

# Doubly compressed diffusion LMS over adaptive networks

Ibrahim El Khalil Harrane, Rémi Flamary, Cédric Richard

### ▶ To cite this version:

Ibrahim El Khalil Harrane, Rémi Flamary, Cédric Richard. Doubly compressed diffusion LMS over adaptive networks. 2016 50th Asilomar Conference on Signals, Systems and Computers, Nov 2016, Pacific Grove, France. pp.987-991, 10.1109/ACSSC.2016.7869515 . hal-03633804

## HAL Id: hal-03633804 https://hal.science/hal-03633804

Submitted on 7 Apr 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Doubly compressed diffusion LMS over adaptive networks

Ibrahim El Khalil Harrane, Rémi Flamary, Cédric Richard

Université Côte d'Azur, OCA, CNRS, France

ibrahim.harrane@oca.eu, remi.flamary@unice.fr, cedric.richard@unice.fr

*Abstract*—Diffusion LMS is an efficient strategy for solving distributed optimization problems with cooperating agents. Nodes are interested in estimating the same parameter vector and exchange information with their neighbors to improve their local estimates. Successful implementation of such applications relies on a substantial amount of communication resources. In this paper, we introduce diffusion LMS strategies that offer significantly reduced communication load without compromising performance. We perform analyses in the mean and mean-square sense of these algorithms. Simulations results are provided to confirm the theoretical findings.

#### I. INTRODUCTION

Distributed adaptation over network has become an active research area with the increasing popularity of wireless sensor networks and, more recently, with the advent of the internet of things. An accessible overview of results in the field can be found in [1]–[5]. The interconnected agents continuously collect data, learn and adapt by carrying out data mining tasks on their own measurements as well as their neighbors. Albeit outperformed by centralized strategies, diffusion strategies are more robust and resilient to agents and links failures. Furthermore, their scalability and flexibility make them attractive for addressing inference problems in a collaborative manner. Among other possible strategies [6]-[8], diffusion LMS plays a central role with its enhanced efficiency and low complexity, and has been extensively studied in the literature in single task [1]-[3] and multitask frameworks [9]-[13]. It was also considered in other contexts such as nonlinear system identification [14] and dictionary learning [15].

Consider an interconnected network of N nodes. The aim of each node is to estimate an  $L \times 1$  unknown vector  $\boldsymbol{w}^o$  from collected measurements. Node k has access to local streaming measurements  $\{d_k(i), \boldsymbol{u}_{k,i}\}$  where  $d_k(i)$  is a scalar zero-mean reference signal, and  $\boldsymbol{u}_{k,i}$  is an  $L \times 1$  zero-mean regression vector with covariance matrix  $\boldsymbol{R}_{uk} = \mathbb{E}\{\boldsymbol{u}_{k,i}\boldsymbol{u}_{k,i}^{\top}\} > 0$ . The data at agent k and time i are assumed to be related via the linear regression model:

$$d_k(i) = \boldsymbol{u}_{k,i}^{\dagger} \boldsymbol{w}^o + v_k(i) \tag{1}$$

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

$$J_k(\boldsymbol{w}) = E\{|d_k(i) - \boldsymbol{u}_{k,i}^\top \boldsymbol{w}|^2\}$$
(2)

Diffusion LMS strategies seek the minimizer of the following aggregate cost function:

$$J^{\text{glob}}(\boldsymbol{w}) = \sum_{k=1}^{N} J_k(\boldsymbol{w})$$
(3)

in a cooperative manner. Let  $w_{k,i}$  denote the estimate of the minimizer of (3) at node k and time instant i. The general structure of diffusion LMS in its Adapt-then-Combine (ATC) form is given by:

$$\boldsymbol{\psi}_{k,i} = \boldsymbol{w}_{k,i-1} - \mu_k \sum_{\ell \in \mathcal{N}_k} c_{\ell k} \hat{\nabla}_w J_\ell(\boldsymbol{w}_{k,i-1})$$
(4)

$$\boldsymbol{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \boldsymbol{\psi}_{\ell,i} \tag{5}$$

where  $\hat{\nabla}_w J_\ell(\boldsymbol{w}_{k,i-1}) = -\boldsymbol{u}_{\ell,i}[d_\ell(i) - \boldsymbol{u}_{\ell,i}^\top \boldsymbol{w}_{k,i-1}]$ ,  $\mathcal{N}_k$  denotes the neighborhood of node k including k, and  $\mu_k$  is a positive step-size. Nonnegative coefficients  $\{a_{\ell k}\}$  and  $\{c_{\ell k}\}$  define a left-stochastic matrix  $\boldsymbol{A}$  and a right-stochastic matrix  $\boldsymbol{C}$ , respectively. Diffusion LMS provides a scalable estimation framework with high flexibility for ad-hoc deployment. Nevertheless, the requirement for all nodes to exchange their current estimates  $\boldsymbol{w}_{k,i-1}$ ,  $\hat{\nabla}_w J_\ell(\boldsymbol{w}_{k,i-1})$  and  $\boldsymbol{\psi}_{k,i}$  with their neighbors at each iteration imposes a substantial burden on communication and energy ressources. Therefore, reducing the communication cost while maintaining the benefits of cooperation is of major importance for systems with limited energy budget such as wireless sensor networks.

