

# Multitask diffusion adaptation over networks with common latent representations

Jie Chen, *Member, IEEE*, Cédric Richard, *Senior Member, IEEE*,

Ali H. Sayed, *Fellow Member, IEEE*

## Abstract

Online learning with streaming data in a distributed and collaborative manner can be useful in a wide range of applications. This topic has been receiving considerable attention in recent years with emphasis on both single-task and multitask scenarios. In single-task adaptation, agents cooperate to track an objective of common interest, while in multitask adaptation agents track multiple objectives simultaneously. Regularization is one useful technique to promote and exploit similarity among tasks in the latter scenario. This work examines an alternative way to model relations among tasks by assuming that they all share a common latent feature representation. As a result, a new multitask learning formulation is presented and algorithms are developed for its solution in a distributed online manner. We present a unified framework to analyze the mean-square-error performance of the adaptive strategies, and conduct simulations to illustrate the theoretical findings and potential applications.

## Index Terms

Multitask learning, distributed optimization, common latent subspace, online adaptation, diffusion strategy, collaborative processing, performance analysis.

## I. INTRODUCTION

Multi-agent networks usually consist of a large number of interconnected agents or nodes. Interconnections between the agents allow them to share information and collaborate in order to solve complex tasks collectively. Examples abound in the realm of social, economic and biological networks. Distributed algorithms over such networks offer a valuable alternative to centralized solutions with useful properties such as scalability, robustness, and decentralization. When endowed with adaptation abilities, these algorithms enable agents to continuously learn

Copyright (c) 2016 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

The work of J. Chen was supported in part by the NSFC grant 61671382. The work of C. Richard was supported in part by the Agence Nationale pour la Recherche, France, (ODISSEE project, ANR-13-ASTR-0030). The work of A. H. Sayed was supported in part by NSF grants CIF-1524250 and ECCS-1407712. A short and preliminary version of this work appears in the conference publication [1]. J. Chen is with CIAIC of School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an, 710072, China (e-mail: [dr.jie.chen@ieee.org](mailto:dr.jie.chen@ieee.org)). C. Richard is with the Université de Nice Sophia-Antipolis, CNRS, France (e-mail: [cedric.richard@unice.fr](mailto:cedric.richard@unice.fr)), in collaboration with Morpheme team (INRIA Sophia-Antipolis). A. H. Sayed is with the department of electrical engineering, University of California, Los Angeles, CA 9005-1594, USA (email: [sayed@ee.ucla.edu](mailto:sayed@ee.ucla.edu)).

and adapt in an online manner to concept drifts in their data streams [2], [3]. Broadly, distributed strategies for online parameter estimation can be applied to single-task or multi-task scenarios. In the first case, agents cooperate with each other to estimate a single parameter vector of interest, such as tracking a common target. Reaching consensus among the agents is critical for successful inference in these problems. In the multitask case, the agents cooperate to estimate multiple parameter vectors simultaneously, such as tracking a collection of targets moving in formation [4].

Extensive studies have been conducted on adaptive distributed strategies for single-task problems. Existing techniques include incremental [5]–[8], consensus [9]–[11], and diffusion strategies [1], [2], [12]–[17]. Incremental techniques require determining a cyclic path that runs across all nodes, which is generally a challenging (NP-hard) task to perform. Besides, feature makes the incremental strategies sensitive to link failures and problematic for adaptation. Consensus techniques aim to reach an agreement among nodes on the estimate of interest via local information exchanges, but they have been shown [2], [3] to suffer from instability problems when used in the context of adaptive networks due to an inherent asymmetry in the update equations. Diffusion techniques, on the other hand, have been shown to have superior stability and performance ranges [18] than consensus-based implementations. For these reasons, we shall focus on diffusion-type implementations in this paper.

Besides single-task scenarios, there are also applications where it is desirable to estimate multiple parameter vectors at the same time, rather than promote consensus among all agents [19]. For example, geosensor networks that monitor dynamic spatial fields, such as temperature or windspeed variations in geographic environments, require node-specific estimation problems that are able to take advantage of the spatial correlation between the measurements of neighboring nodes [20], [21]. A second example is the problem of collaborative target tracking where agents track several objects simultaneously [4], [19]. Motivated by these applications, there have been several variations of distributed strategies to deal with multitask scenarios as well. Existing strategies mostly depend on how the tasks relate to each other and on exploiting some prior information. In a first scenario, nodes are grouped into clusters, and each cluster of nodes is interested in estimating its own parameter vector. Although clusters may generally have distinct though related estimation tasks to perform, the nodes may still be able to capitalize on inductive transfer between clusters to improve their estimation accuracy. Multitask diffusion strategies were developed to perform estimation under these conditions [4], [22]. One useful way to do so is to employ regularization. A couple of other useful works have also addressed variations of this scenario where the only available information is that clusters may exist in the network but nodes do not know which other nodes share the same estimation task [23]–[25]. In [26], the authors use multitask diffusion adaptation with a node clustering strategy to identify a model between the gait information and electroencephalographic signals. In [27], the authors consider the framework in [4] to devise a distributed strategy that allows each node in the network to locally adapt inter-cluster cooperation weights. The authors in [28] promote cooperation between clusters with  $\ell_1$ -norm co-regularizers. They derive a closed-form expression of the proximal operator, and introduce a strategy that also allows each node to automatically set its inter-cluster cooperation weights. The works in [29], [30] propose alternative node clustering strategies. In a

second scenario, it is assumed that there are parameters of global interest to all nodes in the network, a collection of parameters of common interest within sub-groups of nodes, and a set of parameters of local interest at each node. A diffusion strategy was developed to perform estimation under these conditions [31], [32]. Likewise, in the works [33]–[35], distributed algorithms are derived to estimate node-specific parameter vectors that lie in a common latent signal subspace. In another work [36], the diffusion LMS algorithm is extended to deal with structured criteria built upon groups of variables, leading to a flexible framework that can encode various structures in the parameters. An unsupervised strategy to differentially promote or inhibit collaboration between nodes depending on their group is also introduced.

Alternatively, in recent years, there has been an increasing interest in modeling relations between tasks by assuming that all tasks share a common feature representation in a latent subspace [37]–[39]. The authors in [38] proposed a non-convex method based on Alternating Structure Optimization (ASO) for identifying the task structure. A convex relaxation of this approach was developed in [40]. In [39], the authors showed the equivalence between ASO, clustered multitask learning [41], [42] and their convex relaxations. The efficiency of such task relationships has been demonstrated in these works for clustering and classification problems. In our preliminary work [1], we introduced this framework within the context of distributed online adaptation over networks. Useful applications can be envisaged. First, consider the case where the common subspace is spanned by certain selected columns of the identity matrix. This means that a subset of the entries of the parameter vector to be estimated are common to all nodes while no further restriction is imposed on the other entries. Another example concerns beamforming for antenna arrays with a generalized side-lobe canceller (GSC). The latent subspace corresponds to the space where interfering signals reside [43]. A third example deals with cooperative spectrum sensing in cognitive radios, where the common latent subspace characterizes common interferers [31].

Drawing on these motivations, this paper deals with distributed learning and adaptation over multitask networks with common latent representation subspaces. Algorithms are designed accordingly, and their performance analyzed. The contributions of this work include the following main aspects:

- We formulate a new multitask estimation problem, which assumes that all tasks share a common latent subspace representation in addition to node-specific contributions. Additional constraints can be incorporated if needed. This work contrasts with earlier works [4], [28], where the inductive transfer between learning tasks is promoted by regularizers. It also differs from [31], which considers direct models by stacking local and global variables in an augmented parameter vector. Moreover, the work [38] uses a similar inductive transfer model but the common latent subspace is unknown and embedded into a joint estimation process. Our work is the first one to introduce an online estimation algorithm over networks. Estimating the common latent subspace of interest within this context is a challenging perspective.
- We explain how this formulation can be tailored to fit individual application contexts by considering additional model constraints. We illustrate this fact by considering two convex optimization problems and the associated distributed online algorithms. The first algorithm is a generalization in some sense of the diffusion LMS

algorithm, which can be retrieved by defining the low-dimensional common latent subspace as the whole parameter space. The second algorithm uses  $\ell_2$ -norm regularization to account for the multitask nature of the problem. This opens the way to other regularization schemes depending on the application at hand.

- We present a unified framework for analyzing the performance of these algorithms. This framework also allows to address the performance analysis of the multitask algorithms in [4], [19], [44], [45] in a generic manner, though these analyses were performed independently of each other in these works.

The rest of the paper is organized as follows. Section II introduces the multitask estimation problem considered in this paper. Then, two distributed learning strategies are derived in Section III by imposing different constraints on common and node-specific representation subspaces. Section IV provides a general framework for analyzing distributed algorithms of this form. In Section V, experiments are conducted to illustrate the characteristics of these algorithms. Section VI concludes the paper and connects our work with several other learning strategies.

**Notation.** Normal font  $x$  denotes scalars. Boldface small letters  $\mathbf{x}$  denote vectors. All vectors are column vectors. Boldface capital letters  $\mathbf{X}$  denote matrices. The asterisk  $(\cdot)^*$  denotes complex conjugation for scalars and complex-conjugate transposition for matrices. The superscript  $(\cdot)^\top$  represents transpose of a matrix or a vector, and  $\|\cdot\|$  is the  $\ell_2$ -norm of its matrix or vector argument.  $\text{Re}\{\cdot\}$  and  $\text{Im}\{\cdot\}$  denote the real and imaginary parts of their complex argument, respectively. Matrix trace is denoted by  $\text{trace}(\cdot)$ . The operator  $\text{col}\{\cdot\}$  stacks its vector arguments on the top of each other to generate a connected vector. The operator  $\text{diag}\{\cdot\}$  formulates a (block) diagonal matrix with its arguments. Identity matrix of size  $N \times N$  is denoted by  $\mathbf{I}_N$ . Kronecker product is denoted by  $\otimes$ , and expectation is denoted by  $\mathbb{E}\{\cdot\}$ . We denote by  $\mathcal{N}_k$  the set of node indices in the neighborhood of node  $k$ , including  $k$  itself, and  $|\mathcal{N}_k|$  its set cardinality.

## II. MATCHED SUBSPACE ESTIMATION OVER MULTITASK NETWORKS

### A. Multitask estimation problems over networks

Consider a connected network composed of  $N$  nodes. The problem is to estimate an  $L \times 1$  unknown vector  $\mathbf{w}_k^o$  at each node  $k$  from collected measurements. At each time  $n$ , node  $k$  has access to local streaming measurements  $\{d_k(n), \mathbf{x}_{k,n}\}$ , where  $d_k(n)$  is a scalar zero-mean reference signal, and  $\mathbf{x}_{k,n}$  is a  $1 \times L$  zero-mean row regression vector with covariance matrix  $\mathbf{R}_{\mathbf{x},k} = \mathbb{E}\{\mathbf{x}_{k,n}^* \mathbf{x}_{k,n}\} > 0$ . The data at agent  $k$  and time  $n$  are assumed to be related via the linear model:

$$d_k(n) = \mathbf{x}_{k,n} \mathbf{w}_k^o + z_k(n) \quad (1)$$

where  $\mathbf{w}_k^o$  is an unknown complex parameter vector, and  $z_k(n)$  is a zero-mean i.i.d. noise with variance  $\sigma_{z,k}^2 = \mathbb{E}\{|z_k(n)|^2\}$ . The noise signal  $z_k(n)$  is assumed to be independent of any other signal. Let  $J_k(\mathbf{w})$  be a differentiable convex cost function at agent  $k$ . In this paper, we shall consider the mean-square-error criterion:

$$J_k(\mathbf{w}) = \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n} \mathbf{w}|^2\} \quad (2)$$

It is clear from (1) that each  $J_k(\mathbf{w})$  is minimized at  $\mathbf{w}_k^o$ . We refer to each parameter  $\mathbf{w}_k^o$  to estimate (or model in a more general sense) as a task. Depending on whether the minima of all  $J_k(\mathbf{w})$  are achieved at the same  $\mathbf{w}_k^o$  or not, the distributed learning problem can be single-task or multitask oriented [4].

With single-task networks, all agents aim at estimating the same parameter vector  $\mathbf{w}^o$  shared by the entire network, that is,

$$\mathbf{w}_k^o = \mathbf{w}^o \quad (3)$$

for all  $k \in \{1, \dots, N\}$ . Several popular collaborative strategies, such as diffusion LMS [1], [2], [13], [14], were derived to address this problem by seeking the minimizer of the following aggregate cost function:

$$J^{\text{glob}}(\mathbf{w}) = \sum_{k=1}^N J_k(\mathbf{w}) \quad (4)$$

in a distributed manner. Since the individual costs (2) admit the same solution,  $\mathbf{w}^o$  is also the solution of (4). It has been shown that using proper cooperative strategies to solve (4) can improve the estimation performance [2], [3].

With multitask networks, each agent aims at determining a local parameter vector  $\mathbf{w}_k^o$ . It is assumed that some similarities or relations exist among the parameter vectors of neighboring agents so that cooperation can still be meaningful, namely,

$$\mathbf{w}_k^o \sim \mathbf{w}_\ell^o \text{ if } \ell \in \mathcal{N}_k \quad (5)$$

where the symbol  $\sim$  refers to a similarity relationship in some sense, which can be exploited to enhance performance. Depending on the problem characteristics, this property can be promoted in several ways, e.g., by introducing some regularization term, or by assuming a common latent structure. Networks may also be structured into clusters where agents within each cluster estimate the same parameter vector [4], [44].

### B. Node-specific subspace constraints

Although agents aim to estimate distinct minimizers  $\mathbf{w}_k^o$ , exploiting relationships between solutions can make cooperation among agents beneficial. Regularization is one popular technique for introducing prior information about the solution. It can improve estimation accuracy though it may introduce bias [4], [19], [46]. In this paper, we explore an alternative strategy that assumes that the hypothesis spaces partially overlap. Specifically, we assume that each  $\mathbf{w}_k^o$  can be expressed in the form:

$$\mathbf{w}_k^o = \Theta \mathbf{u}^o + \boldsymbol{\epsilon}_k^o \quad (6)$$

where  $\Theta \mathbf{u}^o$  is common to all nodes with  $\Theta$  denoting an  $L \times M$  matrix with known entries and  $\mathbf{u}^o$  an unknown  $M \times 1$  parameter vector (common to all nodes), and where  $\boldsymbol{\epsilon}_k^o$  is an unknown node-specific component. We assume that matrix  $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M]$  is full-rank with  $M \leq L$ . Overcomplete sets of column vectors  $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$  may be advantageous in some scenarios but this usually requires to impose further constraints such as sparsity over  $\mathbf{u}^o$ . We shall not discuss this case further in order to focus on the main points of the presentation. Model (6) means

that all tasks share the same parameter vector  $\Theta \mathbf{u}^o$ , which lies in the subspace spanned by columns of  $\Theta$ . This subspace representation can be useful in several applications. For instance, consider the case where  $\Theta$  is composed of selected columns of the identity matrix  $\mathbf{I}_L$ . This means that a subset of the entries of  $\mathbf{w}_k^o$  are common to all agents while no further assumptions are imposed on the other entries. This situation is a natural generalization of the single-task scenario. Another example concerns beamforming problems with a generalized sidelobe canceller (GSC), where  $\Theta$  acts as a blocking matrix to cancel signal components that lie in the constraint space [43]. In machine learning, formulation (6) is referred to as the alternating structure optimization (ASO) problem [38], [39]. The subspace  $\Theta$  is, however, learnt simultaneously via a non-convex optimization procedure. In what follows, we shall assume that  $\Theta$  is known by each agent.

