

TRANSIENT ANALYSIS OF CLUSTERED MULTITASK DIFFUSION RLS ALGORITHM

Wei Gao[†] Jie Chen^{*} Cédric Richard[‡] Wentao Shi^{*} Qunfei Zhang^{*}

[†]School of Computer Science and Communication Engineering, Jiangsu University, Zhenjiang 212013, China

^{*}School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an 710072, China

[‡]Université Côte d'Azur, OCA, CNRS, 06108 Nice, France

Email: wei_gao@ujs.edu.cn drjie.chen@ieee.org cedric.richard@unice.fr swt@nwpu.edu.cn zhangqf@nwpu.edu.cn

ABSTRACT

In this paper, we propose a novel clustered multitask diffusion RLS (MT-DRLS) algorithm over network to further improve the performance of its counterpart, the multitask diffusion LMS (MT-DLMS) algorithm. Its transient behavior is investigated, in the mean and mean-square error sense. Simulation results illustrate the significant improvement of the MT-DRLS over the MT-DLMS in terms of convergence rate and steady-state error, as well as the accuracy of the theoretical findings.

Index Terms— Multitask adaptive learning, diffusion RLS, transient analysis, distributed estimation, multitask network.

1. INTRODUCTION

During the past decade, distributed detection, estimation, and tracking problems have attracted substantial attention in the context of adaptive networks with diffusion strategies [1–3]. Particularly, it has been found that multitask networks have a wider range of modeling ability than single-task networks. Several multitask strategies for adaptation and learning over networks have been recently proposed based on the diffusion least-mean-squares (DLMS) algorithm [4–11]. To be specific, the DLMS algorithm for multitask networks was first proposed in [4, 5], and studied in asynchronous networks in [6]. A new multitask learning formulation using a common latent representation was presented in [7], as well as a unified framework to analyze its performance. Both ℓ_1 -norm regularization and $\ell_{\infty,1}$ -norm regularization were introduced into multitask networks in [9] and [10], respectively. Recently, the performance of multitask DLMS algorithm has been analyzed in the presence of communication delays [11]. An overview of multitask learning over networks and its applications is available in [12].

The diffusion recursive least-squares (DRLS) algorithm was also extensively studied in [13–20], due to the superior performance of the RLS compared to the LMS [21, 22]. The DRLS algorithm with incremental update-then-combine diffusion strategy was initially proposed in [13], with an analysis of its steady-state performance. The diffusion bias-compensated RLS algorithm was presented in [14] to reduce residual bias, and the DRLS algorithm was considered in [15, 19] to reduce communication cost. Variants of the DRLS algorithm were successively devised to improve the performance in the context of sparse systems [16], noisy links [17], and robustness against impulsive interferences [18]. More recently, a transient analysis of DRLS algorithm was presented in [20]. To the best of our knowledge, the multitask DRLS algorithm has not been

considered so far except in [23]. This motivates us to derive in this paper the clustered multitask diffusion RLS (MT-DRLS) algorithm with adapt-then-combine (ATC) diffusion strategy. Furthermore, analytical models are derived to characterize its transient behavior in the mean and mean-square error sense. Simulation results illustrate the superiority of the MT-DRLS algorithm over the MT-DLMS and DRLS algorithms. The accuracy of the resulting transient analytical models is also investigated.

Notation: The matrix trace is denoted by $\text{tr}\{\cdot\}$. The notation \otimes denotes Kronecker product. Identity matrix of size $N \times N$ is denoted by \mathbf{I}_N , and $\mathbf{1}_N$ denotes an all-one vector of length N . The operator $\text{bdiag}\{\cdot\}$ formulates a (block) diagonal matrix with its arguments, and $\text{col}\{\cdot\}$ stacks its vector arguments on the top of each other to generate a column vector. The notation $\|\mathbf{x}\|_{\Sigma}^2$ denotes the squared norm of \mathbf{x} weighted by any positive semi-definite matrix Σ , i.e., $\|\mathbf{x}\|_{\Sigma}^2 = \mathbf{x}^T \Sigma \mathbf{x}$. \mathcal{N}_k denotes the neighborhood of node k , including k . $\mathcal{C}(k)$ and $\mathcal{C}(k)^-$ denote the cluster of nodes to which node k belongs, including k and excluding k , respectively.

