

# Penalty-Based Multitask Estimation with Non-Local Linear Equality Constraints

Fei Hua<sup>\*†‡</sup>, Roula Nassif<sup>‡</sup>, Cédric Richard<sup>‡</sup>, Haiyan Wang<sup>\*†</sup>

<sup>†</sup>School of Marine Science and Technology, Northwestern Polytechnical University (NPU), Xi'an 710072, China

<sup>\*</sup>Key Laboratory of Ocean Acoustics and Sensing (NPU), Ministry of Industry and Technology, Xi'an 710072, China

<sup>‡</sup>Laboratoire Lagrange, Université Côte d'Azur, OCA, CNRS, Nice 06108, France

Email: {fei.hua, roula.nassif}@oca.eu cedric.richard@unice.fr hywang@nwpu.edu.cn

**Abstract**—We consider distributed estimation problems over multitask networks where the parameter vectors at distinct agents are coupled via a set of linear equality constraints. Unlike previous existing works, the current work assumes that each constraint involves agents that are not necessarily one-hop neighbors. At each time instant, we assume that each agent has access to the instantaneous estimates of its one-hop neighbors and to the past estimates of its multi-hop neighbors through a multi-hop relay protocol. A distributed penalty-based algorithm is then derived and its performance analyses in the mean and in the mean-square-error sense are provided. Simulation results show the effectiveness of the strategy and validate the theoretical models.

## I. INTRODUCTION

Distributed estimation is used in a wide range of applications including communication [1], spectrum sensing [2], distributed localization [3], and power system monitoring [4]. Several useful distributed solutions, such as incremental strategies [5], diffusion strategies [6]–[11], and consensus strategies [12], [13] have been proposed in the literature to address single-task problems where all agents in the network collaborate to estimate a common parameter vector from noisy measurements. Among them, diffusion strategies are advantageous in terms of stability range, robustness, and performance [8]–[10].

In many applications, however, it happens that the agents in the network have to infer multiple parameter vectors simultaneously. Networks of this type are referred to as multitask networks [14], [15]. Multitask diffusion strategies were derived by exploiting prior information on the relationships among tasks. For example, appropriate regularization terms can be used to promote similarities between the tasks [14], [16], [17]. In [18], a diffusion-based algorithm is proposed to solve node-specific estimation problems where each node consists of a set of local parameters and a set of network global parameters. In [19], [20], the parameter space is decomposed into two orthogonal subspaces. The relations among tasks are modeled by assuming that they all share one of the subspaces. In some applications, such as the network flow problem [21], the basis pursuit problem [22], and the interference management problem [23], the parameter vectors may be related via a set of linear equality constraints. Distributed projection-based [24] and penalty-based [25] estimation algorithms were proposed to solve multitask estimation problems where each agent is interested in estimating its own parameter vector and where the parameter vectors at neighboring agents are related according to a set of linear equality constraints. In the current work, we consider

The work of F. Hua was partly supported by China Scholarship Council and NSFC grant 61471298. The work of R. Nassif and C. Richard was supported in part by ANR and DGA grant ANR-13-ASTR-0300 (ODISSEE project). The work of H. Wang was partly supported by NSFC under grants 61571365 and 61671386.

a more general multitask scenario where each equality constraint involves agents that are not necessarily one-hop neighbors.

Let  $N$  denote the number of agents in the network and  $P$  the total number of constraints. We are interested in devising a distributed adaptive solution to solve the following optimization problem:

$$\underset{\mathbf{w}_1, \dots, \mathbf{w}_N}{\text{minimize}} \quad J^{\text{glob}}(\mathbf{w}_1, \dots, \mathbf{w}_N) \triangleq \sum_{k=1}^N J_k(\mathbf{w}_k), \quad (1a)$$

$$\text{subject to} \quad \sum_{\ell \in \mathcal{I}_p} \mathbf{D}_{p\ell} \mathbf{w}_\ell + \mathbf{b}_p = \mathbf{0}, \quad p = 1, \dots, P \quad (1b)$$

Each agent  $k$  seeks to estimate its parameter vector  $\mathbf{w}_k \in \mathbb{R}^{M_k \times 1}$ , and has knowledge of its local cost  $J_k(\cdot)$  and the set of linear equality constraints that it is involved in. Each constraint is indexed by  $p$ , and defined by the  $L_p \times M_\ell$  matrices  $\mathbf{D}_{p\ell}$ , the  $L_p \times 1$  vector  $\mathbf{b}_p$ , and the set  $\mathcal{I}_p$  of agent indices involved in the  $p$ -th constraint. The previous works [24], [25] assumed that  $\mathcal{I}_p \subseteq \mathcal{N}_k$  for all  $k \in \mathcal{I}_p$  with  $\mathcal{N}_k$  denoting the one-hop neighborhood of agent  $k$  that consists of all agents that are connected to  $k$  by an edge. In this paper, we relax this assumption by considering scenarios where the constraints involve agents that are not necessarily one-hop neighbors. In order to derive a distributed solution relying solely on local interactions between neighbors, we shall employ multi-hop relay protocols to enable non-neighboring agents to share their estimates in order to satisfy their constraints. A penalty-based distributed estimation algorithm is derived and its stochastic behavior in the mean and in the mean-square-error sense is analyzed. Simulations are conducted to illustrate the effectiveness of the proposed algorithm and to validate the theoretical models.

**Notations:** All vectors are column vectors. The all-one vector of length  $N$  is denoted by  $\mathbf{1}_N$  and the identity matrix of size  $N$  is denoted by  $\mathbf{I}_N$ . The  $(k, \ell)$ -th block of a block matrix is denoted by  $[\cdot]_{k,\ell}$ . The operator  $\text{col}\{\cdot\}$  stacks its vector entries on top of each other. The symbol  $\otimes$  refers to the Kronecker product. The symbol  $\text{vec}(\cdot)$  denotes the vectorization operator that stacks the columns of a matrix on top of each other. The symbol  $\mathcal{N}_k^{(h)}$  refers to the  $h$ -hop neighborhood of agent  $k$ , that is,  $\ell \in \mathcal{N}_k^{(h)}$  means that the smallest number of hops from agent  $k$  to agent  $\ell$  is equal to  $h$ .