Before proceeding further, we clarify the difference between model (6) addressed here and in our preliminary work [1], and the model studied in [31], [32], [35], [47], [48]. In these last works, the authors consider particular information access models where global and local components are assumed to be related to distinct regressors. The centralized problem can then be formulated by stacking the global and local regressors, and by considering a parameter vector augmented accordingly. In our work, motivated by applications of the latent space model in batch-mode learning, we address the problem where the parameter vectors to be estimated lie in global and local latent subspaces. We do not need to distinguish explicitly between global and local regressors. Instead, as shown in the sequel, some extra conditions are needed so that model (6) is identifiable. Among other possibilities, we shall investigate two strategies where constraints on  $\Theta$  and  $\epsilon_k^o$  are imposed.

Replacing (6) into (2), the global cost function is expressed as a function of a common parameter  $\mathbf{u}$  and node-specific perturbations  $\{\epsilon_k\}_{k=1}^N$ :

$$J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) = \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n}(\Theta \mathbf{u} + \epsilon_k)|^2\} \quad (7)$$

We expect the estimation of  $\mathbf{w}_k^o$  by each agent to benefit from the cooperative estimation of  $\mathbf{u}$ . Problem (7) is still insufficient for estimating the tasks  $\{\mathbf{w}_k^o\}$ . This is because the decomposition  $\mathbf{w}_k = \Theta \mathbf{u} + \epsilon_k$  is not unique. Indeed, given any optimum solution  $\{\bar{\mathbf{u}}, \bar{\epsilon}_k\}$ , and any  $\mathbf{s} = \Theta \mathbf{x}$ , we can generate another optimum solution by considering the shift  $\{\bar{\mathbf{u}} - \mathbf{x}, \bar{\epsilon}_k + \mathbf{s}\}$ . This ambiguity prevents us from deriving collaboration strategies based on  $\mathbf{u}$ . From the point of view of convex analysis, the Hessian matrix of (7) is rank deficient and no unique solution exists.

### III. PROBLEM FORMULATIONS AND SOLUTION ALGORITHMS

Problem (7) can be modified to make it well-determined and more meaningful. In this section, among other possibilities, we investigate two strategies that consist of imposing further constraints and derive the corresponding distributed algorithms. These two formulations guarantee the uniqueness of the solution and have clear interpretations.

### A. Node-specific subspace constraints

We restrict the node-specific components  $\{\epsilon_k\}_{k=1}^N$  to lie in the complementary subspace to  $\text{span}(\Theta)$ . The problem can be formulated as:

$$\begin{aligned} \min_{\mathbf{u}, \{\epsilon_k\}_{k=1}^N} J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) \\ \text{subject to } \epsilon_k \in \text{span}(\Theta_{\perp}), \quad \forall k = 1, \dots, N \end{aligned} \quad (8)$$

where the  $L - M$  columns of matrix  $\Theta_{\perp}$  span the complementary subspace to  $\text{span}(\Theta)$ , that is,  $\Theta^* \Theta_{\perp} = \mathbf{0}$ . We write:

$$\epsilon_k = \Theta_{\perp} \xi_k \quad (9)$$

where  $\xi_k$  is a column vector of size  $(L - M)$ . Now, replacing (9) into (8), the optimization problem becomes unconstrained and the objective function is given by:

$$\begin{aligned} J^{\text{glob}}(\mathbf{u}, \{\xi_k\}_{k=1}^N) \\ &= \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n}(\Theta \mathbf{u} + \Theta_{\perp} \xi_k)|^2\} \\ &= \sum_{k=1}^N \mathbb{E}\{|d_k(n)|^2\} + \mathbf{u}^* \Theta^* \left( \sum_{k=1}^N \mathbf{R}_{x,k} \right) \Theta \mathbf{u} + \sum_{k=1}^N \xi_k^* \Theta_{\perp}^* \mathbf{R}_{x,k} \Theta_{\perp} \xi_k + 2 \text{Re} \left\{ \mathbf{u}^* \Theta^* \sum_{k=1}^N \mathbf{R}_{x,k} \Theta_{\perp} \xi_k \right\} \\ &\quad - 2 \text{Re} \left\{ \sum_{k=1}^N \mathbf{p}_{dx,k}^* \Theta \mathbf{u} \right\} - 2 \text{Re} \left\{ \sum_{k=1}^N \mathbf{p}_{dx,k}^* \Theta_{\perp} \xi_k \right\} \end{aligned} \quad (10)$$

where  $\mathbf{R}_{x,k} = \mathbb{E}\{\mathbf{x}_{k,n}^* \mathbf{x}_{k,n}\}$  is the covariance matrix of  $\mathbf{x}_{k,n}$ , and  $\mathbf{p}_{dx,k} = \mathbb{E}\{d_k(n) \mathbf{x}_{k,n}^*\}$  is the covariance vector between the input data  $\mathbf{x}_{k,n}$  and the reference output data  $d_k(n)$ .

*Lemma 1:* Problem (8) has a unique solution with respect to  $\mathbf{u}$  and  $\{\epsilon_k\}_{k=1}^N$  if the perturbations  $\{\epsilon_k\}_{k=1}^N$  lie in a subspace orthogonal to  $\text{span}(\Theta)$ .  $\blacksquare$

Proof of Lemma 1 is provided in Appendix A. We shall now derive a distributed algorithm to seek the minimizer of (8). Focusing on the terms that depend on  $\mathbf{u}$  in (10), and setting parameters  $\xi_k$  to their optimum values  $\xi_k^o$ , we consider first the global cost function over the variable  $\mathbf{u}$ :

$$\begin{aligned} J_u^{\text{glob}}(\mathbf{u}) &= \sum_{k=1}^N \left( \mathbf{u}^* \Theta^* \mathbf{R}_{x,k} \Theta \mathbf{u} + 2 \text{Re} \left\{ \mathbf{u}^* \Theta^* \mathbf{R}_{x,k} \Theta_{\perp} \xi_k^o \right\} - 2 \text{Re} \left\{ \mathbf{p}_{dx,k}^* \Theta \mathbf{u} \right\} + g_k(\xi_k^o) \right) \\ &= \sum_{k=1}^N J_{u,k}(\mathbf{u}) \end{aligned} \quad (11)$$

where  $g_k(\xi_k^o)$  collects all the terms depending only on  $\xi_k^o$  in (10). The term  $\sum_{k=1}^N \mathbb{E}\{|d_k(n)|^2\}$  is discarded because it is constant with respect to the arguments  $\mathbf{u}$  and  $\{\xi_k\}_{k=1}^N$ . Since  $J_u^{\text{glob}}(\mathbf{u})$  has a unique minimizer for all nodes over the network, we can use a single-task adapt-then-combine (ATC) diffusion strategy to estimate  $\mathbf{u}^o$  [13], [15].

We introduce a right-stochastic matrix  $\mathbf{C}$  with nonnegative entries  $c_{\ell k}$  such that:

$$\sum_{k=1}^N c_{\ell k} = 1, \quad \text{and} \quad c_{\ell k} = 0 \text{ if } k \notin \mathcal{N}_{\ell} \quad (12)$$

With each node  $k$ , we associate the local cost over the variable  $\mathbf{u}$ :

$$J_{u,k}^{\text{loc}}(\mathbf{u}) = \sum_{\ell \in \mathcal{N}_k} c_{\ell k} J_{u,\ell}(\mathbf{u}) \quad (13)$$

Observe that  $\sum_{k=1}^N J_{u,k}^{\text{loc}}(\mathbf{u}) = J_u^{\text{glob}}(\mathbf{u})$  because matrix  $\mathbf{C}$  is right-stochastic. Since  $J_u^{\text{glob}}(\mathbf{u})$  is quadratic with respect to  $\mathbf{u}$ , it can be expressed at each node  $k$  as follows:

$$\begin{aligned} J_u^{\text{glob}}(\mathbf{u}) &= J_{u,k}^{\text{loc}}(\mathbf{u}) + \sum_{\ell \neq k} J_{u,\ell}^{\text{loc}}(\mathbf{u}) \\ &= J_{u,k}^{\text{loc}}(\mathbf{u}) + \sum_{\ell \neq k} \|\mathbf{u} - \mathbf{u}^o\|_{\nabla^2 J_{u,\ell}^{\text{loc}}}^2 \end{aligned} \quad (14)$$

where  $\nabla^2 J_{u,\ell}^{\text{loc}}$  denotes the Hessian matrix of  $J_{u,\ell}^{\text{loc}}(\mathbf{u})$  with respect to  $\mathbf{u}$ , and  $\|\mathbf{u}\|_{\Sigma}^2$  is the squared norm of  $\mathbf{u}$  weighted by any positive semi-definite matrix  $\Sigma$ , i.e.,  $\|\mathbf{u}\|_{\Sigma}^2 = \mathbf{u}^* \Sigma \mathbf{u}$ . Following an argument based on the Rayleigh-Ritz characterization of eigenvalues [13, Sec. 3.1], we approximate  $\nabla^2 J_{u,\ell}^{\text{loc}}$  by a multiple of the identity matrix, so that  $\|\mathbf{u} - \mathbf{u}^o\|_{\nabla^2 J_{u,\ell}^{\text{loc}}}^2 \approx b_{\ell k} \|\mathbf{u} - \mathbf{u}^o\|^2$ .

Minimizing (14) in two successive steps yields:

$$\phi_{k,n} = \mathbf{u}_{k,n-1} - \mu \nabla J_{u,k}^{\text{loc}}(\mathbf{u}_{k,n-1}) \quad (15)$$

$$\mathbf{u}_{k,n} = \phi_{k,n} + \mu \sum_{\ell \neq k} b_{\ell k} (\mathbf{u}^o - \mathbf{u}_{k,n-1}) \quad (16)$$

where  $\mu$  is a positive step size. Its choice to ensure stability of the algorithm will be elaborated on later in Sec. IV. Now, note the following. First, iteration (16) requires knowledge of  $\mathbf{u}^o$ , which is not available. Each node  $\ell$  has a readily available, however, an approximation for  $\mathbf{u}^o$ , which is  $\phi_{k,n}$ . Therefore, we replace  $\mathbf{u}^o$  by  $\phi_{k,n}$  in (16). Second,  $\phi_{k,n}$  at node  $k$  is generally a better estimate for  $\mathbf{u}^o$  than  $\mathbf{u}_{k,n-1}$  since it is obtained by incorporating information from the neighbors through (15). Therefore, we replace  $\mathbf{u}_{k,n-1}$  by  $\phi_{k,n}$  in (16). Then, absorbing coefficients  $b_{\ell k}$  into another set of nonnegative coefficients that satisfies:

$$\sum_{\ell=1}^N a_{\ell k} = 1, \quad \text{and} \quad a_{\ell k} = 0 \text{ if } \ell \notin \mathcal{N}_k, \quad (17)$$

which means that matrix  $\mathbf{A}$  with entries  $a_{\ell k}$  is left-stochastic, using an instantaneous approximation of the gradient, and limiting the summation in (16) to the neighbors of node  $\ell$  (see [13], [15] for more details on a similar derivation in the context of single-task diffusion strategies), we can update  $\mathbf{u}_{k,n}$  as follows:

$$\phi_{k,n} = \mathbf{u}_{k,n-1} + \mu \sum_{\ell \in \mathcal{N}_k} c_{\ell k} \Theta^* \mathbf{x}_{\ell,n}^* [d_{\ell}(n) - \mathbf{x}_{\ell,n}(\Theta \mathbf{u}_{k,n-1}) - \mathbf{x}_{\ell,n}(\Theta \perp \xi_{\ell,n-1})] \quad (18)$$

$$\mathbf{u}_{k,n} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \phi_{\ell,n} \quad (19)$$

where  $\xi_{k,n-1}$  is an estimate for the unknown minimizer  $\xi_k^o$ , to be evaluated as explained further ahead in (21).



Focusing on the terms that depend on  $\{\xi_k\}_{k=1}^N$  in (10), and setting parameter  $\mathbf{u}$  to its optimum value  $\mathbf{u}^o$ , we consider the global cost function over the variables  $\xi_k$ :

$$\begin{aligned} J_{\xi}^{\text{glob}}(\{\xi_k\}_{k=1}^N) &= \sum_{k=1}^N \left( \xi_k^* \Theta_{\perp}^* \mathbf{R}_{x,k} \Theta_{\perp} \xi_k + 2 \operatorname{Re} \left\{ \xi_k^* \Theta_{\perp}^* \mathbf{R}_{x,k} \Theta \mathbf{u}^o \right\} - 2 \operatorname{Re} \left\{ \mathbf{p}_{dx,k}^* \Theta_{\perp} \xi_k \right\} \right) + g'_k(\mathbf{u}^o) \\ &= \sum_{k=1}^N J_{\xi,k}(\xi_k) \end{aligned} \quad (20)$$

where  $g'_k(\xi_k^o)$  collects all the terms depending only on  $\mathbf{u}^o$  in (10). Now since the parameters  $\xi_k$  are node-specific, if no further constraints are imposed, they can be updated independently of each other via an LMS-type update:

$$\xi_{k,n} = \xi_{k,n-1} + \mu \Theta_{\perp}^* \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} (\Theta \mathbf{u}_{k,n-1} + \Theta_{\perp} \xi_{k,n-1})] \quad (21)$$

At each time instant  $n$ , node  $k$  updates its parameters  $\mathbf{u}_{k,n-1}$  and  $\xi_{k,n-1}$  using (18)–(19) and (21), respectively. The local estimate  $\mathbf{w}_{k,n}$  is then given by:

$$\mathbf{w}_{k,n} = \Theta \mathbf{u}_{k,n} + \Theta_{\perp} \xi_{k,n} \quad (22)$$

It is interesting to note that we can rewrite the algorithm without using the auxiliary variables  $\mathbf{u}_{k,n}$  and  $\{\xi_{k,n}\}_{k=1}^N$ , by substituting the relations:

$$\mathbf{u}_{k,n} = (\Theta^* \Theta)^{-1} \Theta^* \mathbf{w}_{k,n} \quad (23)$$

$$\xi_{k,n} = (\Theta_{\perp}^* \Theta_{\perp})^{-1} \Theta_{\perp}^* \mathbf{w}_{k,n} \quad (24)$$

into (18)–(19) and (21), respectively. Selecting  $\mathbf{C} = \mathbf{I}_N$  to avoid exchanging raw data and node-specific components, we can implement the update of  $\mathbf{w}_{k,n-1}$  to an intermediate value  $\psi_{k,n}$  as follows:

$$\begin{aligned} \psi_{k,n} &\stackrel{(a)}{=} \Theta \phi_{k,n} + \Theta_{\perp} \xi_{k,n} \\ &\stackrel{(b)}{=} \Theta \mathbf{u}_{k,n-1} + \Theta_{\perp} \xi_{k,n-1} + \mu \left[ (\Theta \Theta^* + \Theta_{\perp} \Theta_{\perp}^*) \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} (\Theta \mathbf{u}_{k,n-1} + \Theta_{\perp} \xi_{k,n-1})] \right] \\ &= \mathbf{w}_{k,n-1} + \mu \mathbf{S}_{\Theta} \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1}] \end{aligned} \quad (25)$$

with  $\mathbf{S}_{\Theta} = \Theta \Theta^* + \Theta_{\perp} \Theta_{\perp}^*$ . For step (a), we use (22) with the intermediate value  $\phi_{k,n}$  of  $\mathbf{u}_{k,n}$  in (18) and  $\xi_{k,n}$ . Step (b) follows from their adaptation steps (18) and (21). Now substituting (19) in (22) to aggregate the intermediate estimates of  $\mathbf{u}_{k,n}$  from the neighbors of node  $k$ , we arrive at the combination step:

$$\begin{aligned} \mathbf{w}_{k,n} &\stackrel{(22)}{=} \Theta \mathbf{u}_{k,n} + \Theta_{\perp} \xi_{k,n} \\ &\stackrel{(c)}{=} \Theta \sum_{\ell \in \mathcal{N}_k} a_{\ell k} (\Theta^* \Theta)^{-1} \Theta^* \psi_{\ell,n} + \Theta_{\perp} (\Theta_{\perp}^* \Theta_{\perp})^{-1} \Theta_{\perp}^* \psi_{k,n} \\ &= \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{P}_{\Theta} \psi_{\ell,n} + \mathbf{P}_{\Theta_{\perp}} \psi_{k,n} \end{aligned} \quad (26)$$

where  $\mathbf{P}_\Theta = \Theta(\Theta^*\Theta)^{-1}\Theta^*$  and  $\mathbf{P}_{\Theta_\perp} = \mathbf{I}_L - \mathbf{P}_\Theta$  are the projection matrices over subspaces  $\Theta$  and  $\Theta_\perp$ . For step (c), we use (23)–(24) with the intermediate estimate  $\psi_{k,n}$ . Finally, we arrive at the ATC strategy summarized in Algorithm 1.