In recent years, several strategies have been proposed to address this issue. They can be divided into two categories. On the one hand, some authors propose to restrict the number of active links between neighbors at each time instant [5]. On the other hand, there are authors that recommend to reduce the communication load by projecting parameter vectors onto lower dimensional spaces [16], or transmitting only partial parameter vectors [17]–[19]. In [17]–[19], the combination step (5) is redefined as:

$$\boldsymbol{w}_{k,i} = a_{kk} \boldsymbol{\psi}_{k,i} + \sum_{\ell \in \mathcal{N}_k \setminus \{k\}} a_{\ell k} \left( \boldsymbol{H}_{\ell,i} \boldsymbol{\psi}_{\ell,i} + [\boldsymbol{I} - \boldsymbol{H}_{\ell,i}] \boldsymbol{\psi}_{k,i} \right)$$
(6)

where  $H_{\ell,i}$  is a diagonal entry-selection matrix with M ones and L-M zeros on its diagonal. This means that the nodes can use the entries of their own intermediate estimates in lieu of the ones from the neighbors that have not been communicated. Matrix  $H_{\ell,i}$  can be deterministic, or can randomly select Mentries from all entries. The literature has mainly focused on the reformulation (6) of the combination step (4) of the diffusion LMS. This means that, with the so-called *partial-diffusion LMS*, only the intermediate estimates  $\psi_{k,i}$  are partially shared by neighboring nodes. However, it is also of interest to consider the adaptation step as it accounts for a major part of information exchanges.

Our main contribution in this paper is to consider and compress every transmitted information over the network, namely, local estimates  $w_{i,k}$  and estimated gradients  $\hat{\nabla}_w J_k(w_{\ell,i-1})$ . In a preliminary step, we study the case where only the local estimates in (4) are partially transmitted. Next, we study the case where both the estimates and the estimated gradients are partially transmitted. This method allows a control of network communication flows by setting the numbers M and  $M_{\nabla}$ of selected entries in  $w_{i,k}$  and  $\hat{\nabla}_w J_k(w_{\ell,i-1})$ , respectively. Finally, we carry out a theoretical analysis of the algorithms in the mean and mean-square sense, and we perform numerical experiments to confirm the theoretical findings.

*Notation:* Boldface small letters denote vectors. All vectors are column vectors. Boldface capital letters denote matrices. The  $(k, \ell)$ -th entry of a matrix is denoted by  $(\cdot)_{k\ell}$ , and the  $(k, \ell)$ -th block of a block matrix is denoted by  $[\cdot]_{k\ell}$ . Matrix trace is denoted by trace $\{\cdot\}$ . The expectation operator is denoted by  $\mathbb{E}\{\cdot\}$ . The identity matrix of size N is denoted by  $\mathbb{I}_N$ , and the all-one vector of length N is denoted by  $\mathbb{1}_N$ . 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 cardinality. The operator col $\{\cdot\}$  stacks its vector arguments on the top of each other to generate a connected vector. The notation diag $\{a, b\}$  denotes a diagonal matrix with entries a and b. Likewise, the notation diag $\{A, B\}$  denotes a block diagonal matrix with block entries A and B. The other symbols will be defined in the context where they are used.

#### II. COMPRESSED DIFFUSION LMS

We shall now introduce and analyze the stochastic behavior of the compressed diffusion LMS (CD) as an intermediate step towards the doubly-compressed diffusion LMS (DCD). For the sake of simplicity, we shall consider that C is doubly stochastic. We shall set A to the identity matrix  $I_N$ , with the purpose of limiting the network load. The compressed diffusion LMS algorithm is defined as:

$$\boldsymbol{w}_{k,i} = \boldsymbol{w}_{k,i-1} + \mu_k \sum_{\ell \in \mathcal{N}_k} c_{\ell,k} \boldsymbol{u}_{\ell,i} \big[ d_\ell(i) - \boldsymbol{u}_{\ell,i}^\top (\boldsymbol{H}_{k,i} \boldsymbol{w}_{k,i-1} + (\boldsymbol{I}_L - \boldsymbol{H}_{k,i}) \boldsymbol{w}_{\ell,i-1}) \big]$$
(7)

where  $H_{\ell,i} = \text{diag}\{h_{\ell,i}\}$ . We shall assume that  $h_{\ell,i}$  is a  $L \times 1$  binary vector, generated by randomly setting M of its L entries to 1, and the other L - M entries to 0. We shall also assume that all possible outcomes for  $h_{\ell,i}$  are equally likely, and i.i.d. over time and space. This leads to:

$$\mathbb{E}\{\boldsymbol{H}_{\ell,i}\} = \frac{M}{L}\boldsymbol{I}_L \tag{8}$$

and, given any  $L \times L$  matrix  $\Sigma$ ,

$$\mathbb{E}\{\boldsymbol{H}_{\ell,i}\boldsymbol{\Sigma}\boldsymbol{H}_{k,i}\} =$$
(9)  
$$\begin{cases} \frac{M}{L}\left(\left(1 - \frac{M-1}{L-1}\right)\boldsymbol{I}_{L} \odot \boldsymbol{\Sigma} + \frac{M-1}{L-1}\boldsymbol{\Sigma}\right) & \text{if } \ell = k \\ \left(\frac{M}{L}\right)^{2}\boldsymbol{\Sigma} & \text{otherwise} \end{cases}$$

where  $\odot$  denotes the Hadamard entrywise product. Due to the constraint  $\mathbf{1}_{L}^{\top} \boldsymbol{h}_{\ell,i} = M$ , node  $\ell$  receives exactly M entries from its neighbors at each time instant *i*. The missing (L-M)entries are filled in with estimates that are available at node  $\ell$ .

