (33) and corresponds to the  $\ell_1$ -norm regularizer, which enjoys the following properties: On the one hand, it shrinks the estimates of zero-valued parameters to zero when  $[\mathbf{w}_k^o]_m \leq \tau_1$ ; On the other hand, for small nonzero-valued entries, though introducing a bias when shrunk to zero, it lowers the estimates variance and results in a satisfying performance. Since jointly sparse systems are a special case of general sparse ones, for which the  $\ell_1$ -norm regularizer is prescribed, it is expected that our algorithm works well for a large range of  $\tau_1$  values. The same reasoning applies to the reweighted  $\ell_{\infty,1}$ -norm regularizer and  $\tau_2$ .

The effects of parameter  $\epsilon$  in the reweighted  $\ell_{\infty,1}$  regularizer are illustrated in Fig. 14. We set both parameters  $\mu_k$  and  $\lambda_k$  to 0.01. We observe that the best value of  $\epsilon$  resulting in the lowest MSD at steady-state was 0.002, but its effect on the

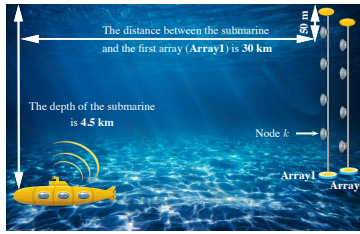


Fig. 15. Experimental setup for the practical application.

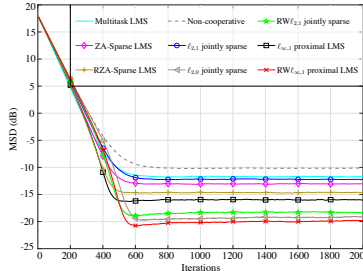


Fig. 16. Comparison of the proposed algorithms with several state-of-the-art algorithms for practical application.

steady-state MSD is weaker than the other parameters.

### C. Practical application

We shall now validate the proposed algorithm in a practical application. This application follows the experimental setup described in [24]. In demodulation and decoding in underwater and wireless communication networks, a standalone low-cost sensor may not be able to decode or demodulate the source signal reliably due to the extra low signal-to-noise ratio (SNR) condition. It is better to adopt multiple sensors to collaboratively recover the information coming from a same source [54], [55] by combining and exchanging information within the sensor network. Practically, it is necessary to estimate the channel impulse responses between the source and each sensor before decoding. On the one hand, it has been shown that real-world underwater channels [27], and wireless communication channel [28], are sparse with large delay spread. On the other hand, the channel supports for neighboring antennas or nodes are approximately the same [28]. Indeed the times of arrival for closely spaced nodes and antennas are quite close, though the tap weights are different [24]. We shall now check that taking the jointly sparse property into consideration can improve the network performance.

Consider the problem of identifying underwater acoustic channels. We set the scenario presented in Fig. 15, where we have one source (the submarine) and two receiving sensor arrays. Both arrays are linear arrays and consist of 10 nodes with the same distance (3.75 m) between each of them. We assume that each node has the ability to communicate and process data. These two arrays are positioned in parallel with interspace 50 m, which gives us the distributed network with 20 nodes. Within each array, the nodes are connected one by one in a chain, and the corresponding pairs of nodes between two arrays are also connected. The underwater acoustic channels to estimate were generated via the BELLHOP model [56], which has been developed for predicting acoustic pressure fields in ocean environments. We used a white signal as input, and the additive noise at each node was white Gaussian. SNR conditions are listed in Table II.

TABLE II  
SNR LEVEL IN DECIBEL (dB) FOR THE PRACTICAL APPLICATION. SINCE IT VARIES ACCORDING TO NODES, WE ENUMERATE THE MAXIMUM, MINIMUM AND MEAN VALUES.

SNR Level	Maximum	Minimum	Mean
SNR	7.89	6.19	7.07

We compared our algorithms with the algorithms considered in Section VI-B1. Parameters used by all them are listed in Table III. The results are depicted in Fig. 16. As can be observed, the proposed proximal multitask diffusion LMS algorithm with reweighted  $\ell_{\infty,1}$  regularizer achieved the best performance in terms of the steady-state MSD.

## VII. CONCLUSION

We considered the problem of estimating a set of parameter vectors in a distributed manner, where the local solutions have the same sparse support. We devised a proximal diffusion algorithm with (reweighted)  $\ell_{\infty,1}$ -norm regularization, with closed-form expressions for the regularizers. We conducted theoretical analyses of the algorithms behavior in the mean and mean-square sense. Simulation results illustrated the effectiveness of the proposed algorithms, as well as the accuracy of theoretical results. Integrating weighted network connection information to enhance this jointly-sparse estimation will be considered in future work.