The first step in (28) is an adaptation step where node  $k$  uses the data realizations  $\{d_k(n), \mathbf{x}_{k,n}\}$  to update its existing estimate  $\mathbf{w}_{k,n-1}$  to an intermediate value  $\psi_{k,n}$ . All other nodes in the network are performing a similar step. The second step in (29) is an aggregation step. To update its intermediate estimate to  $\mathbf{w}_{k,n}$ , each node  $k$  combines the existing estimates of its neighbors in the common latent subspace  $\Theta$  to build up a common representation, and refines it with a node-specific value in  $\Theta_\perp$ . In the special case when  $\mathbf{A} = \mathbf{I}_N$ , so that no information exchange is performed, the ATC strategy reduces to a non-cooperative solution where each node  $k$  runs its own individual descent algorithm.

Matrix  $\mathbf{S}_\Theta$  in the adaptation step (28) is positive-definite. It arises from the calculation of the gradient of (10) with respect to  $\mathbf{u}$  and  $\xi_k$ . The algorithm can be simplified by replacing  $\mathbf{S}_\Theta$  by  $\mathbf{I}_L$  in (28) without compromising the convergence of the method (as analyzed further ahead in Section IV). We then arrive at the recursion:

$$\psi_{k,n} = \mathbf{w}_{k,n-1} + \mu \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1}] \quad (27)$$

Strictly speaking, observe that  $\mathbf{S}_\Theta = \mathbf{I}_L$  if, and only if, the columns of  $\Theta$  and  $\Theta_\perp$  form an orthonormal basis of  $\mathbb{R}^L$ . Note that the adaptation step (27) is the LMS solution for minimizing the cost in (10) with respect to  $\mathbf{w}_k$ .

Before leaving this section, we would like to point out that the algorithm described in [31], which addresses direct models by stacking global and local variables in an augmented parameter vector, may be used to solve problem (8), provided that an appropriate variable change is performed in order to make the latent variables  $\mathbf{u}_{k,n}$  and  $\xi_{k,n}$  explicit in  $\mathbf{w}_{k,n}$ . The resulting algorithm has the same performance as Algorithm 1 defined by (28), (29), but, obviously, they do not have the same form since they do not operate in the same domain. This structural difference has a major consequence for Algorithm 1. As already explained, it can be further tuned by replacing the matrix  $\mathbf{S}_\Theta$  in (28) by any positive definite matrix while ensuring convergence of the method. This extra degree of freedom will be taken into account in the analysis of the algorithm, where the only condition on  $\mathbf{S}_\Theta$  is to be positive definite. We will also show that setting  $\mathbf{S}_\Theta$  to  $\mathbf{I}_L$ , besides simplifying Algorithm 1, can greatly improve its performance.

### B. Node-specific subspace constraints with norm-bounded projections

The second formulation we consider is to relax the constraint that node-specific components  $\{\epsilon_k\}_{k=1}^N$  must lie in  $\text{span}(\Theta_\perp)$ . We now assume that they are norm-bounded in some sense. The problem is formulated as follows:

$$\begin{aligned} & \min_{\mathbf{u}, \{\epsilon_k\}_{k=1}^N} J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) \\ & \text{subject to } \sum_{k=1}^N \|\mathbf{P}_\Theta \epsilon_k\|^2 \leq \nu_1, \quad \sum_{k=1}^N \|\mathbf{P}_{\Theta_\perp} \epsilon_k\|^2 \leq \nu_2 \end{aligned} \quad (30)$$

---

**Algorithm 1:** ATC diffusion LMS with node-specific subspace constraints
 

---

**Parameters:** Preset

- positive step-size  $\mu$  for all agents;
- left-stochastic combination matrix  $\mathbf{A}$ ;
- full-rank matrix  $\Theta$  with columns  $\{\theta_1, \dots, \theta_M\}$ .

**Initialization:** Set initial weights  $\mathbf{w}_{k,0} = \mathbf{0}$  for all  $k \in \{1, \dots, N\}$ .

**Algorithm:** At each time instant  $n \geq 1$ , and for each agent  $k$ , update  $\mathbf{w}_{k,n}$  as:

$$\boldsymbol{\psi}_{k,n} = \mathbf{w}_{k,n-1} + \mu \mathbf{S}_\Theta \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1}] \quad (28)$$

$$\mathbf{w}_{k,n} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{P}_\Theta \boldsymbol{\psi}_{\ell,n} + \mathbf{P}_{\Theta_\perp} \boldsymbol{\psi}_{k,n} \quad (29)$$


---

Since the objective function and the constraints are convex in  $(\mathbf{u}, \{\epsilon_k\}_{k=1}^N)$ , the constrained problem (30) can be formulated as a regularized optimization problem that consists of minimizing a global cost of the form [49]:

$$J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) = \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n}(\Theta \mathbf{u} + \epsilon_k)|^2\} + \eta_1 \sum_{k=1}^N \|\mathbf{P}_\Theta \epsilon_k\|^2 + \eta_2 \sum_{k=1}^N \|\mathbf{P}_{\Theta_\perp} \epsilon_k\|^2 \quad (31)$$

where  $\eta_1$  and  $\eta_2$  are positive regularization parameters that are related to the bounds  $\nu_1$  and  $\nu_2$ .

*Lemma 2:* Problem (30) has a unique solution with respect to  $\mathbf{u}$  and  $\{\epsilon_k\}_{k=1}^N$ . ■

Proof of Lemma 2 is provided in Appendix B. Other norms such as the general  $\ell_{p,q}$ -norm may be used with  $\epsilon_k$  in (30), depending on the application. Some form of regularization on  $\mathbf{u}$  may also be included. However, using the  $\ell_2$ -norm with  $\epsilon_k$  in (30) enables us to solve the problem with respect to  $\mathbf{w}_k$ , without using the auxiliary variables  $\mathbf{u}$  and  $\{\epsilon_k\}_{k=1}^N$ . Indeed, let us rewrite (31) as follows:

$$J^{\text{glob}}(\mathbf{u}, \{\mathbf{w}_k\}_{k=1}^N) = \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_k|^2\} + \eta_1 \sum_{k=1}^N \|\mathbf{P}_\Theta(\mathbf{w}_k - \Theta \mathbf{u})\|^2 + \eta_2 \sum_{k=1}^N \|\mathbf{P}_{\Theta_\perp} \mathbf{w}_k\|^2 \quad (32)$$

The optimality condition relative to  $\mathbf{u}$  gives:

$$\sum_{k=1}^N \Theta^* \mathbf{P}_\Theta (\mathbf{w}_k^o - \Theta \mathbf{u}^o) = \mathbf{0} \quad (33)$$

from which the optimal parameter vector  $\mathbf{u}^o$  can be expressed as:

$$\mathbf{u}^o = \frac{1}{N} \sum_{k=1}^N (\Theta^* \Theta)^{-1} \Theta^* \mathbf{w}_k^o \quad (34)$$

Substituting (34) into (32), and using that  $\mathbf{P}_\Theta$  is Hermitian and idempotent (i.e.,  $\mathbf{P}_\Theta = \mathbf{P}_\Theta^2$ ), yields:

$$J^{\text{glob}}(\{\mathbf{w}_k\}_{k=1}^N) = \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_k|^2\} + \eta_1 \sum_{k=1}^N \left\| \mathbf{P}_\Theta \mathbf{w}_k - \frac{1}{N} \sum_{\ell=1}^N \mathbf{P}_\Theta \mathbf{w}_\ell \right\|^2 + \eta_2 \sum_{k=1}^N \|\mathbf{P}_{\Theta_\perp} \mathbf{w}_k\|^2 \quad (35)$$

Node  $k$  can apply a steepest-descent iteration to minimize the cost in (35) with respect to  $\{\mathbf{w}_k\}_{k=1}^N$ . Computing the gradient vector of (35) we get:

$$\nabla J^{\text{glob}} = \left[ (\mathbf{R}_{x,k} \mathbf{w}_k - \mathbf{p}_{dx}) + \eta_1 \left( \mathbf{P}_\Theta \mathbf{w}_k - \frac{1}{N} \sum_{\ell=1}^N \mathbf{P}_\Theta \mathbf{w}_\ell \right) + \eta_2 \mathbf{P}_{\Theta_\perp} \mathbf{w}_k \right]^* \quad (36)$$

Starting from an initial condition  $\mathbf{w}_{k,0}$ , we arrive at the steepest descent iteration:

$$\begin{aligned} \mathbf{w}_{k,n} &= \mathbf{w}_{k,n-1} - \mu \left[ (\mathbf{R}_{x,k} \mathbf{w}_{k,n-1} - \mathbf{p}_{dx}) + \eta_2 \mathbf{P}_{\Theta_{\perp}} \mathbf{w}_{k,n-1} \right] \\ &\quad - \mu \eta_1 \left( \mathbf{P}_{\Theta} \mathbf{w}_{k,n-1} - \frac{1}{N} \sum_{\ell=1}^N \mathbf{P}_{\Theta} \mathbf{w}_{\ell,n-1} \right) \end{aligned} \quad (37)$$

This iteration indicates that the update term involves adding two correction terms to  $\mathbf{w}_{k,n-1}$ . Among many other forms, we can implement the update in two successive steps by adding one correction term at a time:

$$\boldsymbol{\psi}_{k,n} = \mathbf{w}_{k,n-1} - \mu \left[ (\mathbf{R}_{x,k} \mathbf{w}_{k,n-1} - \mathbf{p}_{dx}) + \eta_2 \mathbf{P}_{\Theta_{\perp}} \mathbf{w}_{k,n-1} \right] \quad (38)$$

$$\mathbf{w}_{k,n} = \boldsymbol{\psi}_{k,n} - \mu \eta_1 \left( \mathbf{P}_{\Theta} \mathbf{w}_{k,n-1} - \frac{1}{N} \sum_{\ell=1}^N \mathbf{P}_{\Theta} \mathbf{w}_{\ell,n-1} \right) \quad (39)$$

Step (38) updates  $\mathbf{w}_{k,n-1}$  to an intermediate value  $\boldsymbol{\psi}_{k,n}$ . We now revise (38)–(39) to achieve a diffusion LMS type algorithm. The intermediate value  $\boldsymbol{\psi}_{\ell,n}$  at node  $\ell$  is generally expected to be a better estimate for  $\mathbf{w}_{\ell}^o$  than  $\mathbf{w}_{\ell,n-1}$  since it is updated by the first step (38). Therefore, we replace  $\mathbf{w}_{\ell,n-1}$  by  $\boldsymbol{\psi}_{\ell,n}$  in the second step (39) as follows to get:

$$\begin{aligned} \mathbf{w}_{k,n} &= \boldsymbol{\psi}_{k,n} - \mu \eta_1 \left( \mathbf{P}_{\Theta} \boldsymbol{\psi}_{k,n} - \frac{1}{N} \sum_{\ell=1}^N \mathbf{P}_{\Theta} \boldsymbol{\psi}_{\ell,n} \right) \\ &= (\boldsymbol{\psi}_{k,n} - \mathbf{P}_{\Theta} \boldsymbol{\psi}_{k,n}) + \left( (1 - \mu \eta_1) \mathbf{P}_{\Theta} \boldsymbol{\psi}_{k,n} + \sum_{\ell=1}^N \frac{\mu \eta_1}{N} \boldsymbol{\psi}_{\ell,n} \right) \end{aligned} \quad (40)$$

Observe that  $\mathbf{P}_{\Theta_{\perp}} \boldsymbol{\psi}_{k,n} = \boldsymbol{\psi}_{k,n} - \mathbf{P}_{\Theta} \boldsymbol{\psi}_{k,n}$ , and introduce the coefficients  $a_{\ell k} = \frac{\mu \eta_1}{N}$  for  $\ell \neq k$ , and  $a_{kk} = 1 - \mu \eta_1 + \frac{\mu \eta_1}{N}$ . We get:

$$\mathbf{w}_{k,n} = \mathbf{P}_{\Theta_{\perp}} \boldsymbol{\psi}_{k,n} + \sum_{\ell=1}^N a_{\ell k} \mathbf{P}_{\Theta} \boldsymbol{\psi}_{\ell,n} \quad (41)$$

Considering that each node in the network can only share information with its neighbors, and using instantaneous approximations for  $\mathbf{R}_{x,k}$  and  $\mathbf{p}_{dx}$ , we arrive at:

$$\boldsymbol{\psi}_{k,n} = (\mathbf{I}_L - \mu \eta_2 \mathbf{P}_{\Theta_{\perp}}) \mathbf{w}_{k,n-1} + \mu \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1}] \quad (42)$$

$$\mathbf{w}_{k,n} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{P}_{\Theta} \boldsymbol{\psi}_{\ell,n} + \mathbf{P}_{\Theta_{\perp}} \boldsymbol{\psi}_{k,n} \quad (43)$$

with  $a_{kk} = 1 - \mu \eta_1 + \frac{\mu \eta_1}{|\mathcal{N}_k|}$  and  $a_{\ell k} = \frac{\mu \eta_1}{|\mathcal{N}_k|}$  for  $\ell \in \mathcal{N}_k$  and  $\ell \neq k$ . Note that, for sufficiently small step-sizes  $\mu_k$ , these coefficients are nonnegative and satisfy  $\sum_{\ell=1}^N a_{\ell k} = 1$  for all  $k$ . We will treat these coefficients as free parameters that can be chosen by the designer according to these conditions (i.e., nonnegative coefficients that add up to one on each column of matrix  $\mathbf{A}$ ). We summarize this statement in Algorithm 2.

Algorithms 1 and 2 employ the same aggregation step in (29) and (45). Node  $k$  combines the intermediate estimates of its neighbors in the common subspace  $\Theta$  without affecting the local contribution in the complementary subspace  $\Theta_{\perp}$ . The norm constraint (30) in  $\Theta_{\perp}$  leads to a leaky-LMS alike term in the adaptation step (44).

---

**Algorithm 2:** ATC diffusion LMS with node-specific subspace constraints (norm-bounded projections)

---

**Parameters:** Preset

- positive step-size  $\mu$  for all agents;
- full-rank matrix  $\Theta$  with columns  $\{\theta_1, \dots, \theta_M\}$ .