Assumption 1 The regression vectors  $u_{k,i}$  arise from a zero-mean random process that is temporally white and spatially independent. A direct consequence of this assumption is that  $u_{k,i}$  is independent of  $w_{\ell,j}$  for all  $\ell$  and j < i.

Assumption 2 The matrices  $H_{i,k}$  arise from a random process that is temporally white, spatially independent, and independent of any other process.

We introduce the  $L \times 1$  error vectors:

$$\widetilde{\boldsymbol{w}}_{k,i} = \boldsymbol{w}^o - \boldsymbol{w}_{k,i} \tag{10}$$

and we collect them from across all nodes into the vectors:

$$\widetilde{\boldsymbol{w}}_i = \operatorname{col}\{\widetilde{\boldsymbol{w}}_{1,i}, \widetilde{\boldsymbol{w}}_{2,i}, \dots, \widetilde{\boldsymbol{w}}_{N,i}\}$$
(11)

Let  $\boldsymbol{R}_{u_{\ell},i} = \boldsymbol{u}_{\ell,i} \boldsymbol{u}_{\ell,i}^{\top}$ . We also introduce:

$$\mathcal{M} = \operatorname{diag}\{\mu_1 I_L, \mu_2 I_L, \dots, \mu_N I_L\}$$
(12)

$$\boldsymbol{\mathcal{R}}_{i} = \operatorname{diag}\left\{\sum_{\ell \in \mathcal{N}_{1}} c_{\ell,1} \boldsymbol{R}_{u_{\ell},i}, \dots, \sum_{\ell \in \mathcal{N}_{N}} c_{\ell,N} \boldsymbol{R}_{u_{\ell},i}\right\}$$
(13)

$$\mathcal{R}_{u,i} = \operatorname{diag}\{R_{u_1,i}, R_{u_2,i}, \dots, R_{u_N,i}\}$$
(14)

$$\mathcal{H}_i = \operatorname{diag}\{\boldsymbol{H}_{1,i}, \boldsymbol{H}_{2,i}, \dots, \boldsymbol{H}_{N,i}\}$$
(15)

$$\mathcal{C} = C \otimes I_L \tag{16}$$

where  $\otimes$  denotes the Kronecker product. Finally, we introduce the  $N \times N$  block matrix  $\mathcal{R}_{m,i}$  with each block defined as:

$$[\mathcal{R}_{I-H,i}]_{k\ell} = c_{\ell k} \mathbf{R}_{u_{\ell},i} (\mathbf{I}_L - \mathbf{H}_{k,i})$$
(17)

Using recursion (7) and definitions (10), (11), we write:

$$\widetilde{\boldsymbol{w}}_i = \boldsymbol{\mathcal{B}}_i \widetilde{\boldsymbol{w}}_{i-1} - \boldsymbol{\mathcal{G}} \boldsymbol{s}_i \tag{18}$$

where

$$\mathcal{B}_i = I_{NL} - \mathcal{M}\mathcal{R}_i \mathcal{H}_i - \mathcal{M}\mathcal{R}_{I-H,i}$$
(19)

$$\mathcal{G} = \mathcal{M}\mathcal{C}^{\top}$$
(20)

$$s_i = col\{u_{1,i}v_1(i), u_{2,i}v_2(i), \dots, u_{N,i}v_N(i)\}$$
 (21)

A. Mean convergence

Taking expectations of both sides of recursion (18), using Assumptions 1 and 2, and  $\mathbb{E}{s_i} = 0$ , we find:

$$\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i}\} = [\boldsymbol{I}_{NL} - \mathbb{E}\{\mathcal{M}\mathcal{R}_{i}\mathcal{H}_{i}\} - \mathbb{E}\{\mathcal{M}\mathcal{R}_{I-H,i}\}]\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i-1}\} - \mathbb{E}\{\mathcal{M}\mathcal{C}^{\top}\boldsymbol{s}_{i}\} = \left[\boldsymbol{I}_{NL} - \frac{M}{L}\mathcal{M}\mathcal{R} - (1 - \frac{M}{L})\mathcal{M}\mathcal{C}^{\top}\mathcal{R}_{u}\right]\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i-1}\}$$
(22)

where

$$\boldsymbol{\mathcal{R}} = \mathbb{E}\{\boldsymbol{\mathcal{R}}_i\} = \text{diag}\{\boldsymbol{R}_1, \dots, \boldsymbol{R}_N\}$$
(23)

$$\mathcal{R}_{u} = \mathbb{E}\{\mathcal{R}_{u,i}\} = \operatorname{diag}\{\mathcal{R}_{u_{1}}, \mathcal{R}_{u_{2}}, \dots, \mathcal{R}_{u_{N}}\}$$
(24)

with

$$\boldsymbol{R}_{k} = \sum_{\ell \in \mathcal{N}_{k}} c_{\ell,k} \boldsymbol{R}_{u_{\ell}}$$
(25)