## II. PROBLEM FORMULATION AND PENALTY-BASED SOLUTION

Consider a strongly connected network of  $N$  agents. At each time instant  $i$ , each agent  $k$  has access to a zero-mean scalar observation  $d_k(i)$ , and to a zero-mean regression vector  $\mathbf{x}_k(i) \in \mathbb{R}^{M_k \times 1}$  with positive definite covariance matrix  $\mathbf{R}_{\mathbf{x},k} = \mathbb{E}\{\mathbf{x}_k(i)\mathbf{x}_k^\top(i)\}$ . The observations  $\{d_k(i), \mathbf{x}_k(i)\}$  are assumed to satisfy a linear regression model:

$$d_k(i) = \mathbf{x}_k^\top(i) \mathbf{w}_k^o + z_k(i), \quad (2)$$

where  $\mathbf{w}_k^o$  is an  $M_k \times 1$  unknown parameter vector to be estimated by agent  $k$ , and  $z_k(i)$  is a zero-mean noise with variance  $\sigma_{z,k}^2$  assumed to be spatially and temporally independent. In order to estimate  $\mathbf{w}_k^o$ , we associate with agent  $k$  the mean-square-error cost which is strongly convex and second-order differentiable :

$$J_k(\mathbf{w}_k) = \mathbb{E}|d_k(i) - \mathbf{x}_k^\top(i)\mathbf{w}_k|^2. \quad (3)$$

Let us collect the parameter vectors  $\mathbf{w}_k$  and  $\mathbf{w}_k^o$  from across the network into the following vectors of length  $M = \sum_{k=1}^N M_k$ :

$$\mathbf{w} \triangleq \text{col}\{\mathbf{w}_1, \dots, \mathbf{w}_N\}, \quad \mathbf{w}^o \triangleq \text{col}\{\mathbf{w}_1^o, \dots, \mathbf{w}_N^o\}.$$

Problem (1) can be written equivalently as:

$$\underset{\mathbf{w}}{\text{minimize}} \quad \sum_{k=1}^N \mathbb{E}|d_k(i) - \mathbf{x}_k^\top(i)\mathbf{w}_k|^2, \quad (4a)$$

$$\text{subject to} \quad \mathbf{D}\mathbf{w} + \mathbf{b} = \mathbf{0}. \quad (4b)$$

where  $\mathbf{D}$  is a  $P \times N$  block matrix with block entries  $\mathbf{D}_{p\ell}$  and  $\mathbf{b}$  is a  $P \times 1$  block column vector with block entries  $\mathbf{b}_p$ . We shall assume that  $P < N$  and that  $\mathbf{D}$  is full row rank, so that (4b) has at least one solution.

The positive-definite quadratic problem (4) has a unique global minimum given by:

$$\mathbf{w}^* = \mathbf{w}^o - \mathcal{R}_x^{-1} \mathbf{D}^\top (\mathcal{D} \mathcal{R}_x^{-1} \mathcal{D}^\top)^{-1} (\mathcal{D} \mathbf{w}^o + \mathbf{b}), \quad (5)$$

where  $\mathcal{R}_x \triangleq \text{diag}\{\mathbf{R}_{x,1}, \dots, \mathbf{R}_{x,N}\}$ .

Instead of using (5), we are interested in an adaptive distributed solution that is able to learn from streaming data and that relies on local interactions between neighboring agents. Penalty methods offer a simple way for tackling constrained optimization problems. These methods consist of approximating the constrained problem (4) into an unconstrained one by adding to the objective function a penalty term that penalizes any violation of the constraints:

$$\underset{\mathbf{w}}{\text{minimize}} \quad J_\eta^{\text{glob}}(\mathbf{w}) \triangleq \sum_{k=1}^N J_k(\mathbf{w}_k) + \eta \|\mathcal{D}\mathbf{w} + \mathbf{b}\|^2, \quad (6)$$

with  $\eta > 0$  a scalar parameter that controls the relative importance of adhering to the constraints. Increasing  $\eta$  improves the approximation (6) in quality [26]–[29], i.e.,  $\mathbf{w}^o(\eta)$  gets closer to  $\mathbf{w}^*$ . The above problem is strongly convex for any  $\eta$  and its closed form solution parameterized by  $\eta$  is given by:

$$\mathbf{w}^o(\eta) = (\mathcal{R}_x + \eta \mathcal{D}^\top \mathcal{D})^{-1} (\mathcal{R}_x \mathbf{w}^o - \eta \mathcal{D}^\top \mathbf{b}). \quad (7)$$