2. NETWORK MODEL AND MT-DRLS ALGORITHM

2.1. Clustered Multitask Network Model

Consider a connected network consisting of K nodes, indexed with $k = 1, \dots, K$. At time instant $n \geq 0$, each node k has access to a random data pair $\{d_{k,n}, \mathbf{x}_{k,n}\}$, which is assumed to be generated by an optimal weight vector $\mathbf{w}_k^* \in \mathbb{R}^L$ at node k via the linear regression model:

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

where $d_{k,n} \in \mathbb{R}$ is the zero-mean desired signal, $\mathbf{x}_{k,n} \in \mathbb{R}^L$ denotes the regression vector with a positive-definite covariance matrix $\mathbf{R}_{x,k} = \mathbb{E}\{\mathbf{x}_{k,n} \mathbf{x}_{k,n}^T\}$, and $z_{k,n}$ is a zero-mean temporally and spatially independent noise with variance $\sigma_{z,k}^2$. We assume that all nodes are grouped into Q clusters, corresponding to Q estimation tasks. The unknown optimal weight vector \mathbf{w}_k^* are constrained to be identical within each cluster, namely, $\mathbf{w}_k^* = \mathbf{w}_{C_q}^*$ for all $k \in C_q$, where C_q denotes the cluster q , i.e., the index set of nodes in the q -th cluster. However, it is assumed that similarities exist among the neighboring clusters, i.e., $\mathbf{w}_{C_p}^* \sim \mathbf{w}_{C_q}^*$ if C_p and C_q are connected with $p \neq q$, where \sim represents a similarity relationship in some sense. Clusters C_p and C_q are connected provided that there exists at least one communication link connecting a node from one cluster to a node in the other cluster.

2.2. Clustered Multitask Diffusion RLS Algorithm

In the context of clustered multitask networks, the objective is to estimate the unknown parameter vectors $\{\mathbf{w}_{C_q}^*\}_{q=1}^Q$. For each node k

This work was supported in part by the National Natural Science Foundation of China Grants (62171205, 62171380).

in cluster $\mathcal{C}(k)$ and nodes of other clusters connecting to node k , we first introduce an intermediate similarity promotion equation to promote the similarities of weight vectors between neighboring clusters:

$$\mathbf{w}_{k,n-\frac{1}{2}} = \mathbf{w}_{k,n-1} - \gamma \sum_{\ell \in \mathcal{N}_k \setminus \mathcal{C}(k)} \frac{\rho_{k\ell} + \rho_{\ell k}}{2} (\mathbf{w}_{k,n-1} - \mathbf{w}_{\ell,n-1}) \quad (2)$$

with the set difference \setminus and the strength parameter $\gamma \geq 0$, and where $\mathbf{w}_{k,n-1}$ denotes the local estimate of \mathbf{w}_k^* . Here, the non-negative weight coefficients $\rho_{k\ell}$ are chosen to satisfy the conditions [5, 6]:

$$\sum_{\ell \in \mathcal{N}_k \setminus \mathcal{C}(k)} \rho_{k\ell} = 1, \text{ and } \begin{cases} \rho_{k\ell} > 0, & \text{if } \ell \in \mathcal{N}_k \setminus \mathcal{C}(k), \\ \rho_{k\ell} \geq 0, & \\ \rho_{k\ell} = 0, & \text{otherwise.} \end{cases} \quad (3)$$

We collect the above coefficients into the random $K \times K$ right-stochastic matrix Θ with (k, ℓ) -th entry $\rho_{k\ell}$. Note that the weight coefficients between pairs of nodes can be assumed to be symmetric or asymmetric. We use the symmetric case here for simplicity.

Let $\psi_{k,n}$ denote the intermediate estimate of \mathbf{w}_k^* . Since the weight vector $\mathbf{w}_{k,n-\frac{1}{2}}$ that considers the similarities of weight vectors between neighboring clusters is a good guess for the intermediate estimate $\psi_{k,n}$ as a prior information, we then consider a special case of local least-squares problem based on the collected data at node k only for current time instant n as in [18, 20]:

$$\psi_{k,n} = \arg \min_{\psi_k \in \mathbb{R}^L} \left\{ \|\psi_k - \mathbf{w}_{k,n-\frac{1}{2}}\|_{\Lambda_{k,n}}^2 + (d_{k,n} - \mathbf{x}_{k,n}^\top \psi_k)^2 \right\} \quad (4)$$

with the positive-definite weighting matrix $\Lambda_{k,n}$ and the time-averaged autocorrelation matrix $\Phi_{k,n}$ of input data for node k at time instant n , namely,