**Initialization:** Set initial weights  $\mathbf{w}_{k,0} = \mathbf{0}$  for all  $k = 1, \dots, N$ .

**Algorithm:** For each instant  $n \geq 1$ , and for each agent  $k$ , update  $\mathbf{w}_{k,n-1}$ :

$$\boldsymbol{\psi}_{k,n} = (\mathbf{I}_L - \mu\eta_2 \mathbf{P}_{\Theta^\perp}) \mathbf{w}_{k,n-1} + \mu \mathbf{x}_{k,n}^* [d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1}] \quad (44)$$

$$\mathbf{w}_{k,n} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{P}_\Theta \boldsymbol{\psi}_{\ell,n} + \mathbf{P}_{\Theta^\perp} \boldsymbol{\psi}_{k,n} \quad (45)$$


---

Let us now examine two special cases of Algorithm 2. First, in the case where  $\Theta = \mathbf{0}$ , problem (31) reduces to a regularized least-mean squares problem with  $\mathbf{w}_k = \boldsymbol{\epsilon}_k$ . That is, the algorithm reduces to the non-cooperative leaky-LMS algorithm. On the other hand, if  $\Theta = \mathbf{I}_L$ , the algorithm reduces to diffusion LMS.

Before leaving this section, we briefly discuss the complexity of Algorithms 1 and 2. Both algorithms have the same adapt-then-combine structure as the diffusion LMS except that each node needs to project data on  $\Theta$  and  $\Theta^\perp$ . This means that each node  $k$  only needs to update the  $L \times 1$  parameter vectors  $\boldsymbol{\psi}_{k,n}$  and  $\mathbf{w}_{k,n}$  at each time instant. Next, each node  $k$  needs to transmit  $\mathbf{w}_{k,n}$  to its  $|\mathcal{N}_k| - 1$  neighbors. A projection performed by a matrix-vector product has a computational complexity of  $\mathcal{O}(L \log_2 L)$  [50]. All the other operations performed by each node have a complexity of  $\mathcal{O}(L)$ .

#### IV. PERFORMANCE AND CONVERGENCE ANALYSES

In this section, we examine the convergence properties and network performance of the proposed adaptive strategies. We shall first describe a convergence framework for a family of distributed algorithms, where Algorithms 1 and 2 are special cases. Quantities specifically related to Algorithms 1 or 2 will be distinguished by superscripts <sup>(1)</sup> and <sup>(2)</sup>, respectively.

In order to perform the analysis, we collect information from across the network into block vectors and matrices. Let us denote by  $\mathbf{w}_n$  and  $\mathbf{w}^o$  the block weight vector at instant  $n$  and the block optimum weight vector, both of size  $LN \times 1$ , that is

$$\mathbf{w}_n = \text{col}\{\mathbf{w}_{1,n}, \dots, \mathbf{w}_{N,n}\} \quad (46)$$

$$\mathbf{w}^o = \text{col}\{\mathbf{w}_1^o, \dots, \mathbf{w}_N^o\} \quad (47)$$

We denote the difference between the optimum  $\mathbf{w}_k^o$  and the instantaneous estimate  $\mathbf{w}_{k,n}$  by:

$$\mathbf{v}_{k,n} = \mathbf{w}_k^o - \mathbf{w}_{k,n} \quad (48)$$

We collect the weight error vectors  $\mathbf{v}_{k,n}$  from across all nodes into the block weight error vector:

$$\mathbf{v}_n = \text{col}\{\mathbf{v}_{1,n}, \dots, \mathbf{v}_{N,n}\} \quad (49)$$

*Assumption 1:* (Independent inputs) The regression vectors  $\mathbf{x}_{k,n}$  arise from a stationary random process that is temporally stationary, white, and independent over space with  $\mathbf{R}_{x,k} = \mathbb{E}\{\mathbf{x}_{k,n}^* \mathbf{x}_{k,n}\} > 0$ . A direct consequence of this condition is that  $\mathbf{x}_{k,n}$  is independent of  $\mathbf{v}_{\ell,m}$  for all  $\ell$  and  $m \leq n$ .

#### A. Mean weight behavior analysis

The estimation error in (28) and (44) can be rewritten as a function of  $\mathbf{v}_{k,n}$ :

$$d_k(n) - \mathbf{x}_{k,n} \mathbf{w}_{k,n-1} = z_k(n) + \mathbf{x}_{k,n} \mathbf{v}_{k,n-1} \quad (50)$$

In what follows, we first show that the weight error update relations for both Algorithms 1 and 2 are of the form:

$$\mathbf{v}_n = \mathbf{B}_n \mathbf{v}_{n-1} - \mu \mathbf{g}_n - \mathbf{r}, \quad (51)$$

with  $\mathbf{B}_n$  an  $LN \times LN$  time-dependent matrix,  $\mathbf{g}_n$  an  $LN \times 1$  zero-mean time-dependent vector, and  $\mathbf{r}$  a constant  $LN \times 1$  vector. Consequently, it will be possible to represent their mean weight behavior in the form of a state-transition equation with a bounded driving term:

$$\mathbb{E}\{\mathbf{v}_n\} = \mathbf{B} \mathbb{E}\{\mathbf{v}_{n-1}\} - \mathbf{r} \quad (52)$$

with  $\mathbf{B} = \mathbb{E}\{\mathbf{B}_{n-1}\}$ . Let  $\mathbf{H}_{x,n}$  be the block diagonal matrix of size  $LN \times LN$ , and  $\mathbf{p}_{zx,n}$  the vector of length  $LN \times 1$ , defined as follows:

$$\mathbf{H}_{x,n} \triangleq \text{diag}\{\mathbf{x}_{1,n}^* \mathbf{x}_{1,n}, \dots, \mathbf{x}_{N,n}^* \mathbf{x}_{N,n}\} \quad (53)$$

$$\mathbf{p}_{zx,n} \triangleq \text{col}\{z_1(n) \mathbf{x}_{1,n}^*, \dots, z_N(n) \mathbf{x}_{N,n}^*\} \quad (54)$$

The expectation of  $\mathbf{H}_{x,n}$  and  $\mathbf{p}_{zx,n}$  are given by:

$$\mathbf{H}_x \triangleq \mathbb{E}\{\mathbf{H}_{x,n}\} = \text{diag}\{\mathbf{R}_{x,1}, \dots, \mathbf{R}_{x,2}\} \quad (55)$$

$$\mathbf{p}_{zx} \triangleq \mathbb{E}\{\mathbf{p}_{zx,n}\} = \mathbf{0} \quad (56)$$

1) *Mean weight behavior of Algorithm 1:* Define the intermediate weight error vector  $\tilde{\boldsymbol{\psi}}_{k,n}$ :

$$\tilde{\boldsymbol{\psi}}_{k,n} = \mathbf{w}_k^o - \boldsymbol{\psi}_{k,n} \quad (57)$$

and collect these vectors from across all nodes into the block weight error vector:

$$\tilde{\boldsymbol{\psi}}_n = \text{col}\{\tilde{\boldsymbol{\psi}}_{1,n}, \dots, \tilde{\boldsymbol{\psi}}_{N,n}\} \quad (58)$$

Subtracting  $\mathbf{w}_k^o$  from both sides of the update relation (28), and using relation (50), leads to the update equation for  $\tilde{\boldsymbol{\psi}}_n$ :

$$\tilde{\boldsymbol{\psi}}_n = (\mathbf{I}_{LN} - \mu \mathbf{D}_{S_{\ominus}} \mathbf{H}_{x,n}) \mathbf{v}_{n-1} - \mu \mathbf{D}_{S_{\ominus}} \mathbf{p}_{zx,n} \quad (59)$$

where  $D_{S_\Theta} = \text{diag}\{S_\Theta, \dots, S_\Theta\}$  is an  $LN \times LN$  block diagonal matrix with  $S_\Theta$  as diagonal entries. Let  $\mathcal{A} = A \otimes I_L$ . Defining  $D_{P_\Theta}$  and  $D_{P_{\Theta_\perp}}$  as the  $LN \times LN$  block diagonal matrices with  $P_\Theta$  and  $P_{\Theta_\perp}$  as diagonal entries, respectively, equation (29) can be written in vector form as:

$$\mathbf{w}_n = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) \boldsymbol{\psi}_n \quad (60)$$

Subtracting  $\mathbf{w}^o$  from both sides of the above expression, we have:

$$\mathbf{v}_n = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) \tilde{\boldsymbol{\psi}}_n - (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}} - I_{LN}) \mathbf{w}^o \quad (61)$$

Combining this equation with (59), the weight error update relation can be written in a single expression:

$$\mathbf{v}_n = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) [(I_{LN} - \mu D_{S_\Theta} \mathbf{H}_{x,n}) \mathbf{v}_{n-1} - \mu D_{S_\Theta} \mathbf{p}_{zx,n}] - (\mathcal{A}^\top - I_{LN}) D_{P_\Theta} \mathbf{w}^o \quad (62)$$

Now we denote several terms in the weight error expression (62) by:

$$\mathbf{B}_n^{(1)} = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) (I_{LN} - \mu D_{S_\Theta} \mathbf{H}_{x,n}) \quad (63)$$

$$\mathbf{g}_n^{(1)} = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) D_{S_\Theta} \mathbf{p}_{zx,n} \quad (64)$$

$$\mathbf{r}^{(1)} = (\mathcal{A}^\top - I_{LN}) D_{P_\Theta} \mathbf{w}^o, \quad (65)$$

and the associated expected values:

$$\begin{aligned} \mathbf{B}^{(1)} &\triangleq \mathbb{E}\{\mathbf{B}_n^{(1)}\} \\ &= (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) (I_{LN} - \mu D_{S_\Theta} \mathbf{H}_x) \end{aligned} \quad (66)$$

$$\mathbf{g}^{(1)} \triangleq \mathbb{E}\{\mathbf{g}_n^{(1)}\} = \mathbf{0} \quad (67)$$

With the above notation, the weight error update relation (62) can be written as:

$$\mathbf{v}_n = \mathbf{B}_n^{(1)} \mathbf{v}_{n-1} - \mu \mathbf{g}_n^{(1)} - \mathbf{r}^{(1)} \quad (68)$$

Taking the expectation on both sides of (68), and using Assumption 1, we arrive at the mean weight behavior for Algorithm 1:

$$\mathbb{E}\{\mathbf{v}_n\} = \mathbf{B}^{(1)} \mathbb{E}\{\mathbf{v}_{n-1}\} - \mathbf{r}^{(1)} \quad (69)$$

2) *Mean weight behavior of Algorithm 2:* Subtracting  $\mathbf{w}_k^o$  from both sides of the update relation (44), and using relation (50), yields:

$$\tilde{\boldsymbol{\psi}}_n = (\mathbf{I} - \mu \eta_2 D_{P_{\Theta_\perp}} - \mu \mathbf{H}_{x,n}) \mathbf{v}_{n-1} - \mu (\mathbf{p}_{zx,n} - \eta_2 D_{P_{\Theta_\perp}} \mathbf{w}^o) \quad (70)$$

Subtracting  $\mathbf{w}^o$  from both sides of (45), we have:

$$\mathbf{v}_n = (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) \tilde{\boldsymbol{\psi}}_n - (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}} - I_{LN}) \mathbf{w}^o \quad (71)$$

Combining this equation with (70), the weight error update relation can be written in a single expression:

$$\begin{aligned} \mathbf{v}_n &= (\mathcal{A}^\top D_{P_\Theta} + D_{P_{\Theta_\perp}}) [(\mathbf{I}_{LN} - \mu \eta_2 D_{P_{\Theta_\perp}} - \mu \mathbf{H}_{x,n}) \mathbf{v}_{n-1} \\ &\quad - \mu (\mathbf{p}_{zx,n} - \eta_2 D_{P_{\Theta_\perp}} \mathbf{w}^o)] - (\mathcal{A}^\top - I_{LN}) D_{P_\Theta} \mathbf{w}^o \end{aligned} \quad (72)$$

where we used the fact that  $\mathbf{I}_{LN} = \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}}$ . Next, we denote several terms in the weight error expression (72) by:

$$\mathbf{B}_n^{(2)} = (\mathcal{A}^\top \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}})(\mathbf{I}_{LN} - \mu\eta_2 \mathbf{D}_{\mathbf{P}_{\Theta_\perp}} - \mu \mathbf{H}_{x,n}) \quad (73)$$

$$\mathbf{g}_n^{(2)} = (\mathcal{A}^\top \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}}) \mathbf{p}_{zx,n} \quad (74)$$

$$\mathbf{r}^{(2)} = (\mathcal{A}^\top - \mathbf{I}_{LN}) \mathbf{D}_{\mathbf{P}_\Theta} \mathbf{w}^o - \mu\eta_2 (\mathcal{A}^\top \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}}) \mathbf{D}_{\mathbf{P}_{\Theta_\perp}} \mathbf{w}^o \quad (75)$$

and the associated expected values:

$$\begin{aligned} \mathbf{B}^{(2)} &\triangleq \mathbb{E}\{\mathbf{B}_n^{(2)}\} \\ &= (\mathcal{A}^\top \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}})(\mathbf{I}_{LN} - \mu\eta_2 \mathbf{D}_{\mathbf{P}_{\Theta_\perp}} - \mu \mathbf{H}_x) \end{aligned} \quad (76)$$

$$\mathbf{g}^{(2)} \triangleq \mathbb{E}\{\mathbf{g}_n^{(2)}\} = \mathbf{0} \quad (77)$$

With the above notation, the weight error update relation (72) can be written as:

$$\mathbf{v}_n = \mathbf{B}_n^{(2)} \mathbf{v}_{n-1} - \mu \mathbf{g}_n^{(2)} - \mathbf{r}^{(2)} \quad (78)$$

Taking the expectation on both sides of (78), and using Assumption 1, we get the mean weight behavior of Algorithm 2:

$$\mathbb{E}\{\mathbf{v}_n\} = \mathbf{B}^{(2)} \mathbb{E}\{\mathbf{v}_{n-1}\} - \mathbf{r}^{(2)} \quad (79)$$

3) *Stability in the mean:* The mean-weight error recursions (69) and (79) are of the same form as (52). The convergence of such recursive state-transition equations, with a bounded driving term, is determined by the stability of matrix  $\mathbf{B}$ . Algorithm parameters should be chosen to satisfy the mean stability condition  $\rho(\mathbf{B}) < 1$ , where  $\rho(\cdot)$  denotes spectral radius of its matrix argument. In this case, the bias of the algorithm will be given by:

$$\lim_{n \rightarrow \infty} \mathbb{E}\{\mathbf{v}_n\} = -(\mathbf{I}_{LN} - \mathbf{B})^{-1} \mathbf{r} \quad (80)$$

We shall now establish two results that provide ranges for selecting the step size  $\mu$  to ensure convergence in the mean for each algorithm.

*Theorem 1:* (Stability in the mean for Algorithm 1) Assume data model (1) and Assumption 1 hold. We select a doubly stochastic matrix  $\mathbf{A}$ . Assume  $\{\Theta, \Theta_\perp\}$  forms an orthonormal basis of  $\mathbb{R}^L$ . Then, for any initial condition, Algorithm 1 asymptotically converges in the mean if the step-size satisfies:

$$0 < \mu < \frac{2}{\max_k \lambda_{\max}(\mathbf{R}_{x,k})} \quad (81)$$

where  $\lambda_{\max}(\cdot)$  denotes the maximum eigenvalue of its matrix argument.