APPENDIX A  
INTERPRETATIONS ABOUT THE (REWEIGHTED)  $\ell_{\infty,1}$ -NORM

The  $\ell_{\infty,1}$ -norm of matrix  $\mathbf{W}_k$  is defined as  $g_1(\mathbf{w}_k)$  of (11). We focus on the interpretation of  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty}$ . From (8) we have that  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty} = \max\{|\mathbf{w}_k]_m|, |\mathbf{w}_\ell^*]_m|\}$  with  $\ell \in \mathcal{N}_k^-$ . Denote the maximal value of  $|\mathbf{w}_\ell^*]_m|$  for  $\ell \in \mathcal{N}_k^-$  as  $[\mathbf{w}_k^o]_m$ . Then  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty}$  writes to  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty} = \max\{|\mathbf{w}_k]_m|, [\mathbf{w}_k^o]_m\}$ . Now we consider the following three cases:

- Case a: The  $m$ -th row  $\bar{\mathbf{w}}_{k,m}$  of matrix  $\mathbf{W}_k$  corresponds to the zero-valued entries. In this case,  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty} = \max\{|\mathbf{w}_k]_m|, 0\} = \max\{|\mathbf{w}_k]_m|\}$ . Minimizing  $g_1(\mathbf{w}_k)$  over  $[\mathbf{w}_k]_m$  attracts the  $m$ -th entry  $[\mathbf{w}_k]_m$  to zero.
- Case b: The  $m$ -th row  $\bar{\mathbf{w}}_{k,m}$  corresponds to the non-zero valued entries, and  $|\mathbf{w}_k]_m| \leq [\mathbf{w}_k^o]_m$ . In this case,  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty} = \max\{|\mathbf{w}_k]_m|, [\mathbf{w}_k^o]_m\} = [\mathbf{w}_k^o]_m$ . Minimizing  $g_1(\mathbf{w}_k)$  over  $[\mathbf{w}_k]_m$  will not penalize  $[\mathbf{w}_k]_m$ .
- Case c: The  $m$ -th row  $\bar{\mathbf{w}}_{k,m}$  corresponds to the non-zero valued entries, and  $|\mathbf{w}_k]_m| > [\mathbf{w}_k^o]_m$ . In this case,  $\|\bar{\mathbf{w}}_{k,m}\|_{\infty} = \max\{|\mathbf{w}_k]_m|, [\mathbf{w}_k^o]_m\} = |\mathbf{w}_k]_m|$ . Minimizing  $g_1(\mathbf{w}_k)$  over  $[\mathbf{w}_k]_m$  will penalize  $[\mathbf{w}_k]_m$  until  $|\mathbf{w}_k]_m| \leq [\mathbf{w}_k^o]_m$ , which then becomes Case b.

Thus, the  $\ell_{\infty,1}$ -norm promotes the similarity between  $|\mathbf{w}_k]_m|$  and  $[\mathbf{w}_k^o]_m$  for all  $m$ . Similar interpretations can be obtained for the reweighted  $\ell_{\infty,1}$ -norm of (12).

APPENDIX B  
DERIVATION OF (36)

The initial problem (24) leads to (25) and (27) in Case 1 and Case 2, respectively. For ease of presentation, we define:

$$\bar{J}_1([\mathbf{w}_k]_m) \triangleq [\mathbf{w}_k^o]_m + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2 \quad (98)$$

$$\bar{J}_2([\mathbf{w}_k]_m) \triangleq |\mathbf{w}_k]_m| + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k]_m - [\psi_{k,n+1}]_m)^2. \quad (99)$$

We select solution (26) or (35) depending on the value taken by costs (98) and (99). To save space, we partially present the derivation. The rest of the derivation can be obtained by following the same routine.

When  $[\psi_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m + \mu_k\lambda_k$ , substituting  $\hat{w} = [\mathbf{w}_k^o]_m$  of Case 1 into (98), we obtain:

$$\bar{J}_1([\mathbf{w}_k^o]_m) = [\mathbf{w}_k^o]_m + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k^o]_m - [\psi_{k,n+1}]_m)^2. \quad (100)$$

Substituting  $\hat{w} = [\psi_{k,n+1}]_m - \mu_k\lambda_k$  of Case 2 into (99), we obtain:

$$\begin{aligned} \bar{J}_2(\hat{w}) &= |\hat{w}| + \frac{1}{2\mu_k\lambda_k} (\hat{w} - [\psi_{k,n+1}]_m)^2 \\ &\leq \bar{J}_2([\mathbf{w}_k^o]_m) = \bar{J}_1([\mathbf{w}_k^o]_m) \end{aligned} \quad (101)$$

since  $\bar{J}_2(\hat{w})$  with  $\hat{w} = [\psi_{k,n+1}]_m - \mu_k\lambda_k$  is the minimal cost. Thus the proximal operator is given by:

$$\hat{w} = [\psi_{k,n+1}]_m - \mu_k\lambda_k \quad (102)$$

for  $[\psi_{k,n+1}]_m \geq [\mathbf{w}_k^o]_m + \mu_k\lambda_k$ .