From (22), we observe that the algorithm (7) asymptotically converges in the mean toward  $w^{o}$  if, and only if,

$$\rho\left(\boldsymbol{I}_{NL} - \frac{M}{L}\boldsymbol{\mathcal{M}}\boldsymbol{\mathcal{R}} - \left(1 - \frac{M}{L}\right)\boldsymbol{\mathcal{M}}\boldsymbol{\mathcal{C}}^{\top}\boldsymbol{\mathcal{R}}_{u}\right) < 1 \qquad (26)$$

where  $\rho(\cdot)$  denotes the spectral radius of its matrix argument. We know that  $\rho(\mathbf{X}) < ||\mathbf{X}||$  for any induced norm. Then:

$$\rho\left(\boldsymbol{I}_{NL} - \frac{M}{L}\mathcal{M}\mathcal{R} - \left(1 - \frac{M}{L}\right)\mathcal{M}\mathcal{C}^{\top}\mathcal{R}_{u}\right) \\
\leq \|\boldsymbol{I}_{NL} - \frac{M}{L}\mathcal{M}\mathcal{R} - \left(1 - \frac{M}{L}\right)\mathcal{M}\mathcal{C}^{\top}\mathcal{R}_{u}\|_{b,\infty} \\
\leq \max_{\ell,k} \|\boldsymbol{I}_{L} - \frac{M}{L}\mu_{k}\boldsymbol{R}_{k} - \left(1 - \frac{M}{L}\right)\mu_{k}c_{\ell k}\boldsymbol{R}_{u_{\ell}}\| \quad (27)$$

where  $\|\cdot\|_{b,\infty}$  denotes the block maximum norm. This leads us to the following condition:

$$\mu_k < \frac{2}{\max_{\ell \in \mathcal{N}_k} \lambda_{\max}\left(\frac{M}{L} \boldsymbol{R}_k + \left(1 - \frac{M}{L}\right) c_{\ell k} \boldsymbol{R}_{u_\ell}\right)}$$
(28)

where  $\lambda_{\max}(\cdot)$  stands for the maximum eigenvalue of its matrix argument. Furthermore, by Weyl's theorem we have:

$$\mu_k < \frac{2}{\frac{M}{L}\lambda_{\max}(\boldsymbol{R}_k) + \left(1 - \frac{M}{L}\right) \max_{\ell \in \mathcal{N}_k} [c_{\ell k} \lambda_{\max}(\boldsymbol{R}_{u_\ell})]}$$
(29)

since  $\mathbf{R}_k$  and  $\mathbf{R}_{u_\ell}$  are Hermitian matrices.

#### B. Mean-square stability

We shall now analyze the mean-square deviation  $\mathbb{E}\{\|\widetilde{\boldsymbol{w}}_i\|_{\boldsymbol{\Sigma}}^2\}$ where  $\boldsymbol{\Sigma}$  denotes a nonnegative  $N \times N$  block diagonal matrix. The choice of matrix  $\boldsymbol{\Sigma}$  determines the type of information extracted from the network and nodes. Using the independence assumption and (18), we find:

$$\mathbb{E}\|\widetilde{\boldsymbol{w}}_{i}\|_{\boldsymbol{\Sigma}}^{2} = \mathbb{E}\{\widetilde{\boldsymbol{w}}_{i-1}^{\top}\boldsymbol{\mathcal{B}}_{i}^{\top}\boldsymbol{\Sigma}\boldsymbol{\mathcal{B}}_{i}\widetilde{\boldsymbol{w}}_{i-1}\} + \mathbb{E}\{\boldsymbol{s}_{i}^{\top}\boldsymbol{\mathcal{G}}^{\top}\boldsymbol{\Sigma}\boldsymbol{\mathcal{G}}\boldsymbol{s}_{i}\} (30)$$

On the one hand, note that the second term on the RHS of (30) can be written as:

$$\mathbb{E}\{\boldsymbol{s}_i^{\top}\boldsymbol{\mathcal{G}}^{\top}\boldsymbol{\Sigma}\boldsymbol{\mathcal{G}}\boldsymbol{s}_i\} = \operatorname{trace}\{\boldsymbol{\mathcal{G}}\boldsymbol{\mathcal{S}}\boldsymbol{\mathcal{G}}^{\top}\boldsymbol{\Sigma}\}$$
(31)

where

$$\boldsymbol{\mathcal{S}} = \operatorname{diag}(\sigma_{v,1}^2 \boldsymbol{R}_{u,1}, \dots, \sigma_{v,N}^2 \boldsymbol{R}_{u,N})$$
(32)

On the other hand, the first term on the RHS of (30) can be expressed as:

$$\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i-1}^{\top}\boldsymbol{\mathcal{B}}_{i}^{\top}\boldsymbol{\Sigma}\boldsymbol{\mathcal{B}}_{i}\widetilde{\boldsymbol{w}}_{i-1}\}=\mathbb{E}\|\widetilde{\boldsymbol{w}}_{i-1}\|_{\boldsymbol{\Sigma}'}^{2}$$
(33)

where the weighting matrix  $\Sigma'$  is defined as

$$\Sigma' = \mathbb{E} \{ \mathcal{B}_i^\top \Sigma \mathcal{B}_i \}$$
  
=  $\Sigma - \frac{M}{L} \Sigma \mathcal{M} \mathcal{R} - (1 - \frac{M}{L}) \Sigma \mathcal{M} \mathcal{C}^\top \mathcal{R}_u$   
-  $\frac{M}{L} \mathcal{R}^\top \mathcal{M} \Sigma - (1 - \frac{M}{L}) \mathcal{R}_u \mathcal{C} \mathcal{M} \Sigma$   
+  $P_1 + P_2 + P_2^\top + P_3$  (34)