*Proof:* The convergence of (69) is determined by the stability of matrix  $\mathbf{B}^{(1)}$ . The required mean stability condition is met by selecting  $\mu$  so that:

$$\rho((\mathcal{A}^\top \mathbf{D}_{\mathbf{P}_\Theta} + \mathbf{D}_{\mathbf{P}_{\Theta_\perp}})(\mathbf{I}_{LN} - \mu \mathbf{S}_\Theta \mathbf{H}_x)) < 1 \quad (82)$$



Let  $\mathbf{x} = \text{col}\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  be any block vector of size  $LN \times 1$ . We have:

$$\|(\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}})\mathbf{x}\|^2 = \sum_{i=1}^N \left\| \sum_{j=1}^N a_{ji} P_{\Theta} \mathbf{x}_j + P_{\Theta_\perp} \mathbf{x}_i \right\|^2 \quad (83)$$

$$= \sum_{i=1}^N \left( \left\| \sum_{j=1}^N a_{ji} P_{\Theta} \mathbf{x}_j \right\|^2 + \|P_{\Theta_\perp} \mathbf{x}_i\|^2 \right) \quad (84)$$

Given that  $\mathbf{A}$  is left stochastic, namely,  $\sum_{j=1}^N a_{ji} = 1$  with  $a_{ji} \geq 0$ , Jensen's inequality guarantees:

$$\left\| \sum_{j=1}^N a_{ji} P_{\Theta} \mathbf{x}_j \right\|^2 \leq \sum_{j=1}^N a_{ji} \|P_{\Theta} \mathbf{x}_j\|^2 \quad (85)$$

Consequently, the quantity in (84) can be upper-bounded as follows:

$$\sum_{i=1}^N \left( \left\| \sum_{j=1}^N a_{ji} P_{\Theta} \mathbf{x}_j \right\|^2 + \|P_{\Theta_\perp} \mathbf{x}_i\|^2 \right) \leq \sum_{i=1}^N \sum_{j=1}^N a_{ji} \|P_{\Theta} \mathbf{x}_j\|^2 + \sum_{i=1}^N \|P_{\Theta_\perp} \mathbf{x}_i\|^2 \quad (86)$$

$$\stackrel{(a)}{=} \sum_{j=1}^N \|P_{\Theta} \mathbf{x}_j\|^2 + \sum_{i=1}^N \|P_{\Theta_\perp} \mathbf{x}_i\|^2 \quad (87)$$

$$= \|\mathbf{x}\|^2 \quad (88)$$

where for step (a) we use that  $\mathbf{A}$  is right stochastic, namely,  $\sum_{i=1}^N a_{ji} = 1$ . We conclude that:

$$\|\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}}\| \leq 1 \quad (89)$$

We know that the spectral radius of any matrix  $\mathbf{X}$  satisfies  $\rho(\mathbf{X}) \leq \|\mathbf{X}\|$ , for any induced norm. Then we have:

$$\rho((\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}})(\mathbf{I}_{LN} - \mu \mathbf{S}_{\Theta} \mathbf{H}_x)) \leq \|\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}}\| \|\mathbf{I}_{LN} - \mu \mathbf{S}_{\Theta} \mathbf{H}_x\| \quad (90)$$

$$\stackrel{(89)}{\leq} \|\mathbf{I}_{LN} - \mu \mathbf{S}_{\Theta} \mathbf{H}_x\| \quad (91)$$

The mean stability condition is thus met by selecting  $\mu$  so that:  $\|\mathbf{I}_{LN} - \mu \mathbf{S}_{\Theta} \mathbf{H}_x\| < 1$ . In the case where  $\{\Theta, \Theta_\perp\}$  forms an orthonormal basis of  $\mathbb{R}^L$ , then  $\mathbf{S}_{\Theta} = \mathbf{I}_L$ . This leads us to the condition in (81). ■

*Theorem 2:* (Stability in the mean for Algorithm 2) Assume data model (1) and Assumption 1 hold. We select a doubly stochastic matrix  $\mathbf{A}$ . Then, for any initial condition, Algorithm 2 asymptotically converges in the mean if the step-size satisfies:

$$0 < \mu < \frac{2}{\max_k \lambda_{\max}(\eta_2 \mathbf{P}_{\Theta_\perp} + \mathbf{R}_{x,k})} \quad (92)$$

*Proof:* The convergence of (79) is determined by the stability of matrix  $\mathbf{B}^{(2)}$ . Considering that:

$$\rho((\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}})(\mathbf{I}_{LN} - \mu \eta_2 D_{P_{\Theta_\perp}} - \mu \mathbf{H}_x)) \leq \|\mathbf{I}_{LN} - \mu \eta_2 D_{P_{\Theta_\perp}} - \mu \mathbf{H}_x\| \quad (93)$$

since  $\|\mathcal{A}^\top D_{P_{\Theta}} + D_{P_{\Theta_\perp}}\| \leq 1$ , the mean stability condition is met by selecting  $\mu$  so that  $\|\mathbf{I}_{LN} - \mu \eta_2 D_{P_{\Theta_\perp}} - \mu \mathbf{H}_x\| < 1$ . This leads us to the condition in (92). Furthermore, by Weyl's theorem, we have  $\lambda_{\max}(\eta_2 \mathbf{P}_{\Theta_\perp} + \mathbf{R}_{x,k}) \leq$

$\eta_2 + \lambda_{\max}(\mathbf{R}_{x,k})$  since  $\mathbf{P}_{\Theta_{\perp}}$  and  $\mathbf{R}_{x,k}$  are Hermitian matrices and  $\lambda_{\max}(\mathbf{P}_{\Theta_{\perp}}) = 1$ . This leads to the sufficient condition:

$$0 < \mu < \frac{2}{\eta_2 + \max_k \lambda_{\max}(\mathbf{R}_{x,k})} \quad (94)$$

■

### B. Mean-square error behavior analysis

We now study the mean-square error behavior of Algorithms 1 and 2. To this end, we consider the general update relation (52) since both algorithms are of this form. From (51), the squared norm  $\|\mathbf{v}_n\|_{\Sigma}^2$  of the weight vector  $\mathbf{v}_n$  weighted by any positive semi-definite matrix  $\Sigma$ , i.e.,  $\|\mathbf{v}_n\|_{\Sigma}^2 = \mathbf{v}_n^* \Sigma \mathbf{v}_n$ , satisfies the following relation:

$$\|\mathbf{v}_n\|_{\Sigma}^2 = \|\mathbf{v}_{n-1}\|_{\mathbf{B}_n^* \Sigma \mathbf{B}_n}^2 - \mu^2 \|\mathbf{g}_n\|_{\Sigma}^2 + \|\mathbf{r}\|_{\Sigma}^2 - 2 \operatorname{Re}\{\mathbf{r}^* \Sigma \mathbf{B}_n \mathbf{v}_{n-1}\} - 2\mu \operatorname{Re}\{\mathbf{g}_n^* \Sigma (\mathbf{B}_n \mathbf{v}_{n-1} - \mathbf{r})\} \quad (95)$$

Under the independence assumption, and considering that  $\mathbf{g}_n$  includes the zero-mean noise term  $z_n$  which is independent of any other signal, taking expectations of both sides of (95) leads to:

$$\mathbb{E}\{\|\mathbf{v}_n\|_{\Sigma}^2\} = \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\Sigma'}^2\} + \mu^2 \operatorname{trace}\{\Sigma \mathbb{E}\{\mathbf{g}_n \mathbf{g}_n^*\}\} + \|\mathbf{r}\|_{\Sigma}^2 - 2 \operatorname{Re}\{\mathbb{E}\{\mathbf{r}^* \Sigma \mathbf{B}_n \mathbf{v}_{n-1}\}\} \quad (96)$$

In the above expression,  $\Sigma$  is any positive semi-definite matrix that the user is free to choose in order to derive different performance metrics, and  $\Sigma' = \mathbb{E}\{\mathbf{B}_n^* \Sigma \mathbf{B}_n\}$ . Let  $\mathbf{G}$  be the expected value of  $\mathbb{E}\{\mathbf{g}_n \mathbf{g}_n^*\}$  in the second term on the RHS of (96). For the two presented algorithms,  $\mathbf{G}$  is respectively given by:

$$\mathbf{G}^{(1)} = (\mathcal{A}^{\top} \mathbf{D}_{\mathbf{P}_{\Theta}} + \mathbf{D}_{\mathbf{P}_{\Theta_{\perp}}}) \mathbf{D}_{\mathbf{S}_{\Theta}} \operatorname{diag}\{\sigma_{z,1}^2 \mathbf{R}_{x,1}, \dots, \sigma_{z,N}^2 \mathbf{R}_{x,N}\} \mathbf{D}_{\mathbf{S}_{\Theta}}^* (\mathcal{A}^{\top} \mathbf{D}_{\mathbf{P}_{\Theta}} + \mathbf{D}_{\mathbf{P}_{\Theta_{\perp}}})^* \quad (97)$$

$$\mathbf{G}^{(2)} = \operatorname{diag}\{\sigma_{z,1}^2 \mathbf{R}_{x,1}, \dots, \sigma_{z,N}^2 \mathbf{R}_{x,N}\} \quad (98)$$

With  $\mathbf{G}$ , equation (96) is expressed as:

$$\begin{aligned} \mathbb{E}\{\|\mathbf{v}_n\|_{\Sigma}^2\} &= \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\Sigma'}^2\} + \mu^2 \operatorname{trace}\{\Sigma \mathbf{G}\} + \|\mathbf{r}\|_{\Sigma}^2 \\ &\quad - 2 \operatorname{Re}\{\mathbf{r}^* \Sigma \mathbf{B} \mathbb{E}\{\mathbf{v}_{n-1}\}\} \end{aligned} \quad (99)$$

Vectorizing matrices  $\Sigma$  and  $\Sigma'$  by  $\boldsymbol{\sigma} = \operatorname{vec}(\Sigma)$  and  $\boldsymbol{\sigma}' = \operatorname{vec}(\Sigma')$ , it can be verified that:

$$\boldsymbol{\sigma}' = \mathbf{K} \boldsymbol{\sigma} \quad (100)$$

where the  $(LN)^2 \times (LN)^2$  matrix  $\mathbf{K}$  is given by:

$$\mathbf{K} = \mathbb{E}\{\mathbf{B}_n^{\top} \otimes \mathbf{B}_n^*\} \approx \mathbf{B}^{\top} \otimes \mathbf{B}^* \quad (101)$$

The above approximation can be used provided that the step size is sufficiently small so that the influence of the second-degree term in  $\mu$  can be neglected [13]. Equation (99) can then be expressed as:

$$\mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\} = \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\mathbf{K}\boldsymbol{\sigma}}^2\} + \mathbf{s}_{n-1}^{\top} \boldsymbol{\sigma} \quad (102)$$

where we use the notation  $\|\cdot\|_{\Sigma}$  and  $\|\cdot\|_{\boldsymbol{\sigma}}$  interchangeably, and

$$\mathbf{s}_{n-1} = \operatorname{vec}\left(\mu^2 \mathbf{G} + \mathbf{r} \mathbf{r}^* - 2 \operatorname{Re}\{\mathbf{B} \mathbb{E}\{\mathbf{v}_{n-1}\} \mathbf{r}^*\}\right) \quad (103)$$

*Theorem 3:* (Mean-square stability) Assume data model (1) and Assumption 1 hold. Assume further that the step size  $\mu$  is sufficiently small to guarantee the stability in the mean of the algorithms, and to ensure that (102) can be used as a reasonable representation for the evolution of the (weighted) mean-square error. Mean-square stability of cooperative algorithms characterized by (51) requires the step-size  $\mu$  to be chosen such that it ensures the stability of matrix  $\mathbf{K}$  (in addition to the mean stability condition  $\rho(\mathbf{B}) < 1$ ).

*Proof:* Iterating (102) starting from  $n = 0$ , we find that

$$\mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\} = \|\mathbf{v}_0\|_{\mathbf{K}^n \boldsymbol{\sigma}}^2 + \sum_{i=1}^n \mathbf{s}_{n-i}^\top \mathbf{K}^{i-1} \boldsymbol{\sigma} \quad (104)$$

with initial condition  $\mathbf{v}_0 = \mathbf{w}^o - \mathbf{w}_0$ . Provided that  $\mathbf{K}$  is stable, the first term on the RHS of (104) converges to zero as  $n \rightarrow \infty$ . We know from (52) that  $\mathbb{E}\{\mathbf{v}_n\}$  is bounded because (52) is a BIBO stable recursion with a bounded driving term  $\mathbf{r}$ . The second term on the RHS of (104) then converges as  $n \rightarrow \infty$ . We conclude that  $\mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\}$  converges to a bounded value as  $n \rightarrow \infty$ , and the algorithm is mean-square stable. ■

*Theorem 4:* (Transient MSD) Consider a sufficiently small step size  $\mu$  to ensure mean and mean-square stabilities. The MSD learning curve  $\zeta_n = \frac{1}{N} \mathbb{E}\{\|\mathbf{v}_n\|^2\}$  of the cooperative algorithms characterized by (51), obtained by setting  $\boldsymbol{\Sigma} = \frac{1}{N} \mathbf{I}_{LN}$ , evolves according to the following recursion for  $n \geq 1$ :

$$\zeta_n = \zeta_{n-1} + \frac{1}{N} [(\boldsymbol{\gamma}_{n-1} + \mathbf{s}_{n-1})^\top \text{vec}(\mathbf{I}_{LN}) - \|\mathbf{v}_0\|_{(\mathbf{I}_{(LN)^2} - \mathbf{K})\mathbf{K}^{n-1}\boldsymbol{\sigma}}^2] \quad (105)$$

$$\boldsymbol{\gamma}_n = \mathbf{K}^\top \boldsymbol{\gamma}_{n-1} + (\mathbf{K} - \mathbf{I}_{(LN)^2})^\top \mathbf{s}_{n-1} \quad (106)$$

with initial conditions  $\zeta_0 = \frac{1}{N} \|\mathbf{v}_0\|^2$  and  $\boldsymbol{\gamma}_0 = \mathbf{0}$ .