Applying a steepest-descent iteration to minimize the cost in (6) with respect to  $\mathbf{w}_k$  and starting from an initial condition  $\mathbf{w}_k(0)$ , we obtain the following algorithm at node  $k$ :

$$\mathbf{w}_k(i+1) = \mathbf{w}_k(i) - \mu \left[ \mathbf{R}_{x,k} \mathbf{w}_k(i) - \mathbf{r}_{dx,k} + \eta \sum_{p \in \mathcal{J}_k} \mathbf{D}_{pk}^\top \left( \sum_{\ell \in \mathcal{I}_p} \mathbf{D}_{p\ell} \mathbf{w}_\ell(i) + \mathbf{b}_p \right) \right], \quad (8)$$

where  $\mathbf{r}_{dx,k} \triangleq \mathbb{E}\{d_k(i)\mathbf{x}_k(i)\}$ ,  $\mathcal{J}_k$  denotes the set of constraint indices involving agent  $k$ , i.e.,  $\mathcal{J}_k \triangleq \{p|k \in \mathcal{I}_p\}$ . In order to evaluate  $\sum_{\ell \in \mathcal{I}_p} \mathbf{D}_{p\ell} \mathbf{w}_\ell(i) + \mathbf{b}_p$  in (8), agent  $k$  needs the estimates  $\mathbf{w}_\ell(i)$  from all agents  $\ell \in \mathcal{I}_p$ . These agents are not necessarily in the one-hop neighborhood of  $k$ , and need to employ multi-hop relay protocols to share their own estimate. We shall assume that the route from agent  $\ell \in \mathcal{I}_p$  to agent  $k$  with the smallest number of relays or hops, which is often the most energy-efficient route [30], is known. Instead of using  $\mathbf{w}_\ell(i)$  since it may not be available, agent  $k$  will use past estimate  $\mathbf{w}_\ell(i-j)$  of agent  $\ell$  where the delay  $j$  depends on the smallest number of hops from agent  $\ell$  to agent  $k$ , denoted by  $h_{\ell k}$ . In

this work, we shall assume that  $j = h_{\ell k} - 1$ . With this multi-hop relay protocol, at each time instant  $i$ , agent  $k$  has access to  $\mathbf{w}_\ell(i+1-h_{\ell k})$ .

Usually, the second-order moments  $\mathbf{R}_{x,k}$  and  $\mathbf{r}_{dx,k}$  in (8) are not available beforehand. We replace them by the instantaneous approximations [31]:

$$\mathbf{R}_{x,k} \approx \mathbf{x}_k(i)\mathbf{x}_k^\top(i), \quad \mathbf{r}_{dx,k} \approx d_k(i)\mathbf{x}_k(i), \quad (9)$$

Replacing  $\mathbf{w}_\ell(i)$  in (8) by  $\mathbf{w}_\ell(i+1-h_{\ell k})$ , and splitting the update iteration into two incremental steps by introducing the intermediate estimate  $\phi_k(i+1)$ , we obtain the following adaptive algorithm at agent  $k$ :

$$\phi_k(i+1) = \mathbf{w}_k(i) + \mu \mathbf{x}_k(i) \left( d_k(i) - \mathbf{x}_k^\top(i)\mathbf{w}_k(i) \right), \quad (10a)$$

$$\mathbf{w}_k(i+1) = \phi_k(i+1)$$

$$- \mu \eta \sum_{p \in \mathcal{J}_k} \mathbf{D}_{pk}^\top \left( \sum_{\ell \in \mathcal{I}_p} \mathbf{D}_{p\ell} \phi_\ell(i+2-h_{\ell k}) + \mathbf{b}_p \right), \quad (10b)$$

where in the second step (10b), we replaced  $\mathbf{w}_\ell(i+1-h_{\ell k})$  by the intermediate estimate  $\phi_\ell(i+2-h_{\ell k})$  which is a better estimate for the solution at agent  $\ell$ . We set  $\phi_\ell(i+2-h_{\ell k}) = \mathbf{0}$  if  $i+2-h_{\ell k} < 0$ , and  $h_{kk} = 1$ . In the first step (10a), which is the adaptation step, node  $k$  uses its own data to update its estimate  $\mathbf{w}_k(i)$  to an intermediate estimate  $\phi_k(i+1)$ . In the second step (10b), which corresponds to the penalization step, node  $k$  collects the intermediate estimates  $\phi_\ell(i+2-h_{\ell k})$  from nodes  $\ell \in \mathcal{I}_p$  for all  $p \in \mathcal{J}_k$ .

### III. STOCHASTIC BEHAVIOR ANALYSIS

In this section, we study the mean and the mean-square-error behavior of algorithm (10) with respect to  $\mathbf{w}^o(\eta)$  (7) and  $\mathbf{w}^*$  (5). Let  $\tilde{\mathbf{w}}_k(i) \triangleq \mathbf{w}_k^o(\eta) - \mathbf{w}_k(i)$  and  $\tilde{\phi}_k(i) \triangleq \mathbf{w}_k^o(\eta) - \phi_k(i)$  denote the error vectors at agent  $k$ . Taking into account the delayed information emerging from the multi-hop protocol, we introduce the extended network error vectors:

$$\tilde{\mathbf{w}}_e(i) \triangleq \text{col}\{\tilde{\mathbf{w}}_1(i), \dots, \tilde{\mathbf{w}}_N(i), \tilde{\phi}_1(i), \dots, \tilde{\phi}_N(i), \dots, \tilde{\phi}_1(i-H+2), \dots, \tilde{\phi}_N(i-H+2)\}, \quad (11)$$

$$\tilde{\phi}_e(i) \triangleq \text{col}\{\tilde{\phi}_1(i), \dots, \tilde{\phi}_N(i), \tilde{\phi}_1(i-1), \dots, \tilde{\phi}_N(i-1), \dots, \tilde{\phi}_1(i-H+1), \dots, \tilde{\phi}_N(i-H+1)\}, \quad (12)$$

where  $H = \max_{k,\ell} \{h_{\ell k}\}$ . We define  $\mathbf{w}^\delta \triangleq \mathbf{w}^o(\eta) - \mathbf{w}^o$ ,  $\mathbf{w}_e^\delta \triangleq \mathbb{1}_H \otimes \mathbf{w}^\delta$  and  $\mathbf{w}^{\delta'} \triangleq \mathbf{w}^* - \mathbf{w}^o$ ,  $\mathbf{w}_e^{\delta'} \triangleq \mathbb{1}_H \otimes \mathbf{w}^{\delta'}$ . Before proceeding, we introduce the following assumption on the regressors.