with

 $\mathbb{E}\{$ 

$$\boldsymbol{P}_1 = \mathbb{E}\{\boldsymbol{\mathcal{H}}_i \boldsymbol{\mathcal{R}}_i^{\top} \boldsymbol{\mathcal{M}} \boldsymbol{\Sigma} \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{R}}_i \boldsymbol{\mathcal{H}}_i\}$$
(35)

$$\boldsymbol{P}_{2} = \mathbb{E}\{\boldsymbol{\mathcal{H}}_{i}\boldsymbol{\mathcal{R}}_{i}^{\top}\boldsymbol{\mathcal{M}}\boldsymbol{\Sigma}\boldsymbol{\mathcal{M}}\boldsymbol{\mathcal{R}}_{I-H,i}\}$$
(36)

$$\boldsymbol{P}_{3} = \mathbb{E}\{\boldsymbol{\mathcal{R}}_{I-H,i}^{\top}\boldsymbol{\mathcal{M}}\boldsymbol{\Sigma}\boldsymbol{\mathcal{M}}\boldsymbol{\mathcal{R}}_{I-H,i}\}$$
(37)

To evaluate  $P_1$ , we write:

$$\boldsymbol{P}_{1} = \mathbb{E}_{\mathcal{H}} \Big\{ \boldsymbol{\mathcal{H}}_{i} \, \mathbb{E}_{\mathcal{R}} \big\{ \boldsymbol{\mathcal{R}}_{i}^{\top} \boldsymbol{\mathcal{M}} \boldsymbol{\Sigma} \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{R}}_{i} \big\} \, \boldsymbol{\mathcal{H}}_{i} \Big\}$$
(38)

Consider first  $\mathbb{E}_{\mathcal{H}}{\{\mathcal{H}_i \Pi \mathcal{H}_i\}}$  where  $\Pi$  denotes any  $NL \times NL$  deterministic matrix. It can be shown that:

$$\begin{aligned} \mathcal{H}_{i} \Pi \mathcal{H}_{i} \\ &= \alpha_{1} \left( \boldsymbol{I}_{N} \otimes \boldsymbol{1}_{LL} \right) \odot \boldsymbol{\Pi} + \alpha_{2} \, \boldsymbol{I}_{NL} \odot \boldsymbol{\Pi} + \alpha_{3} \, \boldsymbol{\Pi} \end{aligned}$$
(39)

where  $\mathbf{1}_{LL}$  denotes an all-one  $L \times L$  matrix, and

$$\alpha_1 = \frac{M}{L} \left( \frac{M-1}{L-1} - \frac{M}{L} \right) \ \alpha_2 = \frac{M}{L} \left( 1 - \frac{M-1}{L-1} \right) \ \alpha_3 = \left( \frac{M}{L} \right)^2$$

Next, consider the inner expectation in the RHS of (38). The evaluation of this expectation depends on higher-order moments of the regression data. While we can continue with the analysis by calculating these terms, it is sufficient for the exposition to focus on the case of sufficiently small step-sizes where a reasonable approximation is [1]:

$$\mathbb{E}\{\mathcal{R}_i^{\top} \mathcal{M} \Sigma \mathcal{M} \mathcal{R}_i\} = \mathcal{R}^{\top} \mathcal{M} \Sigma \mathcal{M} \mathcal{R}$$
(40)

Substituting (39) and (40) into (38) leads to  $P_1$ . Consider  $P_2$ . Using the same higher-order approximation as above, we have:

$$[\mathbf{P}_{2}]_{k\ell} = \frac{M}{L} \mathbf{R}_{k} [\mathbf{M} \Sigma \mathbf{M}]_{kk} c_{\ell k} \mathbf{R}_{u_{\ell}} - \mathbb{E} \{ \mathbf{H}_{k,i} \mathbf{R}_{k,i} [\mathbf{M} \Sigma \mathbf{M}]_{kk} c_{\ell k} \mathbf{R}_{u_{\ell},i} \mathbf{H}_{k,i} \}$$
(41)

The second term in RHS of (41) is of the form:

$$[\varphi(\mathbf{\Pi})]_{k\ell} = \mathbb{E}\{\boldsymbol{H}_{k,i}[\mathbf{\Pi}]_{k\ell}\boldsymbol{H}_{k,i}\}$$
(42)

with  $[\Pi]_{k\ell} = \mathbf{R}_{k,i}[\mathcal{M}\Sigma\mathcal{M}]_{kk}c_{\ell k}\mathbf{R}_{u_{\ell},i}$ . It can be shown that:

$$\varphi(\mathbf{\Pi}) = \alpha_2 \left( \mathbf{1}_{NN} \otimes \boldsymbol{I}_L \right) \odot \mathbf{\Pi} + (\alpha_1 + \alpha_3) \mathbf{\Pi}$$
(43)

This leads to:

$$\boldsymbol{P}_{2} = \frac{M}{L} \boldsymbol{\mathcal{R}} \boldsymbol{\mathcal{M}} \boldsymbol{\Sigma} \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{C}}^{\top} \boldsymbol{\mathcal{R}}_{u} - \varphi (\boldsymbol{\mathcal{R}} \boldsymbol{\mathcal{M}} \boldsymbol{\Sigma} \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{C}}^{\top} \boldsymbol{\mathcal{R}}_{u}) \quad (44)$$