*Proof:*

Comparing (104) at instants  $n$  and  $n - 1$ , we can relate  $\mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\}$  and  $\mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\boldsymbol{\sigma}}^2\}$  as follows:

$$\begin{aligned} \mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\} &= \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\boldsymbol{\sigma}}^2\} - \|\mathbf{v}_0\|_{(\mathbf{I}_{(LN)^2} - \mathbf{K})\mathbf{K}^{n-1}\boldsymbol{\sigma}}^2 + \sum_{i=1}^n \mathbf{s}_{n-i}^\top \mathbf{K}^{i-1} \boldsymbol{\sigma} - \sum_{i=1}^{n-1} \mathbf{s}_{n-1-i}^\top \mathbf{K}^{i-1} \boldsymbol{\sigma} \\ &= \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\boldsymbol{\sigma}}^2\} - \|\mathbf{v}_0\|_{(\mathbf{I}_{(LN)^2} - \mathbf{K})\mathbf{K}^{n-1}\boldsymbol{\sigma}}^2 + \mathbf{s}_{n-1}^\top \boldsymbol{\sigma} + \sum_{i=2}^n \mathbf{s}_{n-i}^\top \mathbf{K}^{i-1} \boldsymbol{\sigma} - \sum_{i=1}^{n-1} \mathbf{s}_{n-1-i}^\top \mathbf{K}^{i-1} \boldsymbol{\sigma} \end{aligned} \quad (107)$$

Introducing the notation

$$\boldsymbol{\gamma}_{n-1} = \left[ \sum_{i=2}^n \mathbf{s}_{n-i}^\top \mathbf{K}^{i-1} - \sum_{i=1}^{n-1} \mathbf{s}_{n-1-i}^\top \mathbf{K}^{i-1} \right]^\top \quad (108)$$

we can reformulate the recursive expression (107) as follows:

$$\mathbb{E}\{\|\mathbf{v}_n\|_{\boldsymbol{\sigma}}^2\} = \mathbb{E}\{\|\mathbf{v}_{n-1}\|_{\boldsymbol{\sigma}}^2\} - \|\mathbf{v}_0\|_{(\mathbf{I}_{(LN)^2} - \mathbf{K})\mathbf{K}^{n-1}\boldsymbol{\sigma}}^2 + (\boldsymbol{\gamma}_{n-1} + \mathbf{s}_{n-1})^\top \boldsymbol{\sigma} \quad (109)$$

$$\boldsymbol{\gamma}_n = \mathbf{K}^\top \boldsymbol{\gamma}_{n-1} + (\mathbf{K} - \mathbf{I}_{(LN)^2})^\top \mathbf{s}_{n-1} \quad (110)$$

with  $\boldsymbol{\gamma}_0 = \mathbf{0}$ . To derive the transient curve for the MSD, replace  $\boldsymbol{\sigma}$  by  $\frac{1}{N} \text{vec}\{\mathbf{I}_{LN}\}$ . ■

*Corollary 1:* (Steady-state MSD) If the step size is chosen sufficiently small to ensure mean and mean-square convergence, then the steady-state MSD, defined as  $\zeta_\infty = \lim_{n \rightarrow \infty} \zeta_n$ , is given by:

$$\zeta_\infty = \frac{1}{N} \mathbf{s}_\infty^\top (\mathbf{I}_{(LN)^2} - \mathbf{K})^{-1} \text{vec}(\mathbf{I}_{LN}) \quad (111)$$

with  $\mathbf{s}_\infty = \lim_{n \rightarrow \infty} \mathbf{s}_n$  determined by (103), using  $\mathbb{E}\{\mathbf{v}_\infty\} = \lim_{n \rightarrow \infty} \mathbb{E}\{\mathbf{v}_n\}$  determined by (80).

*Proof:*

From expression (102), we get:

$$\lim_{n \rightarrow \infty} \mathbb{E}\{\|\mathbf{v}_n\|_{(\mathbf{I}_{(LN)^2} - \mathbf{K})\boldsymbol{\sigma}}^2\} = \mathbf{s}_\infty^\top \boldsymbol{\sigma} \quad (112)$$

Observe that the MSD calculation requires us to choose  $\boldsymbol{\sigma}$  that satisfies:

$$(\mathbf{I}_{(LN)^2} - \mathbf{K}) \boldsymbol{\sigma} = \frac{1}{N} \text{vec}(\mathbf{I}_{LN}) \quad (113)$$

This leads to expression (111). ■

## V. SIMULATIONS

In this section, we report simulation results that illustrate the theoretical results. All agents were initialized with zero parameter vectors  $\mathbf{w}_{k,0} = \mathbf{0}$  for all  $k$ . Simulated curves were obtained by averaging over 100 runs as we obtained sufficiently smooth curves to check the consistency with theoretical results.

### A. Algorithm validation

We considered a network consisting of 12 agents with interconnections shown in Fig. 1(a). The parameter vectors to be estimated were of length  $L = 5$ . The input data  $\mathbf{x}_{k,n}$  were generated from circularly-symmetric zero-mean complex Gaussian distributions. White input data were considered first, by setting:

$$\mathbf{R}_{x,k} = \sigma_{x,k}^2 \mathbf{I}_5 \quad (114)$$

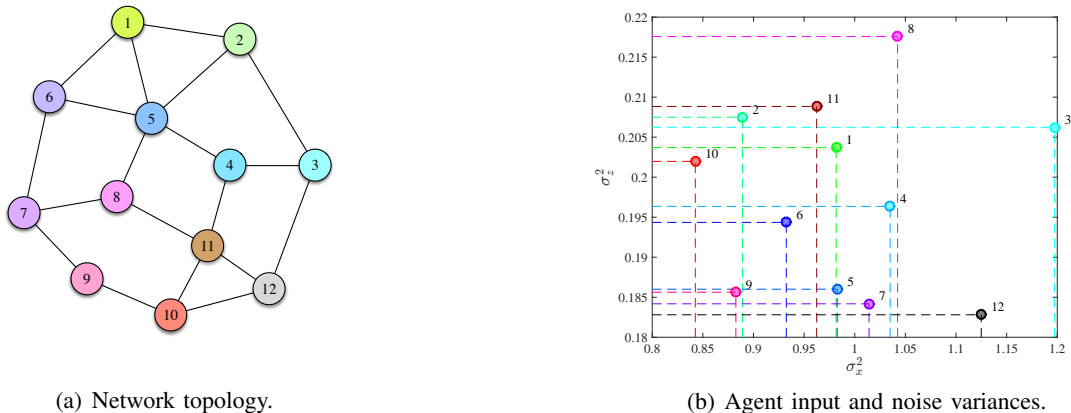
Next, correlated input data, characterized by the following covariance matrix, were considered:

$$\mathbf{R}_{x,k} = \sigma_{x,k}^2 \times \begin{pmatrix} 1 & -.4 + .3j & .2 - .1j & .1 - .05j & .02 + .02j \\ -.4 - .3j & 1 & -.4 + .3j & .2 - .1j & .1 - .05j \\ .2 + .1j & -.4 - .3j & 1 & -.4 + .3j & .2 - .1j \\ .1 + .05j & .2 + .1j & -.4 - .3j & 1 & -.4 + .3j \\ .02 - .02j & .1 + .05j & .2 + .1j & -.4 - .3j & 1 \end{pmatrix} \quad (115)$$

with  $j = \sqrt{-1}$  the imaginary unit. The modeling noises  $z_{k,n}$  were i.i.d. zero-mean circularly-symmetric Gaussian variables, independent of any other signals. The variances  $\sigma_{x,k}^2$  and  $\sigma_{z,k}^2$  were sampled from  $\mathcal{U}(0.8, 1.2)$  and  $\mathcal{U}(0.18, 0.22)$ , respectively. Their values are shown in Fig. 1(b). We considered two sets of subspace basis vectors. The first set is the standard basis:

$$\boldsymbol{\Theta}_1 = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_M], \quad (116)$$

where  $\mathbf{e}_i$  denotes a vector of length  $N$  with 1 at the  $i$ th entry and 0 otherwise. Its orthogonal complementary subspace is spanned by  $\boldsymbol{\Theta}_{1,\perp} = [\mathbf{e}_{M+1}, \dots, \mathbf{e}_L]$ . This setup can be interpreted as a variable selection process for



(a) Network topology.

(b) Agent input and noise variances.

Fig. 1. Network topology and input-noise variances.

information exchange, where the first  $M$  entries of the optimal parameter vectors are identical across the network. Parameter  $M$  was set to 3. The second set of basis vectors is a complex Vandermonde matrix:

$$\Theta_2 = \begin{pmatrix} 1 & 1 & \dots & 1 \\ e^{-j\psi_1} & e^{-j\psi_2} & \dots & e^{-j\psi_M} \\ \dots & \dots & \dots & \dots \\ e^{-j(L-1)\psi_1} & e^{-j(L-1)\psi_2} & \dots & e^{-j(L-1)\psi_M} \end{pmatrix} \quad (117)$$

with  $\psi_k = \frac{2\pi d}{\lambda_o} \sin \theta_k$ . Matrix  $\Theta_2$  can represent the array manifold of a uniform linear array (ULA) with inter-element space  $d$ , operating at wavelength  $\lambda_o$  with impinging signal directions of angles  $\theta_k$ . Parameter  $M$  was set to 3, with  $\theta_1 = \frac{\pi}{6}$ ,  $\theta_2 = \frac{\pi}{4}$ ,  $\theta_3 = \frac{\pi}{3}$  and  $d = \frac{\lambda_o}{2}$ . We considered three settings to validate the theoretical results.

In the first setting, we assumed that model (6) matches the observation data. The entries of the coefficient vectors  $\mathbf{u}^o$  and  $\xi_k^o$  were sampled from the Gaussian distribution  $\mathcal{N}(0, 1)$ . The step-size parameter  $\mu$  for Algorithm 1 was successively set to 0.01 and 0.02. A uniform combination matrix  $\mathbf{A}$  with  $a_{\ell k} = |\mathcal{N}_k|^{-1}$  was used. With  $\Theta_1$ , note that matrix  $\mathbf{S}_\Theta$  is equal to  $\mathbf{I}_5$ . With  $\Theta_2$ , it was successively set to  $\Theta\Theta^* + \Theta_\perp\Theta_\perp^*$  as in (29), and to  $\mathbf{I}_5$ . The transient behavior and the steady-state MSD were determined theoretically. The results with subspace settings  $\Theta_1$  and  $\Theta_2$ , for white and correlated input data, are shown in Fig. 2. It can be observed that setting  $\mathbf{S}_\Theta$  to  $\mathbf{I}_5$  for  $\Theta_2$  leads to a better convergence behavior. For Algorithm 2, we did not set the parameter  $\eta_1$  explicitly but we used the same combination matrix  $\mathbf{A}$  as for Algorithm 1. Parameters  $(\mu, \eta_2)$  were set to  $(0.02, 0.01)$  with white input data. With correlated input data, the following combinations  $(\mu, \eta_2)$  were considered:  $\{(0.01, 0.01); (0.01, 0.02); (0.02, 0.01)\}$ . The results are shown in Fig. 3. The simulation results match the theoretical results, and illustrate the trade-off between the convergence speed and the steady-state MSD. It can also be observed with Algorithm 2 that a small value for  $\eta_2$  is preferable since constraining the norm of node-specific components in the complementary subspace  $\Theta^\perp$  introduces a bias that can degrade the performance. As leaky-LMS, this kind of regularization can improve the stability of the algorithm for some particular problems and practical applications, at the cost of an extra estimation bias. We then considered another scenario in order to illustrate the interest of the extra degree of freedom provided

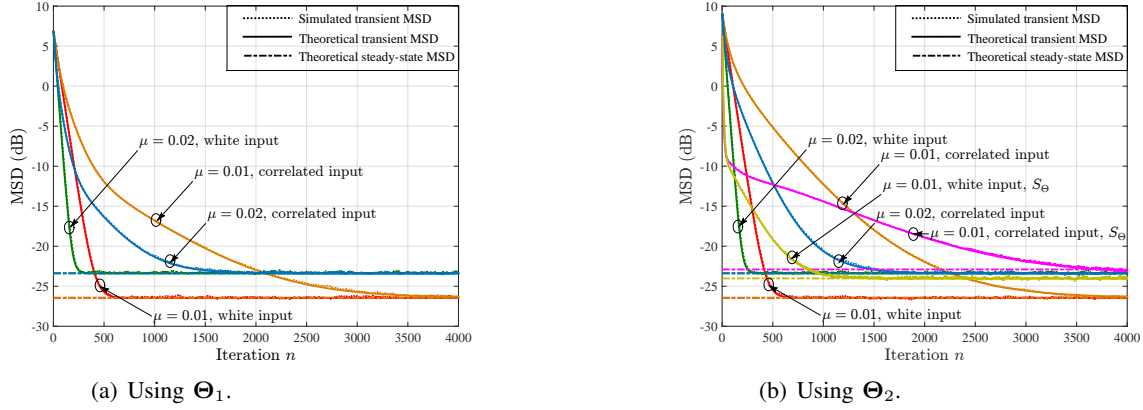


Fig. 2. Learning curves and model validation of Algorithm 1 with different settings.

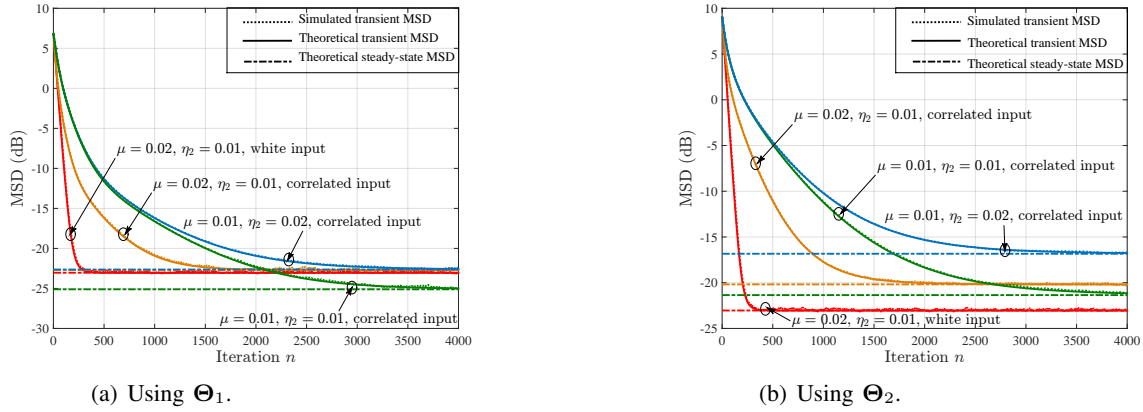


Fig. 3. Learning curves and model validation of Algorithm 2 with different settings.

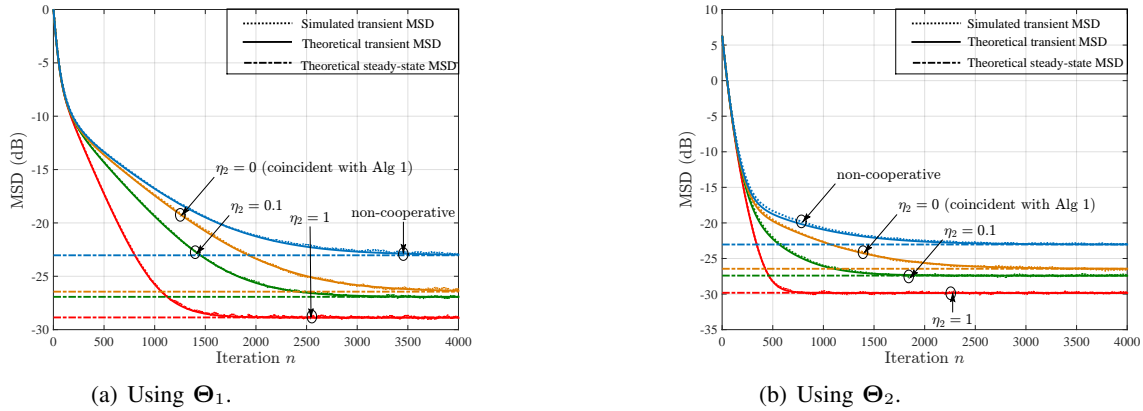


Fig. 4. Learning curves and model validation of the algorithms using  $\xi_k$  with small variances.

by  $\eta_2$  in Algorithm 2. Experimental setups were left unchanged with correlated inputs except for the entries of  $\xi_k^o$ , which were sampled from Gaussian distribution  $\mathcal{N}(0, 0.01)$ . We successively set  $\eta_2$  to 0, 0.1 and 1 in order to progressively constrain the variance of  $\xi_k$ . Note that with  $\eta_2 = 0$ , Algorithm 2 reduces to Algorithm 1. The results with  $\Theta_1$  and  $\Theta_2$  are provided in Fig. 4. The result with non-cooperative LMS is also provided as a reference.