When  $[\mathbf{w}_k^o]_m \leq [\psi_{k,n+1}]_m < [\mathbf{w}_k^o]_m + \mu_k\lambda_k$ , since  $\hat{w}$  of both Case 1 and Case 2 is given by:

$$\hat{w} = [\mathbf{w}_k^o]_m, \quad (103)$$

we arrive at proximal operator (103) directly.

When  $0 < [\psi_{k,n+1}]_m < [\mathbf{w}_k^o]_m$ , substituting  $\hat{w} = [\psi_{k,n+1}]_m$  of Case 1 into (98), we obtain:

$$\bar{J}_1([\psi_{k,n+1}]_m) = [\mathbf{w}_k^o]_m. \quad (104)$$

Substituting  $\hat{w} = [\mathbf{w}_k^o]_m$  of Case 2 into (99), we obtain:

$$\begin{aligned} \bar{J}_2([\mathbf{w}_k^o]_m) &= [\mathbf{w}_k^o]_m + \frac{1}{2\mu_k\lambda_k} ([\mathbf{w}_k^o]_m - [\psi_{k,n+1}]_m)^2 \\ &\geq \bar{J}_1([\psi_{k,n+1}]_m). \end{aligned} \quad (105)$$

Thus the proximal operator is given by:

$$\hat{w} = [\psi_{k,n+1}]_m \quad (106)$$

for  $0 < [\psi_{k,n+1}]_m < [\mathbf{w}_k^o]_m$ .

By following the same routine, we obtain the proximal operator for  $[\psi_{k,n+1}]_m \leq 0$ . Finally, by combining all these results, we arrive at the expression of proximal operator (36) when  $[\mathbf{w}_k^o]_m > 0$ .

APPENDIX C  
DERIVATION OF (53)

Define:

$$\begin{aligned} \bar{J}_3([\mathbf{w}_k]_m) &\triangleq \log\left(1 + \frac{[\mathbf{w}_k^o]_m}{\varepsilon}\right) \\ &\quad + \frac{1}{2\mu_k\lambda_k}([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2 \end{aligned} \quad (107)$$

$$\begin{aligned} \bar{J}_4([\mathbf{w}_k]_m) &\triangleq \log\left(1 + \frac{|[\mathbf{w}_k]_m|}{\varepsilon}\right) \\ &\quad + \frac{1}{2\mu_k\lambda_k}([\mathbf{w}_k]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2. \end{aligned} \quad (108)$$

The initial problem (37) leads to  $\bar{J}_3([\mathbf{w}_k]_m)$  and  $\bar{J}_4([\mathbf{w}_k]_m)$  in Case 1 and Case 2, respectively. We select solution (39), (49) or (51) by comparing costs (107) and (108) as that in Section IV-A. We partially present the derivation to save space.

When  $0 < [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{w}_k^o]_m$ , substituting  $\hat{w} = [\boldsymbol{\psi}_{k,n+1}]_m$  of Case 1 into (107), we obtain:

$$\bar{J}_3([\boldsymbol{\psi}_{k,n+1}]_m) = \log\left(1 + \frac{[\mathbf{w}_k^o]_m}{\varepsilon}\right). \quad (109)$$

Substituting  $\hat{w} = [\mathbf{w}_k^o]_m$  of Case 2 into (108), we obtain:

$$\begin{aligned} \bar{J}_4([\mathbf{w}_k^o]_m) &= \log\left(1 + \frac{[\mathbf{w}_k^o]_m}{\varepsilon}\right) \\ &\quad + \frac{1}{2\mu_k\lambda_k}([\mathbf{w}_k^o]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2 \\ &\geq \bar{J}_3([\boldsymbol{\psi}_{k,n+1}]_m). \end{aligned} \quad (110)$$

This means that the proximal operator is given by:

$$\hat{w} = [\boldsymbol{\psi}_{k,n+1}]_m \quad (111)$$

for  $0 < [\boldsymbol{\psi}_{k,n+1}]_m < [\mathbf{w}_k^o]_m$ .

When  $[\mathbf{w}_k^o]_m \leq [\boldsymbol{\psi}_{k,n+1}]_m < [c_k]_m$ , since  $\hat{w}$  of both Case 1 and Case 2 is given by:

$$\hat{w} = [\mathbf{w}_k^o]_m, \quad (112)$$

we arrive at proximal operator (112) directly.