Finally, using the same higher-order approximation as above with  $P_3$  leads to:

$$[\boldsymbol{P}_3]_{k\ell} = \sum_{m=1}^{N} \mathbb{E}\{[\boldsymbol{\mathcal{R}}_{I-H,i}]_{km}^{\top} [\boldsymbol{\mathcal{M}} \boldsymbol{\Sigma} \boldsymbol{\mathcal{M}}]_{mm} [\boldsymbol{\mathcal{R}}_{I-H,i}]_{m\ell}\}$$
(45)

and

$$P_{3} = (1 - 2\frac{M}{L}) \mathcal{R}_{u} \mathcal{CM} \Sigma \mathcal{M} \mathcal{C}^{\top} \mathcal{R}_{u} + \varphi (\mathcal{R}_{u} \mathcal{CM} \Sigma \mathcal{M} \mathcal{C}^{\top} \mathcal{R}_{u})$$
(46)

Following the same reasoning as in [1] allows to express matrix  $\Sigma'$  in a vector form as

$$\boldsymbol{\sigma}' = \boldsymbol{\mathcal{F}}\boldsymbol{\sigma} \tag{47}$$

where  $\boldsymbol{\sigma} = \operatorname{vec}(\boldsymbol{\Sigma}), \ \boldsymbol{\sigma}' = \operatorname{vec}(\boldsymbol{\Sigma}')$  and:

$$\mathcal{F} = \mathbf{I}_{(NM)^2} - \frac{M}{L} (\mathcal{RM} \otimes \mathbf{I}_{NL}) - (1 - \frac{M}{L}) (\mathcal{R}_u \mathcal{CM} \otimes \mathbf{I}_{NL}) - \frac{M}{L} (\mathbf{I}_{NL} \otimes \mathcal{RM})$$
(48)  
  $- (1 - \frac{M}{L}) (\mathbf{I}_{NL} \otimes \mathcal{R}_u \mathcal{CM}) + \mathbf{Z}_1 + \mathbf{Z}_2 + \mathbf{Z}_{2^{\top}} + \mathbf{Z}_4$ 

where matrices  $Z_j$  are obtained by applying the vec(·) operator on  $P_j$  and using the following property:

$$\operatorname{vec}(\boldsymbol{ABC}) = (\boldsymbol{C}^{\top} \otimes \boldsymbol{A})\operatorname{vec}(\boldsymbol{B})$$
 (49)

The analytical expressions of  $Z_i$  are not provided in this paper due to the lake of space. They will be made available in an extended version of this paper. Substituting (31) and (33) in (30), and applying the vec operation to both sides, we get:

$$\mathbb{E}\|\widetilde{\boldsymbol{w}}_i\|_{\boldsymbol{\sigma}}^2 = \mathbb{E}\|\widetilde{\boldsymbol{w}}_{i-1}\|_{\boldsymbol{\mathcal{F}}\boldsymbol{\sigma}}^2 + [\operatorname{vec}(\boldsymbol{\mathcal{Y}}^{\top})]^{\top}\boldsymbol{\sigma}$$
(50)

where

$$\boldsymbol{\mathcal{Y}} = \boldsymbol{\mathcal{G}} \, \boldsymbol{\mathcal{S}} \, \boldsymbol{\mathcal{G}}^\top \tag{51}$$

#### III. DOUBLY-COMPRESSED DIFFUSION LMS

The compressed diffusion LMS studied before only partially transmit local estimates while the estimated gradient vectors are fully transmitted through the network. We shall now introduce the doubly-compressed diffusion LMS algorithm, which offers an additional compression layer by transmitting partial estimated gradient vectors. Node k receives  $M_{\nabla}$  entries of  $\hat{\nabla}_w J_k(\boldsymbol{w}_{\ell,i-1})$  from its neighbors at each time instant *i*. The missing  $(L - M_{\nabla})$  entries are filled in with the gradient estimates that are available at node k. The algorithm can be formulated as:

$$\boldsymbol{w}_{k,i} = \boldsymbol{w}_{k,i-1} + \mu_k \sum_{\ell \in \mathcal{N}_k} c_{\ell,k} \boldsymbol{Q}_{\ell,i} \boldsymbol{u}_{\ell,i} [d_\ell(i) \\ - \boldsymbol{u}_{\ell,i}^\top (\boldsymbol{H}_{k,i} \boldsymbol{w}_{k,i-1} + (\boldsymbol{I}_L - \boldsymbol{H}_{k,i}) \boldsymbol{w}_{\ell,i-1})] \\ + \mu_k \sum_{\ell \in \mathcal{N}_k} c_{\ell,k} (\boldsymbol{I}_L - \boldsymbol{Q}_{\ell,i}) \boldsymbol{u}_{k,i} [d_k(i) \\ - \boldsymbol{u}_{k,i}^\top (\boldsymbol{H}_{k,i} \boldsymbol{w}_{k,i-1} + (\boldsymbol{I}_L - \boldsymbol{H}_{k,i}) \boldsymbol{w}_{\ell,i-1})]$$
(52)

where  $Q_{\ell,i}$  is a stochastic selection matrix that has the same properties as  $H_{k,i}$  defined in the previous section.

Assumption 3 The matrices  $Q_{i,\ell}$  arise from a random process that is temporally white, spatially independent, and independent of any other process.