**Assumption 1.** *The regressors  $\mathbf{x}_k(i)$  arise from a zero-mean random process that is temporally white and spatially independent.*

This assumption is commonly used in the adaptive filtering literature [31]. It helps to simplify the derivations without constraining the conclusions.

#### A. Mean error behavior analysis

Using recursion (10) and data model (2), the  $H \times 1$  block vector  $\tilde{\mathbf{w}}_e(i)$  can be written as:

$$\tilde{\mathbf{w}}_e(i+1) = \mathcal{B}(i)\tilde{\mathbf{w}}_e(i) - \mu \mathbf{g}(i) + \mu \mathbf{r}(i) + \mu \eta \mathbf{f}, \quad (13)$$

where

$$\mathcal{B}(i) \triangleq \mathcal{H} [\mathbf{I}_{MH} - \mu \mathcal{R}_{x,e}(i)], \quad (14)$$

$$\mathbf{g}(i) \triangleq \mathcal{H} \mathbf{p}_{z,x,e}(i), \quad (15)$$

$$\mathbf{r}(i) \triangleq \mathcal{H} \mathcal{R}_{x,e}(i) \mathbf{w}_e^\delta, \quad (16)$$

$$\mathbf{f} \triangleq \text{col}\left\{ \mathcal{D}^\top (\mathcal{D} \mathbf{w}^o(\eta) + \mathbf{b}), \mathbf{0}_{M \times 1}, \dots, \mathbf{0}_{M \times 1} \right\}, \quad (17)$$

with  $\mathcal{H}$  an  $H \times H$  block matrix,  $\mathcal{R}_{x,e}(i)$  an  $H \times H$  block diagonal matrix, and  $\mathbf{p}_{z_{x,e}}(i)$  an  $H \times 1$  block vector:

$$\mathcal{H} \triangleq [\mathcal{I} - \mu\eta\mathcal{C}_{D,e}], \quad (18)$$

$$\mathcal{R}_{x,e}(i) \triangleq \text{diag}\{\mathcal{R}_x(i), \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\}, \quad (19)$$

$$\mathbf{p}_{z_{x,e}}(i) \triangleq \text{col}\{\mathbf{p}_{z_x}(i), \mathbf{0}_{M \times 1}, \dots, \mathbf{0}_{M \times 1}\}, \quad (20)$$

where

$$\mathcal{I} = \begin{bmatrix} \mathbf{I}_M & \mathbf{0}_{M \times M(H-1)} \\ \mathbf{I}_{M(H-1)} & \mathbf{0}_{M(H-1) \times M} \end{bmatrix}, \quad (21)$$

$$\mathcal{C}_{D,e} = \begin{bmatrix} \mathcal{C}_{D,1} & \mathcal{C}_{D,2} & \dots & \mathcal{C}_{D,H} \\ \hline & \mathbf{0}_{M(H-1) \times MH} & & \end{bmatrix}, \quad (22)$$

$$\mathcal{R}_x(i) = \text{diag}\left\{\mathbf{x}_k(i)\mathbf{x}_k^\top(i)\right\}_{k=1}^N, \quad (23)$$

$$\mathbf{p}_{z_x}(i) = \text{col}\left\{\mathbf{x}_k(i)z_k(i)\right\}_{k=1}^N. \quad (24)$$

and  $\mathcal{C}_{D,h}$  is an  $N \times N$  block matrix with  $(k,l)$ -th block given by:

$$[\mathcal{C}_{D,h}]_{k,l} = \begin{cases} \sum_{p \in \mathcal{J}_k} \mathbf{D}_{pk}^\top \mathbf{D}_{pl} & \text{if } \ell \in \mathcal{I}_p \cap \mathcal{N}_k^{(h)}, \\ \mathbf{0} & \text{otherwise.} \end{cases} \quad (25)$$

Note that  $\sum_{h=1}^H \mathcal{C}_{D,h} = \mathcal{D}^\top \mathcal{D}$ .

Taking the expectation of both sides of (13) and using the fact that  $\mathbb{E}\{\mathbf{g}(i)\} = \mathbf{0}$ , we obtain:

$$\mathbb{E}\tilde{\mathbf{w}}_e(i+1) = \mathcal{B}\mathbb{E}\tilde{\mathbf{w}}_e(i) + \mu\mathbf{r} + \mu\eta\mathbf{f}, \quad (26)$$

where

$$\mathcal{B} \triangleq \mathcal{H}[\mathbf{I}_{MH} - \mu\mathcal{R}_{x,e}], \quad (27)$$

$$\mathbf{r} \triangleq \mathcal{H}\mathcal{R}_{x,e}\mathbf{w}_e^\delta, \quad (28)$$

$$\mathcal{R}_{x,e} \triangleq \text{diag}\{\mathcal{R}_x, \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\}. \quad (29)$$

The mean error vector  $\mathbb{E}\{\tilde{\mathbf{w}}_e(i)\}$  converges as  $i \rightarrow \infty$  if the matrix  $\mathcal{B}$  is stable, i.e., the spectral radius of  $\mathcal{B}$  is less than 1. In this case, the asymptotic mean bias is given by:

$$\mathbb{E}\tilde{\mathbf{w}}_e(\infty) = \lim_{i \rightarrow \infty} \mathbb{E}\tilde{\mathbf{w}}_e(i) = \mu(\mathbf{I}_M - \mathcal{B})^{-1}(\mathbf{r} + \eta\mathbf{f}). \quad (30)$$