$$\Lambda_{k,n} = \Phi_{k,n} - \mathbf{x}_{k,n} \mathbf{x}_{k,n}^\top, \quad (5)$$

$$\Phi_{k,n} = \lambda \Phi_{k,n-1} + \mathbf{x}_{k,n} \mathbf{x}_{k,n}^\top, \quad (6)$$

where $0 \ll \lambda < 1$ is the forgetting factor, and the initial condition is $\Phi_{k,0} = \delta \mathbf{I}_L$ with a small positive value δ . Considering variable substitutions $\mathbf{v} = \psi_k - \mathbf{w}_{k,n-\frac{1}{2}}$ and $b = d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n-\frac{1}{2}}$, problem (4) can be reformulated as:

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathbb{R}^L} \left\{ \mathbf{v}^\top \Lambda_{k,n} \mathbf{v} + (b - \mathbf{x}_{k,n}^\top \mathbf{v})^2 \right\}. \quad (7)$$

Setting the derivative of (7) with respect to \mathbf{v} to zero, we obtain:

$$\Phi_{k,n} \mathbf{v} = b \mathbf{x}_{k,n}. \quad (8)$$

By the definition of \mathbf{v} , applying the matrix inversion lemma to the right hand side (r.h.s.) of (6), then (8) can be rewritten as:

$$\psi_{k,n} = \mathbf{w}_{k,n-\frac{1}{2}} + b \mathbf{P}_{k,n} \mathbf{x}_{k,n} \quad (9)$$

with the definition $\Phi_{k,n} = \mathbf{P}_{k,n}^{-1}$, where the well-known recursion for matrix $\mathbf{P}_{k,n}$ at node k and time instant n is given by [21, 22]:

$$\mathbf{P}_{k,n} = \lambda^{-1} \left(\mathbf{P}_{k,n-1} - \frac{\mathbf{P}_{k,n-1} \mathbf{x}_{k,n} \mathbf{x}_{k,n}^\top \mathbf{P}_{k,n-1}}{\lambda + \mathbf{x}_{k,n}^\top \mathbf{P}_{k,n-1} \mathbf{x}_{k,n}} \right) \quad (10)$$

with the initial condition $\mathbf{P}_{k,0} = \delta^{-1} \mathbf{I}_L$. Specifically, b can be approximated by estimation error at node k and time instant n , i.e.,

$$b = d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n-\frac{1}{2}} \approx d_{k,n} - \mathbf{x}_{k,n}^\top \mathbf{w}_{k,n-1} = e_{k,n}. \quad (11)$$

Substituting (2) and (11) into (9), we arrive at the adaptive update step of clustered MT-DRLS algorithm:

$$\psi_{k,n} = \mathbf{w}_{k,n-1} + \mathbf{P}_{k,n} \mathbf{x}_{k,n} e_{k,n} - \gamma \sum_{\ell \in \mathcal{N}_k \setminus \mathcal{C}(k)} \frac{\rho_{k\ell} + \rho_{\ell k}}{2} (\mathbf{w}_{k,n-1} - \mathbf{w}_{\ell,n-1}). \quad (12)$$

Let us introduce the entire intermediate estimated vector and the block column matrix with individual entries the $L \times L$ identity matrix, which are defined as follows:

$$\psi_n = \text{col}\{\psi_{1,n}, \dots, \psi_{K,n}\} \in \mathbb{R}^{KL}, \quad (13)$$

$$\mathbf{H} = \text{col}\{\mathbf{I}_L, \dots, \mathbf{I}_L\} \in \mathbb{R}^{KL \times L}. \quad (14)$$

We consider the final weighted least-squares problem that requires that each node k communicates with its immediate neighbors within the same cluster [13, 20]:

$$\mathbf{w}_{k,n} = \arg \min_{\mathbf{w} \in \mathbb{R}^L} \left\{ \|\psi_n - \mathbf{H} \mathbf{w}\|_{\Pi_k}^2 \right\} \quad (15)$$

with the node-dependent weighting block diagonal matrix $\Pi_k = \text{bdiag}\{a_{1k} \mathbf{I}_L, \dots, a_{Kk} \mathbf{I}_L\}$. Here, the non-negative combination coefficients $\{a_{\ell k}\}$ are chosen to satisfy [5, 6]:

$$\sum_{\ell \in \mathcal{N}_k \cap \mathcal{C}(k)} a_{\ell k} = 1, \text{ and } \begin{cases} a_{\ell k} > 0, & \text{if } \ell \in \mathcal{N}_k \cap \mathcal{C}(k), \\ a_{\ell k} = 0, & \text{otherwise.} \end{cases} \quad (16)$$