Proceeding as in the previous section, we find that:

$$\widetilde{\boldsymbol{w}}_i = \boldsymbol{\mathcal{B}}_i \widetilde{\boldsymbol{w}}_{i-1} - \boldsymbol{\mathcal{G}}_i \boldsymbol{s}_i \tag{53}$$

where

$$\mathcal{B}_{i} = I_{NL} - \mathcal{MR}_{Q,i}\mathcal{H}_{i} - \mathcal{MQ}'_{i}\mathcal{R}_{u,i}\mathcal{H}_{i} - \mathcal{MR}_{Q(I-H),i} - \mathcal{MR}_{(I-Q)(I-H),i}$$
(54)

$$\boldsymbol{\mathcal{G}}_{i} = \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{Q}}_{i} \boldsymbol{\mathcal{C}}^{\top} + \boldsymbol{\mathcal{M}} \boldsymbol{\mathcal{Q}}_{i}^{\prime}$$
(55)



Fig. 1: (left) Network topology. (right) Variance  $\sigma_{u_k}^2$  of regressors in Experiment 1 (top) and in Experiment 2 (bottom).

with

$$\begin{split} \boldsymbol{\mathcal{R}}_{Q,i} &= \operatorname{diag}\left\{\sum_{\ell \in N_1} c_{\ell 1} \boldsymbol{Q}_{\ell,i} \boldsymbol{R}_{u_{\ell},i}, \dots, \sum_{\ell \in N_N} c_{\ell N} \boldsymbol{Q}_{\ell,i} \boldsymbol{R}_{u_{\ell},i}\right\} \\ \boldsymbol{\mathcal{Q}}_i &= \operatorname{diag}\{\boldsymbol{Q}_{1,i}, \boldsymbol{Q}_{2,i}, \dots, \boldsymbol{Q}_{N,i}\} \\ \boldsymbol{\mathcal{Q}}'_i &= \operatorname{diag}\left\{\sum_{\ell \in N_1} c_{\ell 1} (\boldsymbol{I}_L - \boldsymbol{Q}_{\ell,i}), \dots, \sum_{\ell \in N_N} c_{\ell N} (\boldsymbol{I}_L - \boldsymbol{Q}_{\ell,i})\right\} \\ [\boldsymbol{\mathcal{R}}_{Q(I-H),i}]_{k\ell} &= c_{\ell k} \boldsymbol{Q}_{\ell,i} \boldsymbol{R}_{u_{\ell},i} (\boldsymbol{I}_L - \boldsymbol{H}_{k,i}) \\ [\boldsymbol{\mathcal{R}}_{(I-Q)(I-H),i}]_{k\ell} &= c_{\ell k} (\boldsymbol{I}_L - \boldsymbol{Q}_{\ell,i}) \boldsymbol{R}_{u_k,i} (\boldsymbol{I}_L - \boldsymbol{H}_{k,i}) \end{split}$$

The performance analysis of this algorithm is similar to the analysis of compressed diffusion LMS, although more tedious. The following result can be obtained following the same steps.

#### A. Convergence in mean

Taking expectations of both sides of (53), we obtain:

$$\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i}\} = \left(\boldsymbol{I}_{NL} - \frac{MM_{\nabla}}{L^{2}}\mathcal{M}\mathcal{R} - \frac{M}{L}\left(1 - \frac{M_{\nabla}}{L}\right)\mathcal{M}\mathcal{R}_{u} - \frac{M_{\nabla}}{L}\left(1 - \frac{M}{L}\right)\mathcal{M}\mathcal{C}^{\top}\mathcal{R}_{u}$$
(56)  
$$-\left(1 - \frac{M_{\nabla}}{L}\right)\left(1 - \frac{M}{L}\right)\mathcal{M}\mathcal{R}_{u}\mathcal{C}^{\top}\right)\mathbb{E}\{\widetilde{\boldsymbol{w}}_{i-1}\}$$

From (56), the algorithm is stable as long as the following condition is met:

$$\mu_k < \frac{2}{\lambda_{\max,k}} \tag{57}$$

where

$$\lambda_{k} = \frac{MM_{\nabla}}{L^{2}} \lambda_{\max}(\boldsymbol{R}_{k}) + \frac{M}{L} \left(1 - \frac{M_{\nabla}}{L}\right) \lambda_{\max}(\boldsymbol{R}_{u_{k}}) + \frac{M_{\nabla}}{L} \left(1 - \frac{M}{L}\right) \max_{\ell \in \mathcal{N}_{k}} [c_{\ell k} \lambda_{\max}(\boldsymbol{R}_{u_{\ell}})] + \left(1 - \frac{M_{\nabla}}{L}\right) \left(1 - \frac{M}{L}\right) \max_{\ell \in \mathcal{N}_{k}} [c_{\ell k} \lambda_{\max}(\boldsymbol{R}_{u_{k}})]$$
(58)

Due to lake of space and of the complexity of the analysis in the mean-square sense, this analysis is not reported here. It will be made available in an extended version of this article.