Using similar definition with (11), we find that the extended mean error recursion  $\mathbb{E}\tilde{\mathbf{w}}'_e(i)$  with respect to  $\mathbf{w}^*$  evolves according to:

$$\mathbb{E}\tilde{\mathbf{w}}'_e(i+1) = \mathcal{B}\mathbb{E}\tilde{\mathbf{w}}'_e(i) + \mu\mathbf{r}_s. \quad (31)$$

where  $\mathbf{r}_s \triangleq \mathcal{H}\mathcal{R}_{x,e}\mathbf{w}_e^{\delta'}$ . When  $\mathcal{B}$  is stable, we obtain:

$$\mathbb{E}\tilde{\mathbf{w}}'_e(\infty) = \lim_{i \rightarrow \infty} \mathbb{E}\tilde{\mathbf{w}}'_e(i) = \mu(\mathbf{I}_{MH} - \mathcal{B})^{-1}\mathbf{r}_s. \quad (32)$$

We observe that, when  $\mathbf{w}^o$  satisfies the constraints, we have  $\mathbf{w}^o = \mathbf{w}^o(\eta) = \mathbf{w}^*$ ,  $\mathbf{r} = \mathbf{0}$ ,  $\mathbf{f} = \mathbf{0}$ , and  $\mathbf{r}_s = \mathbf{0}$ . In this case, the asymptotic mean biases  $\mathbb{E}\tilde{\mathbf{w}}_e(\infty)$  and  $\mathbb{E}\tilde{\mathbf{w}}'_e(\infty)$  reduce to zero.

### B. Mean-square-error behavior analysis

We shall evaluate the weighted variance  $\mathbb{E}\{\|\tilde{\mathbf{w}}_e(i+1)\|_\Sigma^2\}$  where  $\Sigma$  is a positive semi-definite matrix that we are free to choose. Let  $\boldsymbol{\sigma} \triangleq \text{vec}(\Sigma)$ . In the following, we use the alternative notation  $\|\mathbf{w}\|_\Sigma^2$  to refer to the same weighted squared norm  $\|\mathbf{w}\|_\Sigma^2$ . Following the same line of reasoning as in [7], [10] for single-task diffusion strategies, and extending the arguments to our multitask scenario, we find:

$$\mathbb{E}\{\|\tilde{\mathbf{w}}_e(i+1)\|_\Sigma^2\} = \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_{\mathcal{F}\boldsymbol{\sigma}}^2\} + [\text{vec}(\mathcal{Y}(i))]^\top \boldsymbol{\sigma}, \quad (33)$$

where  $\mathcal{F}$  is the  $(MH)^2 \times (MH)^2$  matrix given by:

$$\mathcal{F} \triangleq \mathbb{E}\{\mathcal{B}^\top(i) \otimes \mathcal{B}^\top(i)\}. \quad (34)$$

It is sufficient in this work to consider the case of sufficiently small step-sizes where the influence of terms involving higher-order powers of  $\mu$  can be ignored [10], [31]. In this case,  $\mathcal{F}$  can be approximated by  $\mathcal{F} \approx \mathcal{B}^\top \otimes \mathcal{B}^\top$ . Under this approximation, the stability of  $\mathcal{F}$  is ensured if the matrix  $\mathcal{B}$  is stable.

The matrix  $\mathcal{Y}(i)$  in (33) is given by:

$$\mathcal{Y}(i) \triangleq \mu^2 \mathcal{G}^\top + \mu^2 \mathcal{Q}^\top + 2\mu \mathcal{P}^\top(i) + \mu^2 \eta^2 \mathbf{f} \mathbf{f}^\top + 2\mu \eta \mathbf{f} \mathbb{E}\tilde{\mathbf{w}}_e^\top(i) \mathcal{B}^\top + 2\mu^2 \eta \mathbf{f} \mathbf{r}^\top, \quad (35)$$

where  $\mathcal{G}$  and  $\mathcal{Q}$  are  $H \times H$  block matrices given by:

$$\begin{aligned} \mathcal{G} &\triangleq \mathbb{E}\{\mathbf{g}(i)\mathbf{g}^\top(i)\} \\ &= \mathcal{H} \text{diag}\{\mathcal{T}, \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\} \mathcal{H}^\top, \end{aligned} \quad (36)$$

$$\begin{aligned} \mathcal{Q} &\triangleq \mathbb{E}\{\mathbf{r}(i)\mathbf{r}^\top(i)\} \\ &= \mathcal{H} \text{diag}\{\mathcal{T}_1, \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\} \mathcal{H}^\top, \end{aligned} \quad (37)$$

with  $\mathcal{T} \triangleq \text{diag}\{\sigma_{z,k}^2 \mathcal{R}_{x,k}\}_{k=1}^N$ ,  $\mathcal{T}_1 \triangleq \mathbb{E}\{\mathcal{R}_x(i) \mathcal{K}_1 \mathcal{R}_x(i)\}$ , and  $\mathcal{K}_1 \triangleq \mathbf{w}^\delta (\mathbf{w}^\delta)^\top$ . The evaluation of  $\mathcal{T}_1$  depends on higher order moments of the regressors. In the following, we shall evaluate  $\mathcal{T}_1$  when the regressors are zero-mean real Gaussian. For any square matrix  $\mathbf{A}$  and zero-mean Gaussian regressors, we have [32]:

$$\begin{aligned} \mathbb{E}\{\mathbf{x}_k(i)\mathbf{x}_k^\top(i) \mathbf{A} \mathbf{x}_\ell(i)\mathbf{x}_\ell^\top(i)\} &= \mathbf{R}_{x,k} \mathbf{A} \mathbf{R}_{x,\ell} + \\ &\delta_{k,\ell} (\mathbf{R}_{x,k} \mathbf{A}^\top \mathbf{R}_{x,k} + \mathbf{R}_{x,k} \text{tr}(\mathbf{R}_{x,k} \mathbf{A})). \end{aligned} \quad (38)$$