When  $[\boldsymbol{\psi}_{k,n+1}]_m \geq [c_k]_m$ , substituting  $\hat{w} = [\mathbf{w}_k^o]_m$  of Case 1 into (107), we obtain:

$$\begin{aligned} \bar{J}_3([\mathbf{w}_k^o]_m) &= \log\left(1 + \frac{[\mathbf{w}_k^o]_m}{\varepsilon}\right) \\ &\quad + \frac{1}{2\mu_k\lambda_k}([\mathbf{w}_k^o]_m - [\boldsymbol{\psi}_{k,n+1}]_m)^2 \\ &= \bar{J}_4([\mathbf{w}_k^o]_m). \end{aligned} \quad (113)$$

Further, since  $\hat{w} = [\boldsymbol{\psi}_{k,n+1}]_m - \mu_k\lambda_k/(\varepsilon + |[\mathbf{w}_k^{(t)}]_m|)$  of Case 2 is an approximation of the minimizer  $\hat{w}^o$  of (41), and according to the relation  $\bar{J}_4(\hat{w}^o) \leq \bar{J}_4([\mathbf{w}_k^o]_m)$ , we define the proximal operator as:

$$\hat{w} = [\boldsymbol{\psi}_{k,n+1}]_m - \frac{\mu_k\lambda_k}{\varepsilon + |[\mathbf{w}_k^{(t)}]_m|} \quad (114)$$

for  $[\boldsymbol{\psi}_{k,n+1}]_m \geq [c_k]_m$ .

By following the same routine, we obtain the proximal operator for  $[\boldsymbol{\psi}_{k,n+1}]_m \leq 0$ . Finally, by combining all these results, we arrive at expression (53) when  $[\mathbf{w}_k^o]_m > 0$ .

APPENDIX D  
PROOF OF THEOREM 1

Iterating (72) from  $n = 0$ , we obtain:

$$\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\} = \mathbf{B}^{n+1}\mathbb{E}\{\tilde{\mathbf{w}}_0\} - \sum_{j=0}^n \mathbf{B}^j \mathbb{E}\{\boldsymbol{\gamma}_{n+1-j}^{(i)}\} \quad (115)$$

where  $\mathbb{E}\{\tilde{\mathbf{w}}_0\}$  is the initial condition. The convergence of (115) requires that both terms on the RHS to be convergent. For the first term, it requires that spectral radius  $\rho(\mathbf{B}) < 1$  to ensure the convergence. For the second term, it is sufficient to prove that  $\sum_{j=0}^n [\mathbf{B}^j \mathbb{E}\{\gamma_{n+1-j}^{(i)}\}]_m$  is convergent for  $m = 1, \dots, NL$ . A series is absolutely convergent if each term is bounded by a term of an absolutely convergent series [6], [14]. Define  $s_m \triangleq [\mathbf{B}^j \mathbb{E}\{\gamma_{n+1-j}^{(i)}\}]_m$ . Since the block maximum norm of a block vector is larger than or equal to the largest absolute value of its entry, we have:

$$\begin{aligned} |s_m| &\leq \|\mathbf{B}^j \mathbb{E}\{\gamma_{n+1-j}^{(i)}\}\|_{b,\infty} \\ &\leq \|\mathbf{B}\|_{b,\infty}^j \cdot \|\mathbb{E}\{\gamma_{n+1-j}^{(i)}\}\|_{b,\infty} \\ &\leq [\rho(\mathbf{B})]^j \cdot \gamma_{\max}^{(i)}, \end{aligned} \quad (116)$$

where  $\|\cdot\|_{b,\infty}$  is the block maximum norm [7]. The quantity  $\|\mathbb{E}\{\gamma_{n+1-j}^{(i)}\}\|_{b,\infty}$  is finite for all  $j$  and  $n$ , and bounded by some constant  $\gamma_{\max}^{(i)}$ . Actually, from (76), (77), (83), (84) and following the routine for the boundness of  $\gamma_{k,n+1}^{(i)}$ , we obtain:

$$|\mathbb{E}\{\gamma_{k,n+1-j}^{(1)}\}| \preceq \mu_k \lambda_k \mathbb{1}_L \quad (117)$$

$$|\mathbb{E}\{\gamma_{k,n+1-j}^{(2)}\}| \preceq b_k \mathbb{1}_L \quad (118)$$

where the quantity  $b_k$  has been defined in (61).

Thus, we conclude that  $\|\mathbb{E}\{\gamma_{n+1-j}^{(i)}\}\|_{b,\infty}$  is bounded, with

$$\|\mathbb{E}\{\gamma_{n+1-j}^{(1)}\}\|_{b,\infty} \leq \gamma_{\max}^{(1)} \triangleq \sqrt{L} \cdot \max_k \{\mu_k \lambda_k\} \quad (119)$$

$$\|\mathbb{E}\{\gamma_{n+1-j}^{(2)}\}\|_{b,\infty} \leq \gamma_{\max}^{(2)} \triangleq \frac{\sqrt{L}}{\varepsilon} \cdot \max_k \{\mu_k \lambda_k\}. \quad (120)$$

Consequently, convergence of (116) is ensured by condition  $\rho(\mathbf{B}) < 1$ . Step-sizes  $\mu_k$  satisfying (85) ensure the mean stability of the network.