This means that matrix \mathbf{A} with (ℓ, k) -th entry $a_{\ell k}$ is a left-stochastic matrix, i.e., $\mathbf{A}^\top \mathbf{1}_K = \mathbf{1}_K$. Likewise, setting the derivative of (15) with respect to \mathbf{w} to zero, the combination step is given by:

$$\mathbf{w}_{k,n} = \sum_{\ell \in \mathcal{N}_k \cap \mathcal{C}(k)} a_{\ell k} \psi_{\ell,n}. \quad (17)$$

Therefore, the clustered MT-DRLS algorithm with ATC diffusion strategy is consisting of (12) and (17).

3. TRANSIENT PERFORMANCE ANALYSIS

We now perform the transient performance analysis of MT-DRLS algorithm for clustered multitask networks in the mean and mean-square error sense. The weight error vectors for node k at instant n are defined respectively as follows:

$$\tilde{\psi}_{k,n} \triangleq \psi_{k,n} - \mathbf{w}_k^*, \quad \tilde{\mathbf{w}}_{k,n} \triangleq \mathbf{w}_{k,n} - \mathbf{w}_k^*. \quad (18)$$

Let $\tilde{\mathbf{w}}_n$ and \mathbf{w}^* denote the block weight error vector and the block optimal weight vector, all of size $K \times 1$ with blocks of size $L \times 1$, i.e.,

$$\tilde{\mathbf{w}}_n \triangleq \text{col}\{\tilde{\mathbf{w}}_{1,n}, \dots, \tilde{\mathbf{w}}_{K,n}\}, \quad (19)$$

$$\mathbf{w}^* \triangleq \text{col}\{\mathbf{w}_1^*, \dots, \mathbf{w}_K^*\}. \quad (20)$$

We also introduce the following required $K \times K$ block diagonal matrices with each block of size $L \times L$ defined as:

$$\mathbf{R}_{x,n} \triangleq \text{bdiag}\{\mathbf{x}_{1,n} \mathbf{x}_{1,n}^\top, \dots, \mathbf{x}_{K,n} \mathbf{x}_{K,n}^\top\}, \quad (21)$$

$$\Phi_n \triangleq \text{bdiag}\{\Phi_{1,n}, \dots, \Phi_{K,n}\}, \quad (22)$$

$$\mathbf{P}_n \triangleq \text{bdiag}\{\mathbf{P}_{1,n}, \dots, \mathbf{P}_{K,n}\}, \quad (23)$$

$$\mathcal{A} \triangleq \mathbf{A}^\top \otimes \mathbf{I}_L, \quad (24)$$

and the block column vector with individual entries of size $L \times 1$ defined as:

$$\mathbf{s}_{xz,n} \triangleq \text{col}\{z_{1,n}\mathbf{x}_{1,n}, \dots, z_{K,n}\mathbf{x}_{K,n}\}. \quad (25)$$

Moreover, it holds that:

$$\mathbb{E}\{\mathbf{s}_{xz,n}\} = \mathbf{0}_{KL} \quad (26)$$

due to the statistical properties of measurement noise $z_{k,n}$. Before proceeding, we introduce the following independence assumption.

Assumption 1. (Independent Regressors): The regression vectors $\mathbf{x}_{k,n}$ arise from a stationary random process that is temporally stationary, temporally white, and spatially independent with the positive-definite covariance matrix $\mathbf{R}_{x,k}$.

A consequence of Assumption 1 is that $\mathbf{x}_{k,n}$ is independent of $\tilde{\mathbf{w}}_{\ell,m}$ for all ℓ and $m \leq n$. Although not true in general, this assumption is widely used in the theoretical analysis of adaptive filters because it allows to simplify the derivations without constraining the conclusions [21, 22].