#### **IV. SIMULATION RESULTS**

We shall now evaluate the accuracy of the theoretical models with simulations. We performed two experiments to determine the performance and efficiency of both algorithms. First, we considered a small network in order to validate the theoretical models. Then, we considered a larger network and high dimensional measurements in order to test the algorithms. For both experiments, parameter vectors  $w^o$  were generated from



Fig. 2: Theoretical model and simulated curves (left), evolution of the MSD as a function of the compression ratio for compressed diffusion LMS (center), and doubly-compressed diffusion LMS (right).

a zero-mean Gaussian distribution. The input data  $u_{k,i}$  were drawn from zero-mean Gaussian distributions with covariance  $R_{u,k} = \sigma_{u,k}^2 I_L$  reported in Fig. 1. The weighting matrices C were generated using the Metropolis rule [1]. Noises  $v_k(i)$ were zero-mean, i.i.d. and Gaussian distributed with variance  $\sigma_{v,k}^2 = 10^{-3}$ . The step-sizes were set to  $\mu_k = 10^{-3}$ . The simulation results were averaged over 100 Monte Carlo runs.

1) Experiment 1: Due to the high dimensionality of the theoretical models, we considered the network of N = 10 nodes shown in Fig. 1 (left). We set L = 5, M = 3 and  $M_{\nabla} = 1$ . This resulted in compression ratios of  $\frac{10}{8}$  and  $\frac{5}{2}$  for the compressed and the doubly-compressed diffusion LMS, respectively. It can be observed in Fig. 2 (left) that the theoretical models fit accurately the simulated results for both algorithms. Unsurprisingly, diffusion LMS outperformed its compressed counterparts at the expense of a higher communication load.

2) Experiment 2: Since compression is relevant for relatively large data flows, we considered a network with N = 50 agents. Data dimension was set to L = 50. Figure 2 illustrates the performance of the algorithms for different compression ratios. Note that maximum compression ratio that can be reached with the compressed diffusion LMS is  $\frac{100}{55}$  since estimated gradient vectors are fully transmitted. Doubly-compressed diffusion LMS is more flexible since the compression ratio can be controlled via M and  $M_{\nabla}$ .

#### V. CONCLUSION

In order to preserve energy resources, we investigated compression techniques for diffusion LMS that consist of exchanging partial local estimates. We carried out an analysis of the stochastic behavior of the proposed algorithms in the mean and mean-square sense. Finally we provided simulation results to illustrate the accuracy of the theoretical models and the performance of the compression layer.

#### REFERENCES

- A. H. Sayed, "Diffusion adaptation over networks," in *Academic Press Libraray in Signal Processing*, R. Chellapa and S. Theodoridis, Eds. Elsevier, 2014. Also available as arXiv:1205.4220 [cs.MA], May 2012., pp. 322–454.
   A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. J. Towfic, "Diffusion
- [2] Ä. 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.

- [3] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, 2014.
- [4] X. Zhao, S.-Y. Tu, and A. H. Sayed, "Diffusion adaptation over networks under imperfect information exchange and non-stationary data," *IEEE Transactions on Signal Processing*, vol. 60, no. 7, pp. 3460–3475, 2012.
- [5] R. Arablouei, S. Werner, K. Doğançay, and Y.-F. Huang, "Analysis of a reduced-communication diffusion LMS algorithm," *Signal Processing*, vol. 117, pp. 355–361, 2015.
- [6] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multiagent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [7] M. G. Rabbat and R. D. Nowak, "Quantized incremental algorithms for distributed optimization," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 4, pp. 798–808, Apr. 2005.
- [8] C. G. Lopes and A. H. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Transactions on Signal Processing*, vol. 55, no. 8, pp. 4064–4077, 2007.
- [9] J. Chen, C. Richard, A. O. Hero, and A. H. Sayed, "Diffusion lms for multitask problems with overlapping hypothesis subspaces," in *Proc. IEEE MLSP'14*, Reims, France, 2014, pp. 1–6.
- [10] 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.
- [11] —, "Diffusion LMS over multitask networks," *IEEE Transactions on Signal Processing*, vol. 63, no. 11, pp. 2733–2748, 2015.
- [12] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Multitask diffusion adaptation over asynchronous networks," *IEEE Transactions on Signal Processing*, vol. 64, no. 11, pp. 2835–2850, 2016.
- [13] —, "Proximal multitask learning over networks with sparsity-inducing coregularization," *IEEE Transactions on Signal Processing*, vol. 64, no. 23, pp. 6329–6344, 2016.
- [14] W. Gao, J. Chen, C. Richard, and J. Huang, "Diffusion adaptation over networks with kernel least-mean-square," in *Proc. IEEE CAMSAP'15*, Cancún, Mexico, 2015.
- [15] P. Chainais and C. Richard, "Learning a common dictionary over a sensor network," in *Proc. IEEE CAMSAP'13*, Saint Martin, French West Indies, 2013.
- [16] M. O. Sayin and S. S. Kozat, "Compressive diffusion strategies over distributed networks for reduced communication load," *IEEE Transactions* on Signal Processing, vol. 62, no. 20, pp. 5308–5323, 2014.
- [17] R. Arablouei, S. Werner, Y.-F. Huang, and K. Dogancay, "Distributed least mean-square estimation with partial diffusion," *IEEE Transactions* on Signal Processing, vol. 62, no. 2, pp. 472–484, 2014.
- [18] R. Arablouei, K. Dogancay, S. Werner, and Y.-F. Huang, "Adaptive distributed estimation based on recursive least-squares and partial diffusion," *IEEE Transactions on Signal Processing*, vol. 62, no. 14, pp. 3510–3522, 2014.
- [19] V. Vadidpour, A. Rastegarnia, A. Khalili, and S. Sanei, "Partial-diffusion least mean-square estimation over networks under noisy information exchange," arXiv preprint arXiv:1511.09044, 2015.