With step-sizes  $\mu_k$  satisfying (85) to ensure  $\rho(\mathbf{B}) < 1$ , the block maximum norm of the bias can be bounded as  $n \rightarrow \infty$ :

$$\begin{aligned} \lim_{n \rightarrow \infty} \|\mathbb{E}\{\tilde{\mathbf{w}}_{n+1}\}\|_{b,\infty} &\leq \lim_{n \rightarrow \infty} \sum_{j=0}^n \|\mathbf{B}^j\|_{b,\infty} \|\mathbb{E}\{\gamma_{n+1-j}^{(i)}\}\|_{b,\infty} \\ &\leq \frac{\gamma_{\max}^{(i)}}{1 - \|\mathbf{B}\|_{b,\infty}}. \end{aligned} \quad (121)$$

Substituting (119) and (120) into (121), we arrive at (86) and (87), respectively.

## APPENDIX E PROOF OF THEOREM 2

Since  $\Sigma$  is a positive semi-definite matrix, and vector  $\gamma_{n+1}^{(i)}$  is uniformly bounded for all time instant  $n$  and  $i = 1$  or  $2$ , we have:

$$0 \leq \mathbb{E}\{\|\gamma_{n+1}^{(i)}\|_{\Sigma}^2\} \leq \kappa_1^{(i)} \quad (122)$$

for all  $n$ , where  $\kappa_1^{(i)}$  is a positive constant depending on the jointly sparse regularizer. Since  $\gamma_{n+1}^{(i)}$  is uniformly bounded, vector  $\mathbf{B}^\top \Sigma \gamma_{n+1}^{(i)}$  is also bounded for all  $n$ . Denote the bound of the largest component of  $2\mathbf{B}^\top \Sigma \gamma_{n+1}^{(i)}$  in absolute-value sense as  $\tau_{\max}^{(i)}$  for all  $n$ . We obtain:

$$\begin{aligned} |2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \Sigma \gamma_{n+1}^{(i)}\}| &\leq \tau_{\max}^{(i)} \sum_{k=1}^N \sum_{m=1}^L |\mathbb{E}\{\tilde{\mathbf{w}}_{k,n}\}_m| \\ &= \tau_{\max}^{(i)} \cdot \|\mathbb{E}\{\tilde{\mathbf{w}}_n\}\|_1. \end{aligned} \quad (123)$$

Given step-sizes  $\mu_k$  satisfying condition (85) to ensure stability in the mean sense, we conclude that quantity  $\|\mathbb{E}\{\tilde{\mathbf{w}}_n\}\|_1$  is upper bounded by some constant  $\kappa_2^{(i)}$  for all  $n$ . Consequently, (123) becomes:

$$|2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \Sigma \gamma_{n+1}^{(i)}\}| \leq \tau_{\max}^{(i)} \cdot \kappa_2^{(i)}. \quad (124)$$

Define

$$\phi(\tilde{\mathbf{w}}_n, \gamma_{n+1}^{(i)}) \triangleq \mathbb{E}\{\|\gamma_{n+1}^{(i)}\|_{\sigma}^2\} - 2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \Sigma \gamma_{n+1}^{(i)}\}. \quad (125)$$

Using (122)–(125), we have:

$$\begin{aligned} |\phi(\tilde{\mathbf{w}}_n, \gamma_{n+1}^{(i)})| &\leq \mathbb{E}\{\|\gamma_{n+1}^{(i)}\|_{\boldsymbol{\sigma}}^2\} + |2\mathbb{E}\{\tilde{\mathbf{w}}_n^\top \mathbf{B}^\top \boldsymbol{\Sigma} \gamma_{n+1}^{(i)}\}| \\ &\leq \kappa_1^{(i)} + \tau_{\max}^{(i)} \cdot \kappa_2^{(i)} \end{aligned} \quad (126)$$

for all  $n$ . Given a weighting matrix  $\boldsymbol{\Sigma}$ , the positive constant  $\kappa_3^{(i)} \triangleq \kappa_1^{(i)} + \tau_{\max}^{(i)} \cdot \kappa_2^{(i)}$  can be written as a scaled multiple of  $[\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma}$  as:

$$\kappa_3^{(i)} = p \cdot [\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma} \quad (127)$$

with  $p \geq 0$  [6]. Using (95) and (125)–(127), we obtain an upper bound of  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\}$  as:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\} &\leq \mathbb{E}\{\|\tilde{\mathbf{w}}_n\|_{\mathbf{F}\boldsymbol{\sigma}}^2\} + [\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma} + |\phi(\tilde{\mathbf{w}}_n, \gamma_{n+1}^{(i)})| \\ &\leq \mathbb{E}\{\|\tilde{\mathbf{w}}_n\|_{\mathbf{F}\boldsymbol{\sigma}}^2\} + (1+p) \cdot [\text{vec}\{\mathbf{H}\}]^\top \boldsymbol{\sigma}. \end{aligned} \quad (128)$$