Using (38), it can be verified that  $\mathcal{T}_1$  can be expressed as:

$$\begin{aligned} \mathcal{T}_1 &= \mathcal{R}_x \mathcal{K}_1 \mathcal{R}_x + \sum_{k=1}^N \left( \mathcal{S}_k (\mathbf{I}_N \otimes \mathbf{R}_{x,k}) \mathcal{K}_1^\top (\mathbf{I}_N \otimes \mathbf{R}_{x,k}) \mathcal{S}_k + \right. \\ &\quad \left. \mathcal{S}_k (\mathbf{I}_N \otimes \mathbf{R}_{x,k}) \mathcal{Z}_k \mathcal{S}_k \right), \end{aligned} \quad (39)$$

where  $\mathcal{S}_k \triangleq \text{diag}(\mathbf{e}_k^\top) \otimes \mathbf{I}_{M_k}$  is an  $N \times N$  block diagonal matrix and  $\mathbf{e}_k$  is the column vector with a unit entry at position  $k$  and zeros elsewhere.  $\mathcal{Z}_k$  is an  $N \times N$  block matrix with the  $(m,n)$ -th block given by:

$$[\mathcal{Z}_k]_{m,n} = [\text{vec}(\mathbf{R}_{x,k})]^\top \text{vec}([\mathcal{K}_1]_{m,n}) \mathbf{I}_{M_k}. \quad (40)$$

The  $H \times H$  block matrix  $\mathcal{P}(i)$  is given by:

$$\begin{aligned} \mathcal{P}(i) &\triangleq \mathbb{E}\{\mathcal{B}(i)\tilde{\mathbf{w}}_e(i)\mathbf{r}^\top(i)\} \\ &= \mathcal{H} \text{diag}\{\mathcal{T}_2, \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\} \mathcal{H}^\top \end{aligned} \quad (41)$$

with

$$\mathcal{T}_2(i) \triangleq \mathbb{E}\{\mathcal{R}_x(i) \mathcal{K}_2(i) (\mathbf{I}_M - \mu \mathcal{R}_x(i))\}, \quad (42)$$

$$\mathcal{K}_2(i) \triangleq \mathbb{E}\{\tilde{\mathbf{w}}(i)\} (\mathbf{w}^\delta)^\top. \quad (43)$$

Note that, following similar arguments as the one for  $\mathcal{Q}$ ,  $\mathcal{P}(i)$  can be evaluated for zero-mean real Gaussian regressors.

Starting from the initial condition  $\tilde{\mathbf{w}}_e(0)$  in (33) and comparing the expressions of  $\mathbb{E}\{\|\tilde{\mathbf{w}}_e(i+1)\|_\Sigma^2\}$  and  $\mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_\Sigma^2\}$ , we obtain:

$$\begin{aligned} \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i+1)\|_\Sigma^2\} &= \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_\Sigma^2\} + \\ &[\text{vec}(\mathbb{E}\{\tilde{\mathbf{w}}_e(0)\tilde{\mathbf{w}}_e^\top(0)\})]^\top (\mathcal{F} - \mathbf{I}_{(MH)^2}) \boldsymbol{\sigma} + \\ &[\text{vec}(\mathcal{Y}(i))]^\top \boldsymbol{\sigma} + \boldsymbol{\Gamma}(i) \boldsymbol{\sigma}, \end{aligned} \quad (44)$$

where  $\boldsymbol{\Gamma}(i)$  is an  $(MH)^2 \times 1$  vector that evolves according to:

$$\boldsymbol{\Gamma}(i+1) = [\text{vec}(\mathcal{Y}(i))]^\top (\mathcal{F} - \mathbf{I}_{(MH)^2}) + \boldsymbol{\Gamma}(i), \quad (45)$$

with  $\boldsymbol{\Gamma}(0) = \mathbf{0}_{(MH)^2 \times 1}$ .

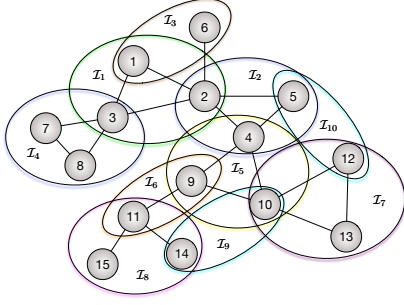


Fig. 1: Multitask MSE network with constraints.

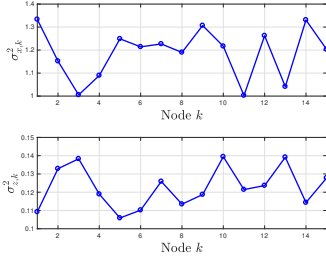


Fig. 2: Regressors and noise variances.

Recursion (33) converges to a steady-state value if the matrix  $\mathcal{F}$  is stable and the mean error vector  $\mathbb{E}\{\tilde{\mathbf{w}}_e(i)\}$  converges. In this case, we obtain from (33):

$$\lim_{i \rightarrow \infty} \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_{(\mathbf{I}_{(MH)^2} - \mathcal{F})}^2\} = [\text{vec}(\mathcal{Y}(\infty))]^\top \boldsymbol{\sigma}, \quad (46)$$

where  $\mathcal{Y}(\infty)$  can be obtained from (30) and (35). Define the network mean-square-deviation (MSD) as:  $\zeta^* = \lim_{i \rightarrow \infty} \frac{1}{N} \mathbb{E}\{\|\tilde{\mathbf{w}}(i)\|^2\}$ . In this case, the steady-state network MSD can be written as:

$$\zeta^* = \lim_{i \rightarrow \infty} \frac{1}{N} \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_{\sigma_{ss}}^2\} \quad (47)$$

with  $\boldsymbol{\sigma}_{ss} = \text{vec}(\text{diag}\{\mathbf{I}_M, \mathbf{0}_{M \times M}, \dots, \mathbf{0}_{M \times M}\})$ , then the network MSD can be obtained from (46) by selecting  $\boldsymbol{\sigma}$  that satisfies:

$$(\mathbf{I}_{(MH)^2} - \mathcal{F})\boldsymbol{\sigma} = \frac{1}{N} \boldsymbol{\sigma}_{ss}. \quad (48)$$