In the second setting, we assumed that the node-specific components  $\epsilon_k^o$  in (6) do not strictly lie in the complementary subspace  $\Theta^\perp$ . To evaluate the robustness of our algorithms and the power of the analytical models, we

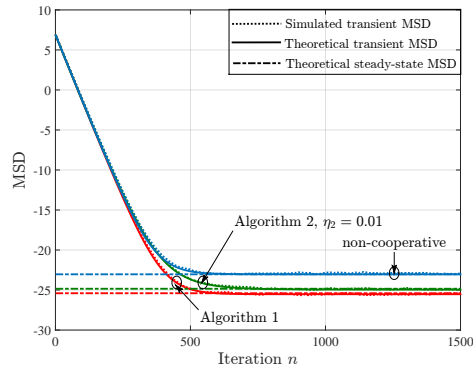


Fig. 5. Learning curves of the algorithms using  $\xi_k$  with small variances.

set:

$$\epsilon_k^o = \Theta \nu_k^o + \Theta_{\perp} \xi_k^o \quad (118)$$

where  $\nu_k^o$  are zero-mean circular Gaussian variables. This setting refers to a non-ideal situation because components  $\Theta(\mathbf{u}^o + \nu_k^o)$  lie in  $\text{span}(\Theta)$  but differ from one node to another. The entries of  $\mathbf{u}^o$  and  $\nu_k^o$  were sampled from Gaussian distributions  $\mathcal{N}(0, 1)$  and  $\mathcal{N}(0, 0.01)$ , respectively. The step-size  $\mu$  was set to 0.01 for Algorithms 1 and 2. Parameter  $\eta_2$  in Algorithm 2 was set to 0.01. Subspace  $\Theta_1$  and white input signals were considered to test the model. The transient behavior and the steady-state MSD were determined theoretically. The simulation results provided in Fig. 5 match the theoretical results, and illustrate that cooperation among nodes can still be beneficial when optimal solutions in the subspace  $\Theta$  are different but close to each other. This is another illustration of the conclusion reached in [19] for single-task diffusion LMS operating in multitask environments.

In the third setting, we exploited the leaky property of Algorithm 2 to promote its use in real applications. It is well known that the (non-cooperative) leaky LMS algorithm introduces an estimation bias compared to the (non-cooperative) LMS, but improves its robustness when applied to practical applications [51]. In particular, it avoids the so-called weight-drift problem of the LMS algorithm [52]. To highlight this phenomenon in the context of diffusion adaptation, we assumed that, say, the last tap/channel of node #1 was failing to work and was providing consistent null-valued readings, i.e.,  $[\mathbf{x}_{n,1}]_5 = 0$  for all  $n$ . We also assumed that, e.g., finite-precision effect was corrupting the combination step (29), or (45), with an additive non-zero mean disturbance  $\mathbf{q}_k$ . The poor conditioning of regressors associated with a non-zero mean disturbance is known to possibly lead to a weight-drift problem. We considered the same experimental setup as in the first experiment with the standard basis  $\Theta_1$ . We picked each entry of the random vectors  $\mathbf{q}_k$  according to the Gaussian distribution  $\mathcal{N}(10^{-4}, 10^{-8})$ . We set  $\eta_2$  to 0.1. All the vectors  $\mathbf{w}_k$  were initialized to  $\mathbf{0}$ . Fig. 6 shows the behavior of the weight vector at node #1 for (a) Algorithm 1 with  $\mathbf{S}_{\Theta} = \mathbf{I}_5$ , and (b) Algorithm 2. We can observe the drift of the 5<sup>th</sup> entry of  $\mathbf{w}_1$  with Algorithm 1. Algorithm 2 alleviates this effect.

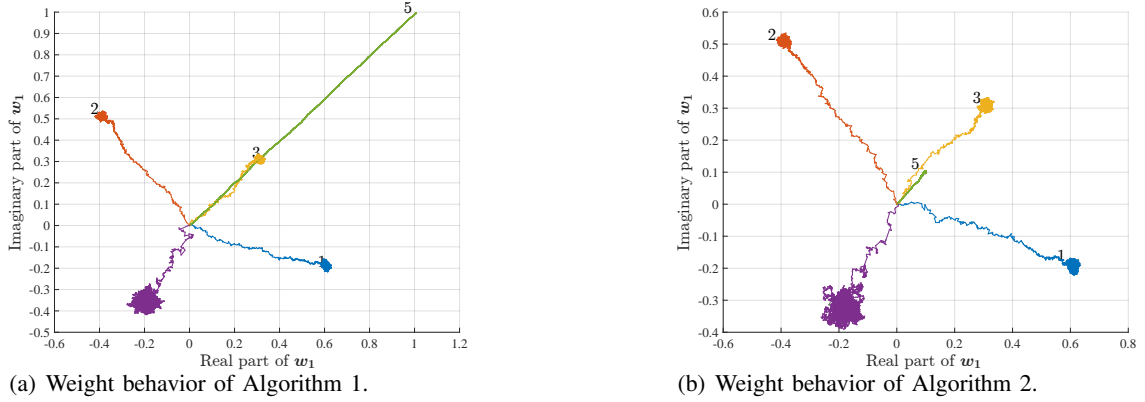


Fig. 6. Weight behavior of Algorithms 1 and 2 with singular inputs and non-zero mean disturbance.

### B. Target localization

We now consider a target localization problem. Cooperative localization with a diffusion strategy was already addressed in the case of a single target [13], and of multiple nearby targets [4]. We focus here on the case where targets lie in a manifold.

To make the presentation clearer, we assumed that the targets were collinear in  $\mathbb{R}^3$ . Their locations were estimated by the network with 100 nodes shown in Fig. 7(a). Each node randomly selected a target to localize. Let  $\mathcal{R}$  be a member of the rotation group  $SO(3)$  defined by the matrix  $\mathbf{R} = \mathbf{R}_x(\theta_x) \mathbf{R}_y(\theta_y) \mathbf{R}_z(\theta_z)$ , where  $\mathbf{R}_x(\theta_x)$ ,  $\mathbf{R}_y(\theta_y)$  and  $\mathbf{R}_z(\theta_z)$  are rotation matrices that rotate vectors by an angle of  $\theta_{x,y,z}$  around  $x$ ,  $y$  and  $z$  axis, respectively. The coordinate vector  $\mathbf{w}_q^o$  of each target  $q$  was generated as follows:

$$\mathbf{w}_q^o = \mathbf{R}_{1,2} \mathbf{u} + \epsilon_q \mathbf{r}_3 \quad (119)$$

where  $\mathbf{R}_{1,2}$  is the matrix composed of the first and second columns of  $\mathbf{R}$ , and  $\mathbf{r}_3$  corresponds to the third column of  $\mathbf{R}$ . As illustrated in Fig. 7(b), this model means that all targets lie on a common line defined by point  $\mathbf{R}_{1,2} \mathbf{u}$  and direction vector  $\mathbf{r}_3$ . Parameter  $\epsilon_q$  characterizes the location of each target  $q$  on this line. We considered the problem of estimating  $\mathbf{u}$  (common to all targets) and the parameters  $\epsilon_q$  for seven targets. We set the angles and the parameter vectors in (119) as follows:

$$\theta_x = \frac{\pi}{6}, \quad \theta_y = \frac{\pi}{3}, \quad \theta_z = \frac{\pi}{4} \quad (120)$$

$$\mathbf{v} = [1 \ 2]^\top \quad (121)$$

$$\epsilon_1 = 0, \epsilon_2 = 1, \epsilon_3 = 3, \epsilon_4 = 4, \epsilon_5 = 7, \epsilon_6 = 7.5, \epsilon_7 = 9 \quad (122)$$

The distance between each agent  $k$  and target  $q$  can be expressed in the inner product form:

$$r_{kq} = \mathbf{x}_{kq} (\mathbf{w}_q^o - \mathbf{p}_k) \quad (123)$$

where  $\mathbf{p}_k$  is the location of agent  $k$ , and  $\mathbf{x}_{kq}$  is the unit-norm row vector pointing from  $\mathbf{p}_k$  to  $\mathbf{w}_q^o$ . We assumed that agents were aware of their location  $\mathbf{p}_k$ . Let  $d_{kq} = r_{kq} + \mathbf{x}_{kq} \mathbf{p}_k$ , that is,  $d_{kq} = \mathbf{x}_{kq} \mathbf{w}_q^o$ . The problem was to



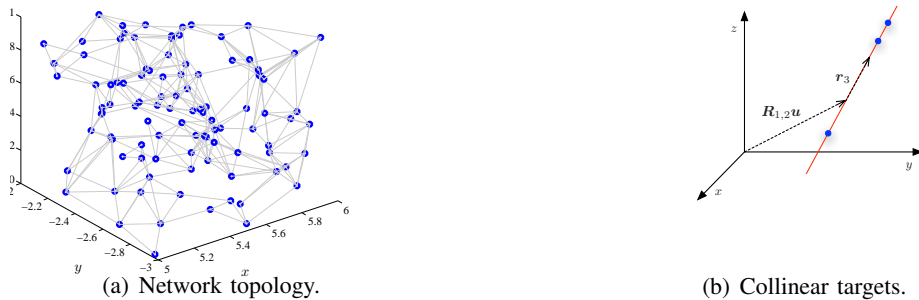


Fig. 7. Network topology and locations of targets.

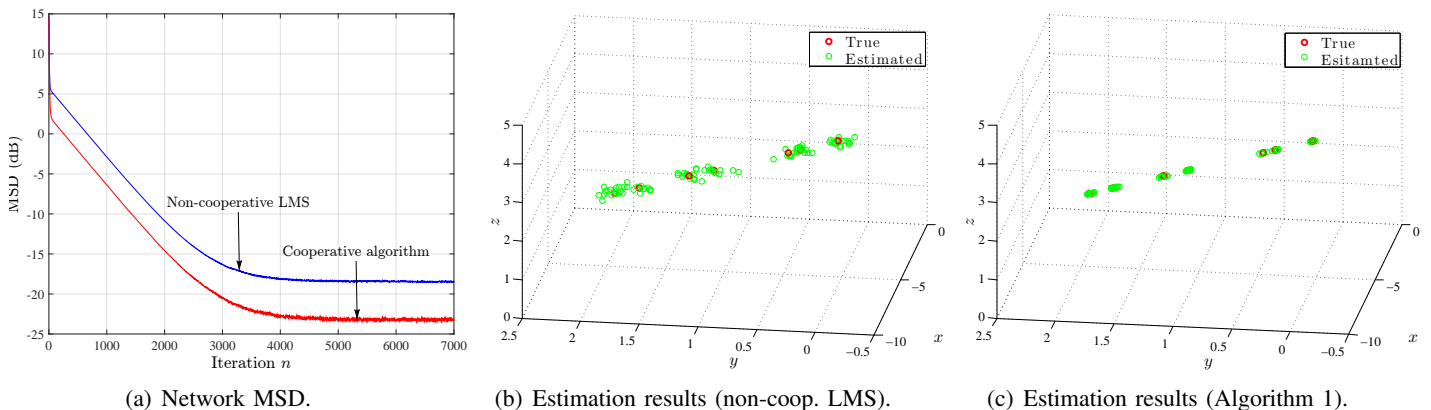


Fig. 8. Estimated network MSD and estimation results for a single realization.

estimate  $\mathbf{w}_q^o$  from noisy streaming measurements  $\{d_{kq}(n), \mathbf{x}_{kq,n}\}$  collected by each agent  $k$ , and governed by the linear model [13]:

$$d_{kq}(n) = \mathbf{x}_{kq,n} \mathbf{w}_q^o + z_{kq}(n)$$

with

$$\mathbf{x}_{kq,n} = [1 - \beta_k(n)] \mathbf{x}_{kq} + \mathbf{x}_{kq}^\perp \text{diag}\{\alpha_{k1}(n), \alpha_{k2}(n)\} \quad (124)$$

with  $z_{kq}(n)$  a zero-mean temporally and spatially i.i.d. Gaussian noise of variance  $\sigma_z^2$ . As shown in (124), the measured direction vector  $\mathbf{x}_{kq,n}$  was assumed to be a noisy realization of the unit-norm vector pointing from  $\mathbf{p}_k$  to  $\mathbf{w}_q^o$ , with  $\mathbf{x}_{kq}^\perp$  a unit-norm orthogonal contribution to  $\mathbf{x}_{kq}$ . Random variables  $\alpha_{k1}(n)$ ,  $\alpha_{k2}(n)$ ,  $\beta_k(n)$  and  $z_k(n)$  were zero-mean Gaussian with standard deviation  $\sigma_{\alpha_1} = \sigma_{\alpha_2} = 0.1$ ,  $\sigma_\beta = 0.001$  and  $\sigma_z = 0.3$ , respectively. We ran the (non-cooperative) LMS algorithm at each node, and Algorithm 1, with  $\Theta = \mathbf{R}_{1,2}$  and  $\Theta_\perp = \mathbf{r}_3$ . The step-size  $\mu$  was set to 0.1. A uniform combination matrix  $\mathbf{A}$  with  $a_{\ell k} = |\mathcal{N}_k|^{-1}$  was used for Algorithm 1, where  $|\mathcal{N}_k|$  denotes the cardinality of the neighborhood of node  $k$ . Figure 8(a) compares the MSD of these strategies. Figures 8(b) and 8(c) show one realization of the target locations estimated with the (noncooperative) LMS algorithm and Algorithm 1. This experiment illustrates the advantage of cooperative strategies over the non-cooperative one.

## VI. CONCLUSION AND PERSPECTIVES

In this paper, we formulated an online multitask adaptation problem that assumes that all tasks share a common latent feature representation, locally refined by node-specific contributions. This model can be extended into interesting directions by imposing new constraints, depending on applications. Based on this principle, we derived two cooperative algorithms and analyzed their performance. Although this work considers that common representation subspaces are known a priori, it paves the way towards more general frameworks.

## APPENDIX A

## PROOF OF LEMMA 1

The uniqueness of the solution of (8) follows from the strict convexity of (10), which is ensured by the positive definiteness of its Hessian matrix. For the quadratic cost (10), the Hessian matrix with respect to the vector of stacked variables  $\text{col}\{\mathbf{u}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_N\}$  is block diagonal [3, App. B], with blocks given by the following matrix  $\mathbf{X}$  and its transpose:

$$\nabla^2 J^{\text{glob}} = \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^\top \end{bmatrix}$$

with

$$\mathbf{X} = \left( \begin{array}{c|ccc} \boldsymbol{\Theta}^* (\sum_{k=1}^N \mathbf{R}_{x,k}) \boldsymbol{\Theta} & \boldsymbol{\Theta}^* \mathbf{R}_{x,1} \boldsymbol{\Theta}_\perp & \dots & \boldsymbol{\Theta}^* \mathbf{R}_{x,N} \boldsymbol{\Theta}_\perp \\ \hline \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,1} \boldsymbol{\Theta} & \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,1} \boldsymbol{\Theta}_\perp & & \mathbf{0} \\ \vdots & & \ddots & \\ \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,N} \boldsymbol{\Theta} & \mathbf{0} & & \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,N} \boldsymbol{\Theta}_\perp \end{array} \right) \quad (125)$$

where  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Theta}_\perp$  have full column rank. The positive definiteness of (125) can be checked by verifying the positive definiteness of each term  $\boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta}$  and of the Schur complement relative to the block diagonal corner of  $\mathbf{X}$ , namely, [53]

$$\text{Schur}(\mathbf{X}) = \sum_{k=1}^N [\boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta} - \boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta}_\perp (\boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,k} \boldsymbol{\Theta}_\perp)^{-1} \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,k} \boldsymbol{\Theta}] \quad (126)$$

where each inverse  $(\boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,k} \boldsymbol{\Theta}_\perp)^{-1}$  exists since  $\boldsymbol{\Theta}_\perp$  has full column rank. Each term inside the summation (126) is positive definite since it is the Schur complement of the block  $\boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta}$  in the positive definite matrix:

$$\begin{pmatrix} \boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta} & \boldsymbol{\Theta}^* \mathbf{R}_{x,k} \boldsymbol{\Theta}_\perp \\ \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,k} \boldsymbol{\Theta} & \boldsymbol{\Theta}_\perp^* \mathbf{R}_{x,k} \boldsymbol{\Theta}_\perp \end{pmatrix} = [\boldsymbol{\Theta} \ \boldsymbol{\Theta}_\perp]^* \mathbf{R}_{x,k} [\boldsymbol{\Theta} \ \boldsymbol{\Theta}_\perp] > 0 \quad (127)$$

This guarantees the positive definiteness of (126). It follows that the cost in (10) is strictly convex and has a unique minimizer.