Iterating (128) from  $n = 0$ , we obtain:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\} &\leq \mathbb{E}\{\|\tilde{\mathbf{w}}_0\|_{\mathbf{F}^{n+1}\boldsymbol{\sigma}}^2\} \\ &\quad + (1+p)[\text{vec}\{\mathbf{H}\}]^\top \sum_{j=0}^n \mathbf{F}^j \boldsymbol{\sigma} \end{aligned} \quad (129)$$

where  $\mathbb{E}\{\|\tilde{\mathbf{w}}_0\|_{\boldsymbol{\sigma}}^2\}$  is the initial condition. The stability of  $\mathbb{E}\{\|\tilde{\mathbf{w}}_{n+1}\|_{\boldsymbol{\sigma}}^2\}$  requires the convergence of terms on the RHS of (129), which is ensured by condition  $\rho(\mathbf{F}) < 1$ . Since  $\mathbf{F} = \mathbf{B}^\top \otimes \mathbf{B}^\top$ , it is enough to select sufficiently small step-sizes  $\mu_k$  satisfying (85) to ensure  $\rho(\mathbf{F}) < 1$ .

#### APPENDIX F PARAMETERS USED BY ALGORITHMS IN SECTION VI-B1

Table III provides all parameters used by the algorithms in Section VI-B1. Notations used for these parameters are the same as those in the original references. Some pairs of columns, standing for different parameters, are merged into a single column for compactness. The corresponding symbols can be distinguished by the symbol “|”. We used a uniform step-size 0.01 for all algorithms.

TABLE III  
PARAMETERS USED BY ALL ALGORITHMS FOR SIMULATION.

Algorithms	$\lambda_k$	$\gamma$	$\varepsilon$	$\eta$	$\tau_1$   $\tau_2$
proposed $\ell_{\infty,1}$ proximal LMS	0.08				0.1
proposed RW $\ell_{\infty,1}$ proximal LMS	0.01	0.01			0.05
ZA-Sparse LMS [6]	0.03				
RZA-Sparse LMS [6]	0.03	0.45			
$\ell_{2,0}$ jointly sparse [24]			8	0.1	
$\ell_{2,1}$ jointly sparse [24]				0.03	
RW $\ell_{2,1}$ jointly sparse [24]			100	1.5	

#### REFERENCES

- [1] S.-Y. Tu and A. H. Sayed, “Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks,” *IEEE Trans. Signal Process.*, vol. 60, no. 12, pp. 6217–6234, Dec. 2012.
- [2] C. G. Lopes and A. H. Sayed, “Diffusion least-mean squares over adaptive networks: Formulation and performance analysis,” *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3122–3136, Jul. 2008.
- [3] F. Cattivelli, C. G. Lopes, and A. H. Sayed, “Diffusion recursive least-squares for distributed estimation over adaptive networks,” *IEEE Trans. Signal Process.*, vol. 56, no. 5, pp. 1865–1877, May 2008.
- [4] L. Li and J. Chambers, “Distributed adaptive estimation based on the APA algorithm over diffusion networks with changing topology,” in *Proc. IEEE SSP*, 2009, pp. 757–760.
- [5] Y. Liu, C. Li, and Z. Zhang, “Diffusion sparse least-mean squares over networks,” *IEEE Trans. Signal Process.*, vol. 60, no. 8, pp. 4480–4485, Aug. 2012.
- [6] P. Di Lorenzo and A. H. Sayed, “Sparse distributed learning based on diffusion adaptation,” *IEEE Trans. Signal Process.*, vol. 61, no. 6, pp. 1419–1433, Mar. 2013.
- [7] A. H. Sayed, “Diffusion adaptation over networks,” in *Academic Press Library in Signal Processing*, R. Chellapa and S. Theodoridis, Eds., vol. 3, pp. 322–454. Elsevier, 2014.
- [8] A. H. Sayed, “Adaptive networks,” *Proc. of the IEEE*, vol. 102, no. 4, pp. 460–497, Apr. 2014.
- [9] A. H. Sayed, *Adaptation, Learning, and Optimization over Networks*, vol. 7, Now Publishers Inc., Hanover, MA, USA, Jul. 2014.
- [10] X. Mao, Y. Gu, and W. Yin, “Walk proximal gradient: An energy-efficient algorithm for consensus optimization,” *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 2048–2060, 2018.
- [11] X. Mao, K. Yuan, Y. Hu, A. H. Sayed, and W. Yin, “Walkman: A communication-efficient random-walk algorithm for decentralized optimization,” *IEEE Trans. Signal Process.*, vol. 68, pp. 2513 – 2528, 2020.