Replacing  $\boldsymbol{\sigma}$  in (46) by  $\frac{1}{N}(\mathbf{I}_{(MH)^2} - \mathcal{F})^{-1} \boldsymbol{\sigma}_{ss}$  we have:

$$\zeta^* = \frac{1}{N} [\text{vec}(\mathcal{Y}(\infty))]^\top (\mathbf{I}_{(MH)^2} - \mathcal{F})^{-1} \boldsymbol{\sigma}_{ss}. \quad (49)$$

Finally, the transient and steady-state behavior of  $\mathbb{E}\{\|\tilde{\mathbf{w}}'_e(i)\|_{\boldsymbol{\sigma}}^2\}$  can be derived from the following relation:

$$\mathbb{E}\{\|\tilde{\mathbf{w}}'_e(i)\|_{\boldsymbol{\sigma}}^2\} = \mathbb{E}\{\|\tilde{\mathbf{w}}_e(i)\|_{\boldsymbol{\sigma}}^2\} + 2\mathbb{E}\{\tilde{\mathbf{w}}_e(i)\}^\top \boldsymbol{\Sigma} \mathbf{w}_{d,e} + \|\mathbf{w}_{d,e}\|_{\boldsymbol{\Sigma}}^2$$

with  $\mathbf{w}_{d,e} = \mathbf{1}_H \otimes (\mathbf{w}^* - \mathbf{w}^o(\eta))$ .

#### IV. SIMULATIONS

In this section, we provide experimental results to illustrate the convergence of algorithm (10) and to validate our theoretical models. We considered a network of 15 nodes with the topology and the constraints shown in Fig. 1. We randomly sampled 10 linear constraints of the form  $\sum_{\ell \in \mathcal{I}_p} d_{p\ell} \mathbf{w}_\ell = b_p \cdot \mathbf{1}_2$ , where the coefficients  $d_{p\ell}$  and  $b_p$  were randomly chosen from  $\{-2, -1, 1, 2\}$ . The regression vectors  $\mathbf{x}_k(i)$  were zero-mean Gaussian with covariance matrix  $\mathbf{R}_{x,k} = \sigma_{x,k}^2 \mathbf{I}_2$ . The noises  $z_k(i)$  were zero-mean i.i.d. Gaussian

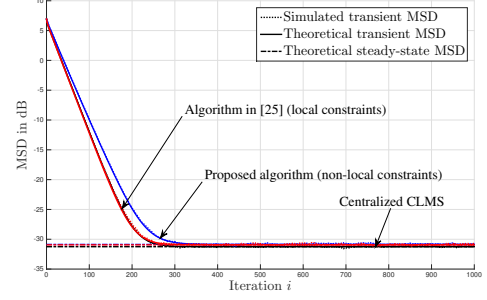


Fig. 3: MSD comparison (perfect model scenario).

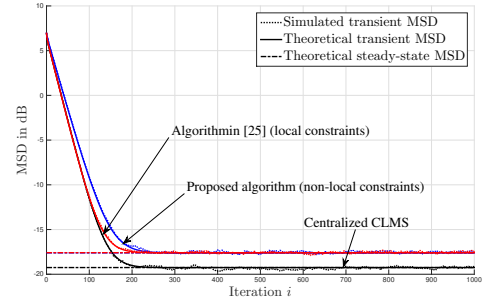


Fig. 4: MSD comparison (imperfect model scenario).

random variables independent of any other signal with variances  $\sigma_{z,k}^2$ . The variances  $\sigma_{x,k}^2$  and  $\sigma_{z,k}^2$  used in the simulations are shown in Fig. 2. The results were averaged over 200 Monte-Carlo runs.

We considered two scenarios: i) the parameter vector  $\mathbf{w}^o = \mathbf{w}_o$  where  $\mathbf{w}_o$  satisfies the constraints, i.e.  $\mathbf{w}^o = \mathbf{w}^*$  (Fig.3); ii)  $\mathbf{w}^o$  does not satisfy the constraints, specifically, we perturbed it as  $\mathbf{w}^o = \mathbf{w}_o + \mathbf{u}$  where  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  (Fig.4). The step-size  $\mu$  was set to 0.02. We compared our algorithm (10) with the centralized CLMS algorithm [24], [33]. Furthermore, for comparison purposes, we assumed additional links connecting nodes 1 to 6, 5 to 12, 10 to 14, and 14 to 15. In this case, the constraints are local and the algorithm derived in [25] can be applied. We set  $\mu = 0.018$  for the centralized CLMS algorithm and for algorithm [25] (with additional links) so that their steady-state MSD match. Observe that the simulation results match well the theoretical models. Furthermore, our algorithm performs well in the mean-square-error compared to the centralized solution. Finally, observe that, as expected, the delays emerging from the multi-hop protocols required in non-local constraints scenarios will lead to a slower convergence rate.

#### V. CONCLUSION

We proposed a distributed multitask algorithm for estimating multiple parameter vectors that are coupled through non-local linear equality constraints. Based on the penalty method, we solved the original constrained problem by approximating it into an unconstrained one. A multi-hop relay protocol was employed in order to deal with the non-local constraints and to devise a distributed algorithm. The stochastic behavior of the algorithm in the mean and in the mean-square-error sense was studied. Simulation results were conducted to show the effectiveness of the proposed method and to validate our theoretical performance analysis.