APPENDIX B  
PROOF OF LEMMA 2

Without loss of generality, assume that  $\eta_1 > \eta_2$ . Otherwise, replace (129) by:

$$J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) = \sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n}(\Theta\mathbf{u} + \epsilon_k)|^2\} + \eta_1 \sum_{k=1}^N \|\epsilon_k\|^2 + (\eta_2 - \eta_1) \sum_{k=1}^N \|\mathbf{P}_{\Theta^\perp} \epsilon_k\|^2 \quad (128)$$

Recalling that  $\mathbf{P}_{\Theta^\perp} = \mathbf{I}_L - \mathbf{P}_\Theta$ , the objective function (31) can be written as follows:

$$J^{\text{glob}}(\mathbf{u}, \{\epsilon_k\}_{k=1}^N) = (\eta_1 - \eta_2) \sum_{k=1}^N \|\mathbf{P}_\Theta \epsilon_k\|^2 + \underbrace{\sum_{k=1}^N \mathbb{E}\{|d_k(n) - \mathbf{x}_{k,n}(\Theta\mathbf{u} + \epsilon_k)|^2\}}_{J_1^{\text{glob}}} + \eta_2 \sum_{k=1}^N \|\epsilon_k\|^2 \quad (129)$$

The uniqueness of the minimizer of (31) follows from its strict convexity. For the quadratic cost in (129), the Hessian of  $J_1^{\text{glob}}$  with respect to the vector of stacked variables  $\text{col}\{\mathbf{u}, \epsilon_1, \dots, \epsilon_N\}$  is again block diagonal, with its blocks determined by the matrix  $\mathbf{Y}$  below and its transpose:

$$\nabla^2 J_1^{\text{glob}} = \begin{bmatrix} \mathbf{Y} & \mathbf{0} \\ \mathbf{0} & \mathbf{Y}^\top \end{bmatrix} \quad (130)$$

with

$$\mathbf{Y} = \left( \begin{array}{c|ccc} \Theta^* (\sum_{k=1}^N \mathbf{R}_{x,k}) \Theta & \Theta^* \mathbf{R}_{x,1} & \dots & \Theta^* \mathbf{R}_{x,N} \\ \mathbf{R}_{x,1} \Theta & \mathbf{R}_{x,1} + \eta_2 \mathbf{I} & & \mathbf{0} \\ \vdots & & \ddots & \\ \mathbf{R}_{x,N} \Theta & \mathbf{0} & & \mathbf{R}_{x,N} + \eta_2 \mathbf{I} \end{array} \right) \quad (131)$$

The positive definiteness of (130) can be checked by verifying the positive definiteness of each term  $\mathbf{R}_{x,k} + \eta_2 \mathbf{I}$  and of the Schur complement relative to the right block diagonal corner in (130), namely, [53]

$$\text{Schur}(\mathbf{Y}) = \sum_{k=1}^N [\Theta^* \mathbf{R}_{x,k} \Theta - \Theta^* \mathbf{R}_{x,k} (\mathbf{R}_{x,k} + \eta_2 \mathbf{I})^{-1} \mathbf{R}_{x,k} \Theta] \quad (132)$$

Since they are positive definite, each covariance matrix  $\mathbf{R}_{x,k}$  can be decomposed as follows:

$$\mathbf{R}_{x,k} = \mathbf{U}_k \text{diag}\{\lambda_{k,1}, \dots, \lambda_{k,L}\} \mathbf{U}_k^* \quad (133)$$

where the  $\lambda_{k,i}$  are the eigenvalues of  $\mathbf{R}_{x,k}$ , which are real and positive, and  $\mathbf{U}_k$  is the corresponding matrix of eigenvectors. Since  $\mathbf{U}_k$  is an orthonormal matrix, each term in the summation (132) can be written as:

$$\begin{aligned} & \Theta^* \mathbf{R}_{x,k} \Theta - \Theta^* \mathbf{R}_{x,k} (\mathbf{R}_{x,k} + \eta_2 \mathbf{I})^{-1} \mathbf{R}_{x,k} \Theta \\ &= \Theta^* \mathbf{U} \text{diag} \left\{ \lambda_{k,1} - \frac{\lambda_{k,1}^2}{\lambda_{k,1} + \eta_2}, \dots, \lambda_{k,L} - \frac{\lambda_{k,L}^2}{\lambda_{k,L} + \eta_2} \right\} \mathbf{U}^* \Theta > 0 \end{aligned} \quad (134)$$

Since  $\Theta$  has full column rank, the above matrix and the Schur complement (132) are positive definite. In addition, the block diagonal matrix  $\text{diag}\{\mathbf{R}_{x,1} + \eta_2 \mathbf{I}, \dots, \mathbf{R}_{x,N} + \eta_2 \mathbf{I}\}$  is positive definite. Finally, since  $(\eta_1 - \eta_2) \sum_{k=1}^N \|\mathbf{P}_\Theta \epsilon_k\|^2$  in (129) is convex, problem (129) is strictly convex and problem (30) has a unique solution.

## REFERENCES

- [1] J. Chen, C. Richard, A. O. Hero, and A. H. Sayed, "Diffusion LMS for multitask problems with overlapping hypothesis subspaces," in *Proc. IEEE Int. Workshop on Machine Learn. for Signal Process. (MLSP)*, Reims, France, Sept. 2014, pp. 1–6.
- [2] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, 2014.
- [3] A. H. Sayed, "Adaptation, learning, and optimization over networks," in *Foundations and Trends in Machine Learning*, vol. 7, pp. 311–801. NOW Publishers, Boston-Delft, Jul. 2014.
- [4] 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.
- [5] D. P. Bertsekas, "A new class of incremental gradient methods for least squares problems," *SIAM J. Optimiz.*, vol. 7, no. 4, pp. 913–926, Nov. 1997.
- [6] M. G. Rabbat and R. D. Nowak, "Quantized incremental algorithms for distributed optimization," *IEEE J. of Sel. Topics Areas Commun.*, vol. 23, no. 4, pp. 798–808, Apr. 2005.
- [7] D. Blatt, A. O. Hero, and H. Gauchman, "A convergent incremental gradient method with constant step size," *SIAM J. Optimiz.*, vol. 18, no. 1, pp. 29–51, Feb. 2007.
- [8] C. G. Lopes and A. H. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Trans. Signal Process.*, vol. 55, no. 8, pp. 4064–4077, Aug. 2007.
- [9] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.
- [10] S. Kar and J. M. F. Moura, "Distributed consensus algorithms in sensor networks: Link failures and channel noise," *IEEE Trans. Signal Process.*, vol. 57, no. 1, pp. 355–369, Jan. 2009.
- [11] K. Srivastava and A. Nedic, "Distributed asynchronous constrained stochastic optimization," *IEEE J. Sel. Topics Signal Process.*, vol. 5, no. 4, pp. 772–790, Aug. 2011.
- [12] A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. Towfic, "Diffusion strategies for adaptation and learning over networks," *IEEE Sig. Process. Mag.*, vol. 30, no. 3, pp. 155–171, May 2013.
- [13] A. H. Sayed, "Diffusion adaptation over networks," in *Academic Press Library in Signal Processing*, R. Chellapa and S. Theodoridis, Eds., pp. 323–454. Academic Press, Elsevier, 2014.
- [14] 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.
- [15] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1035–1048, Mar. 2010.
- [16] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Trans. Signal Process.*, vol. 60, no. 8, pp. 4289–4305, Aug. 2012.
- [17] J. Chen and A. H. Sayed, "Distributed Pareto optimization via diffusion strategies," *IEEE J. Sel. Topics Signal Process.*, vol. 7, no. 2, pp. 205–220, Apr. 2013.
- [18] 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.
- [19] 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.
- [20] R. Abdolee, B. Champagne, and A. H. Sayed, "Estimation of space-time varying parameters using a diffusion LMS algorithm," *IEEE Trans. Signal Process.*, vol. 62, no. 2, pp. 403–418, Jan. 2014.
- [21] Y. Hirata, D. P. Mandic, H. Suzuki, and K. Aihara, "Wind direction modelling using multiple observation points," *Philosophical Transactions of the Royal Society A*, vol. 366, pp. 591–607, 2008.
- [22] 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. (ICASSP)*, Brisbane, Australia, 2015, pp. 3516–3520.

- [23] X. Zhao and A. H. Sayed, "Distributed clustering and learning over networks," *IEEE Trans. Signal Process.*, vol. 63, no. 13, pp. 3285–3300, Jul. 2015.
- [24] J. Chen, C. Richard, and A. H. Sayed, "Adaptive clustering for multitask diffusion networks," in *Proc. European Signal Process. Conf. (EUSIPCO)*, Nice, France, Sept. 2015, pp. 200–204.
- [25] R. Nassif, C. Richard, J. Chen, A. Ferrari, and A. H. Sayed, "Diffusion LMS over multitask networks with noisy links," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Shanghai, China, 2016.
- [26] S. Monajemi, K. Eftaxias, S. Sanei, and Ong S.-H., "An informed multitask diffusion adaptation approach to study tremor in Parkinson's disease," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 7, pp. 1306–1314, Oct. 2016.
- [27] Y. Wang, W. P. Tay, and W. Hu, "Multitask diffusion LMS with optimized inter-cluster cooperation," in *Proc. IEEE Stat. Signal Process. Workshop (SSP)*, Palma de Mallorca, Spain, 2016, pp. 1–5.
- [28] N. Roula, 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.
- [29] S. Monajemi, S. Sanei, Ong S.-H., and A. H. Sayed, "Adaptive regularized diffusion adaptation over multitask networks," in *Proc. IEEE Int. Workshop on Machine Learn. for Signal Process. (MLSP)*, Boston, USA, Sept. 2015, pp. 1–5.
- [30] S. Khawatmi, A. M. Zoubir, and A. H. Sayed, "Decentralized clustering over adaptive networks," in *Proc. European Signal Process. Conf. (EUSIPCO)*, Nice, France, 2015, pp. 2696–2700.
- [31] J. Plata-Chaves, N. Bogdanović, and K. Berberidis, "Distributed diffusion-based LMS for node-specific adaptive parameter estimation," *IEEE Trans. Signal Process.*, vol. 63, no. 13, pp. 3448–3460, Jul. 2015.
- [32] J. Plata-Chaves, H. H. Bahari, M. Moonen, and A. Bertrand, "Unsupervised diffusion-based LMS for node-specific parameter estimation over wireless sensor networks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Shanghai, China, Mar. 2016, pp. 4159–4163.
- [33] A. Bertrand and M. Moonen, "Distributed adaptive node-specific signal estimation in fully connected sensor networks – Part I: sequential node updating," *IEEE Trans. Signal Process.*, vol. 58, no. 10, pp. 5277–5291, Oct. 2010.
- [34] A. Bertrand and M. Moonen, "Distributed adaptive estimation of node-specific signals in wireless sensor networks with a tree topology," *IEEE Trans. Signal Process.*, vol. 59, no. 5, pp. 2196–2210, May 2011.
- [35] N. Bogdanović, J. Plata-Chaves, and K. Berberidis, "Distributed incremental-based LMS for node-specific adaptive parameter estimation," *IEEE Trans. Signal Process.*, vol. 62, no. 20, pp. 5382–5397, Oct. 2014.
- [36] J. Chen, S. K. Ting, C. Richard, and A. H. Sayed, "Group diffusion LMS," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Shanghai, China, March 2016.
- [37] J. Baxter, "A model of inductive bias learning," *J. Artif. Intell. Res.*, vol. 12, no. 1, pp. 149–198, Feb. 2000.
- [38] R. K. Ando and T. Zhang, "A framework for learning predictive structures from multiple tasks and unlabeled data," *J. Mach. Learn. Res.*, vol. 6, pp. 1817–1853, 2005.
- [39] J. Zhou, J. Chen, and J. Ye, "Clustered multi-task learning via alternating structure optimization," in *Proc. Adv. Neural Inf. Process. Syst.*, Granada, Spain, Dec. 2011, pp. 702–710.
- [40] J. Chen, L. Tang, J. Liu, and J. Ye, "A convex formulation for learning shared structures from multiple tasks," in *Proc. Ann. Int. Conf. Machine Learning (ICML)*, Montreal, Canada, Jun. 2009, pp. 137–144.
- [41] L. Jacob, F. Bach, and J.-P. Vert, "Clustered multitask learning: A convex formulation," *Adv. Neural Inf. Process. Syst.*, vol. 21, pp. 745–752, 2009.
- [42] G. Obozinski, B. Taskar, and M. Jordan, "Joint covariate selection and joint subspace selection for multiple classification problems," *Stat. Comput.*, vol. 20, no. 2, pp. 231–252, 2010.
- [43] L. J. Griffiths and C. W. Jim, "An alternative approach to linear constrained adaptive beamforming," *IEEE Trans. Antenn. Propag.*, vol. 30, no. 1, pp. 27–34, Jan. 1982.
- [44] X. Zhao and A. H. Sayed, "Clustering via diffusion adaptation over networks," in *Proc. Int. Workshop Cognitive Inf. Process. (CIP)*, Parador de Baiona, Spain, May 2012, pp. 1–6.

- [45] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Performance analysis of multitask diffusion adaptation over asynchronous networks," in *Proc. Asilomar Conf. Signals, Systems and Computers (ASILOMAR)*, Pacific Grove, CA, USA, Nov. 2014, pp. 788–792.
- [46] J. Chen and C. Richard, "Performance analysis of diffusion LMS in multitask networks," in *Proc. IEEE Int. Workshop on Compt. Adv. in Multi-Sensor Adaptive Process. (CAMSAP)*, Saint Martin, France, Dec. 2013, pp. 137–140.
- [47] J. Plata-Chaves, N. Bogdanović, and K. Berberidis, "Distributed incremental-based RLS for node-specific parameter estimation over adaptive networks," in *Proc. European Signal Process. Conf. (EUSIPCO)*, Marrakech, Morocco, Sept. 2013, pp. 1–5.
- [48] N. Bogdanović, J. Plata-Chaves, and K. Berberidis, "Distributed diffusion-based LMS for node-specific parameter estimation over adaptive networks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Florence, Italy, May 2014, pp. 7223–7227.
- [49] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [50] I. Gohberg and V. Olshevsky, "Fast algorithms with preprocessing for matrix-vector multiplication problems," *J. Complexity*, vol. 10, no. 4, pp. 411–427, 1994.
- [51] A. H. Sayed, *Fundamentals of Adaptive Filtering*, J. Wiley & Sons, Hoboken, NJ, 2003.
- [52] V. H. Nascimento and A. H. Sayed, "Unbiased and stable leakage-based adaptive filters," *IEEE Trans. Signal Process.*, vol. 47, no. 12, pp. 3261–3276, 1999.
- [53] F. Zhang, *The Schur Complement and its Applications*, Springer, 2005.