- [12] J. Chen, C. Richard, and A. H. Sayed, "Multitask diffusion adaptation over networks," *IEEE Trans. Signal Process.*, vol. 62, no. 16, pp. 4129–4144, Aug. 2014.
- [13] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Multitask diffusion LMS with sparsity-based regularization," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 3516–3520.
- [14] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Proximal multitask learning over networks with sparsity-inducing coregularization," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6329–6344, Dec. 2016.
- [15] J. Chen, C. Richard, and A. H. Sayed, "Multitask diffusion adaptation over networks with common latent representations," *IEEE J. Sel. Topics Signal Process.*, vol. 11, no. 3, pp. 563–579, 2017.
- [16] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Diffusion LMS for multitask problems with local linear equality constraints," *IEEE Trans. Signal Process.*, vol. 65, no. 19, pp. 4979–4993, Oct. 2017.
- [17] J. Chen, C. Richard, and A. H. Sayed, "Diffusion LMS over multitask networks," *IEEE Trans. Signal Process.*, vol. 63, no. 11, pp. 2733–2748, Jun. 2015.
- [18] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Multitask diffusion adaptation over asynchronous networks," *IEEE Trans. Signal Process.*, vol. 64, no. 11, pp. 2835–2850, 2016.
- [19] M. J. Piggott and V. Solo, "Stability of adaptive network algorithms in multitask environments," in *Proc. IEEE CDC*, 2017, pp. 1472–1477.
- [20] D. Jin, J. Chen, C. Richard, J. Chen, and A. H. Sayed, "Affine combination of diffusion strategies over networks," *IEEE Trans. Signal Process.*, vol. 68, pp. 2087–2104, Dec. 2020.
- [21] V. C. Gogineni and M. Chakraborty, "Diffusion affine projection algorithm for multitask networks," in *Proc. Asia-Pacific Signal Inf. Process. Assoc.*, 2018, pp. 201–206.
- [22] V. C. Gogineni and M. Chakraborty, "Improving the performance of multitask diffusion APA via controlled inter-cluster cooperation," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 67, no. 3, pp. 903–912, 2020.
- [23] V. C. Gogineni and M. Chakraborty, "Partial diffusion affine projection algorithm over clustered multitask networks," in *Proc. IEEE ISCAS*, 2019, pp. 1–5.
- [24] C. Li, S. Huang, Y. Liu, and Z. Zhang, "Distributed jointly sparse multitask learning over networks," *IEEE Trans. Cybern.*, vol. 48, no. 1, pp. 151–164, Jan. 2018.
- [25] J. A. Bazerque and G. B. Giannakis, "Distributed spectrum sensing for cognitive radio networks by exploiting sparsity," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1847–1862, 2010.
- [26] M. F. Duarte, S. Sarvotham, D. Baron, M. B. Wakin, and R. G. Baraniuk, "Distributed compressed sensing of jointly sparse signals," in *Proc. Asilomar Conf. Signals Syst. Comput.*, Oct./Nov. 2005, pp. 1537–1541.
- [27] M. Kocic, D. Brady, and M. Stojanovic, "Sparse equalization for realtime digital underwater acoustic communications," in *Proc. OCEANS*, San Diego, CA, USA, 1995, pp. 1417–1422.
- [28] M. Masood, L. H. Afify, and T. Y. Al-Naffouri, "Efficient coordinated recovery of sparse channels in massive MIMO," *IEEE Trans. Signal Process.*, vol. 63, no. 1, pp. 104–118, 2015.
- [29] Y. Gu and M. Wang, "Learning distributed jointly sparse systems by collaborative LMS," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, May 2014, pp. 7228–7232.
- [30] S. Banert and R. I. Bot, "A general double-proximal gradient algorithm for d.c. programming," *Math. Program.*, vol. 178, pp. 301–326, Nov. 2019.
- [31] N. Parikh and S. Boyd, "Proximal algorithms," *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 127–239, 2014.
- [32] D. P. Bertsekas, "Incremental proximal methods for large scale convex optimization," *Math. Program.*, vol. 129, no. 2, pp. 163, Jun. 2011.
- [33] W. M. Wee and I. Yamada, "A proximal splitting approach to regularized distributed adaptive estimation in diffusion networks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, May 2013, pp. 5420–5424.