## REFERENCES

- [1] S. Barbarossa, S. Sardellitti, and P. Di Lorenzo, "Distributed detection and estimation in wireless sensor networks," in *Academic Press Library in Signal Processing*, S. Theodoridis and R. Chellappa, Eds. Academic Press, Elsevier, 2013, vol. 2, pp. 329–408.
- [2] Y. Zhang, W. P. Tay, K. H. Li, and D. Gaiti, "Distributed boundary estimation for spectrum sensing in cognitive radio networks," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 11, pp. 1961–1973, 2014.
- [3] F. Meyer, O. Hlinka, H. Wymeersch, E. Riegler, and F. Hlawatsch, "Distributed localization and tracking of mobile networks including noncooperative objects," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 1, pp. 57–71, 2016.
- [4] L. Xie, D.-H. Choi, S. Kar, and H. V. Poor, "Fully distributed state estimation for wide-area monitoring systems," *IEEE Transactions on Smart Grid*, vol. 3, no. 3, pp. 1154–1169, 2012.
- [5] D. P. Bertsekas, "A new class of incremental gradient methods for least squares problems," *SIAM Journal on Optimization*, vol. 7, no. 4, pp. 913–926, 1997.
- [6] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1035–1048, 2010.
- [7] A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. J. Towfic, "Diffusion strategies for adaptation and learning over networks: an examination of distributed strategies and network behavior," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 155–171, 2013.
- [8] S.-Y. Tu and A. H. Sayed, "Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 12, pp. 6217–6234, 2012.
- [9] Z. J. Towfic and A. H. Sayed, "Stability and performance limits of adaptive primal-dual networks," *IEEE Transactions on Signal Processing*, vol. 63, no. 11, pp. 2888–2903, 2015.
- [10] A. H. Sayed, "Diffusion adaptation over networks," in *Academic Press Library in Signal Processing*, S. Theodoridis and R. Chellappa, Eds. Academic Press, Elsevier, 2013, vol. 3, pp. 322–454.
- [11] —, "Adaptation, learning, and optimization over networks," *Foundations and Trends in Machine Learning*, vol. 7, no. 4-5, pp. 311–801, 2014.
- [12] L. Xiao, S. Boyd, and S. Lall, "A scheme for robust distributed sensor fusion based on average consensus," in *Proceedings of the 4th international symposium on Information processing in sensor networks*. IEEE, 2005, p. 9.
- [13] R. Olfati-Saber, J. A. Fax, and R. M. Murray, "Consensus and cooperation in networked multi-agent systems," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 215–233, 2007.
- [14] J. Chen, C. Richard, and A. H. Sayed, "Multitask diffusion adaptation over networks," *IEEE Transactions on Signal Processing*, vol. 62, no. 16, pp. 4129–4144, 2014.
- [15] —, "Diffusion LMS over multitask networks," *IEEE Transactions on Signal Processing*, vol. 63, no. 11, pp. 2733–2748, 2015.
- [16] Y. Wang, W. P. Tay, and W. Hu, "A multitask diffusion strategy with optimized inter-cluster cooperation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 3, pp. 504–517, 2017.
- [17] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Proximal multitask learning over networks with sparsity-inducing coregularization," *IEEE Transactions on Signal Processing*, vol. 64, no. 23, pp. 6329–6344, 2016.
- [18] J. Plata-Chaves, N. Bogdanović, and K. Berberidis, "Distributed diffusion-based LMS for node-specific adaptive parameter estimation," *IEEE Transactions on Signal Processing*, vol. 63, no. 13, pp. 3448–3460, 2015.
- [19] J. Chen, C. Richard, A. O. Hero, and A. H. Sayed, "Diffusion LMS for multitask problems with overlapping hypothesis subspaces," in *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, 2014.
- [20] J. Chen, C. Richard, and A. H. Sayed, "Multitask diffusion adaptation over networks with common latent representations," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 3, pp. 563–579, 2017.
- [21] D. P. Bertsekas, *Network optimization: continuous and discrete models*. Athena Scientific, 1998.
- [22] J. F. Mota, J. M. Xavier, P. M. Aguiar, and M. Puschel, "Distributed basis pursuit," *IEEE Transactions on Signal Processing*, vol. 60, no. 4, pp. 1942–1956, 2012.
- [23] C. Shen, T.-H. Chang, K.-Y. Wang, Z. Qiu, and C.-Y. Chi, "Distributed robust multicell coordinated beamforming with imperfect CSI: An ADMM approach," *IEEE Transactions on signal processing*, vol. 60, no. 6, pp. 2988–3003, 2012.
- [24] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Diffusion LMS for multitask problems with local linear equality constraints," *IEEE Transactions on Signal Processing*, to appear. Also available as arXiv:1610.02943.
- [25] F. Hua, R. Nassif, C. Richard, H. Wang, and J. Huang, "Penalty-based multitask distributed adaptation over networks with constraints," submitted to 2017 Asilomar Conference on Signals, Systems and Computers.
- [26] B. T. Polyak, *Introduction to Optimization*. Optimization Software, New York, 1987.
- [27] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear programming: theory and algorithms*. John Wiley & Sons, 2013.
- [28] D. G. Luenberger and Y. Ye, *Linear and nonlinear programming*. Springer, 2015, vol. 228.
- [29] Z. J. Towfic and A. H. Sayed, "Adaptive penalty-based distributed stochastic convex optimization," *IEEE Transactions on Signal Processing*, vol. 62, no. 15, pp. 3924–3938, 2014.
- [30] W. Dargie and C. Poellabauer, *Fundamentals of wireless sensor networks: theory and practice*. John Wiley & Sons, 2010.
- [31] A. H. Sayed, *Adaptive filters*. John Wiley & Sons, 2011.
- [32] K. B. Petersen and M. S. Pedersen, "The matrix cookbook. version: November 15, 2012," 2012.
- [33] O. L. Frost, "An algorithm for linearly constrained adaptive array processing," *Proceedings of the IEEE*, vol. 60, no. 8, pp. 926–935, 1972.