- [34] S. Vlaski and A. H. Sayed, "Proximal diffusion for stochastic costs with non-differentiable regularizers," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Apr. 2015, pp. 3352–3356.
- [35] S. Vlaski, L. Vandenbergh, and A. H. Sayed, "Diffusion stochastic optimization with non-smooth regularizers," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Mar. 2016, pp. 4149–4153.
- [36] J. Huang and T. Zhang, "The benefit of group sparsity," *The Annals of Statistics*, vol. 38, no. 4, pp. 1978–2004, Aug. 2010.
- [37] L. Meier, S. Van De Geer, and P. Bühlmann, "The group lasso for logistic regression," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 70, no. 1, pp. 53–71, 2008.
- [38] Y. Chen and A. O. Hero, "Recursive  $\ell_{1,\infty}$  group lasso," *IEEE Trans. Signal Process.*, vol. 60, no. 8, pp. 3978–3987, 2012.
- [39] Y. Chen, Y. Gu, and A. O. Hero, "Regularized least-mean-square algorithms," *Arxiv preprint arXiv:1012.5066*, 2010.
- [40] S. N. Negahban and M. J. Wainwright, "Simultaneous support recovery in high dimensions: Benefits and perils of block  $\ell_1/\ell_\infty$ -regularization," *IEEE Trans. Inf. Theory*, vol. 57, no. 6, pp. 3841–3863, Jun. 2011.
- [41] E. Esser, M. Moller, S. Osher, G. Sapiro, and J. Xin, "A convex model for nonnegative matrix factorization and dimensionality reduction on physical space," *IEEE Trans. Image Process.*, vol. 21, no. 7, pp. 3239–3252, 2012.
- [42] E. M. Eksioğlu, "Group sparse RLS algorithms," *International journal of adaptive control and signal processing*, vol. 28, no. 12, pp. 1398–1412, Dec. 2013.
- [43] J. A. Tropp, "Algorithms for simultaneous sparse approximation. Part II: Convex relaxation," *Signal Process.*, vol. 86, no. 3, pp. 589 – 602, 2006.
- [44] David R. Hunter, Kenneth Lange, Departments Of Biomathematics, and Human Genetics, "A tutorial on MM algorithms," *Amer. Statist.*, pp. 30–37, 2004.
- [45] E. J. Candès, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted  $\ell_1$  minimization," *Journal of Fourier Analysis and Applications*, vol. 14, no. 5, pp. 877–905, 2008.
- [46] A. H. Sayed, *Adaptive Filters*, John Wiley & Sons, Inc., 2008.
- [47] D. Jin, J. Chen, C. Richard, and J. Chen, "Model-driven online parameter adjustment for zero-attracting LMS," *Signal Process.*, vol. 152, pp. 373 – 383, Nov. 2018.
- [48] J. Chen, C. Richard, J. C. M. Bermudez, and P. Honeine, "Nonnegative least-mean-square algorithm," *IEEE Trans. Signal Process.*, vol. 59, no. 11, pp. 5225–5235, Nov. 2011.
- [49] J. Arenas-Garcia, A. R. Figueiras-Vidal, and A. H. Sayed, "Mean-square performance of a convex combination of two adaptive filters," *IEEE Trans. Signal Process.*, vol. 54, no. 3, pp. 1078–1090, Mar. 2006.
- [50] V. H. Nascimento, M. T. M. Silva, R. Candido, and J. Arenas-Garcia, "A transient analysis for the convex combination of adaptive filters," in *Proc. IEEE SSP*, Aug. 2009, pp. 53–56.
- [51] M. T. M. Silva, V. H. Nascimento, and J. Arenas-Garcia, "A transient analysis for the convex combination of two adaptive filters with transfer of coefficients," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Mar. 2010, pp. 3842–3845.
- [52] M. Fiedler, "Algebraic connectivity of graphs," *Czechoslovak Mathematical Journal*, vol. 23, no. 2, pp. 298–305, 1973.
- [53] A. Simões and J. Xavier, "FADE: Fast and asymptotically efficient distributed estimator for dynamic networks," *IEEE Trans. Signal Process.*, vol. 67, no. 8, pp. 2080–2092, Apr. 2019.
- [54] H. Zhu, G. B. Giannakis, and A. Cano, "Distributed in-network channel decoding," *IEEE Trans. Signal Process.*, vol. 57, no. 10, pp. 3970–3983, 2009.
- [55] H. Zhu, A. Cano, and G. B. Giannakis, "Distributed consensus-based demodulation: algorithms and error analysis," *IEEE Trans. Wireless Commun.*, vol. 9, no. 6, pp. 2044–2054, 2010.



- [56] M. B. Porter, "The bellhop manual and user's guide: Preliminary draft," tech. rep., Heat, Light, and Sound Research, Inc., LaJolla, CA, USA, 2011  
[Online] Available: <http://oalib.hlsresearch.com/Rays/HLS-2010.pdf>.