3.1. Mean Error Behavior Analysis

With (21) and (22), (6) can be written in the extended form as:

$$\Phi_n = \lambda \Phi_{n-1} + \mathbf{R}_{x,n}. \quad (27)$$

Taking the expectation of both sides, we obtain:

$$\mathbb{E}\{\Phi_n\} = \lambda \mathbb{E}\{\Phi_{n-1}\} + \mathbf{R}_x \quad (28)$$

where the expectation of input correlation matrix $\mathbf{R}_{x,n}$ is given by:

$$\mathbf{R}_x = \mathbb{E}\{\mathbf{R}_{x,n}\} = \text{bdiag}\{\mathbf{R}_{x,1}, \dots, \mathbf{R}_{x,K}\} \in \mathbb{R}^{KL \times KL}. \quad (29)$$

Since matrix Φ_n (or \mathbf{P}_n) only depends on \mathbf{R}_x , the relation (28) is very useful in the sequel. In view of (1) and (18), the a priori estimation error given in (11) can be rewritten as:

$$e_{k,n} = z_{k,n} - \mathbf{x}_{k,n}^\top \tilde{\mathbf{w}}_{k,n-1}. \quad (30)$$

Subtracting \mathbf{w}_k^* from both sides of (12) and (17), respectively, then using (18) and (30), we find that:

$$\begin{aligned} \tilde{\psi}_{k,n} &= \tilde{\mathbf{w}}_{k,n-1} - \mathbf{P}_{k,n} \mathbf{x}_{k,n} \mathbf{x}_{k,n}^\top \tilde{\mathbf{w}}_{n-1} + \mathbf{P}_{k,n} \mathbf{x}_{k,n} z_{k,n} \\ &\quad + \gamma \sum_{\ell \in \mathcal{N}_k \setminus \mathcal{C}(k)} \frac{\rho_{k\ell} + \rho_{\ell k}}{2} (\mathbf{w}_{\ell,n-1} - \mathbf{w}_{k,n-1}), \end{aligned} \quad (31)$$

$$\tilde{\mathbf{w}}_{k,n} = \sum_{\ell \in \mathcal{N}_k \cap \mathcal{C}k} a_{\ell k} \tilde{\psi}_{\ell,n}. \quad (32)$$

Substituting (31) into (32), and using the above definitions (19)–(21) and (23)–(25), the update equation of block weight error vector can be expressed as follows:

$$\tilde{\mathbf{w}}_n = \mathcal{A}[\tilde{\mathbf{w}}_{n-1} - \mathbf{P}_n \mathbf{R}_{x,n} \tilde{\mathbf{w}}_{n-1} + \mathbf{P}_n \mathbf{s}_{xz,n} - \gamma \mathbf{Q}(\mathbf{w}_{n-1} + \mathbf{w}^*)] \quad (33)$$

where

$$\mathbf{Q} = \frac{1}{2} [\text{diag}\{(\Theta + \Theta^\top) \mathbf{1}_K\} - (\Theta + \Theta^\top)] \otimes \mathbf{I}_L. \quad (34)$$

Pre-multiplying both sides of (33) by $\mathbf{P}_n^{-1} \mathcal{A}^{-1}$, using (27) and the relation $\Phi_n = \mathbf{P}_n^{-1}$ based on the definition $\Phi_{k,n} = \mathbf{P}_{k,n}^{-1}$, yields:

$$\Phi_n \mathcal{A}^{-1} \tilde{\mathbf{w}}_n = (\lambda \Phi_{n-1} - \gamma \Phi_n \mathbf{Q}) \tilde{\mathbf{w}}_{n-1} - \gamma \Phi_n \mathbf{Q} \mathbf{w}^* + \mathbf{s}_{xz,n}. \quad (35)$$

The aim of the above manipulations is to separate \mathbf{P}_n and $\mathbf{R}_{x,n}$ in the second term on the r.h.s. of (33). Taking the expectation of both sides of (35), and utilizing the property (26), we then obtain:

$$\begin{aligned} \mathbb{E}\{\Phi_n \mathcal{A}^{-1} \tilde{\mathbf{w}}_n\} &= \mathbb{E}\{(\lambda \Phi_{n-1} - \gamma \Phi_n \mathbf{Q}) \tilde{\mathbf{w}}_{n-1}\} \\ &\quad - \gamma \mathbb{E}\{\Phi_n\} \mathbf{Q} \mathbf{w}^*. \end{aligned} \quad (36)$$

In order to make the analysis mathematical tractable, we need the following approximations [20, 24]:

$$\mathbb{E}\{\Phi_n \mathcal{A}^{-1} \tilde{\mathbf{w}}_n\} \approx \mathbb{E}\{\Phi_n\} \mathcal{A}^{-1} \mathbb{E}\{\tilde{\mathbf{w}}_n\}, \quad (37)$$

$$\mathbb{E}\{\Phi_n \tilde{\mathbf{w}}_n\} \approx \mathbb{E}\{\Phi_n\} \mathbb{E}\{\tilde{\mathbf{w}}_n\}, \quad (38)$$

$$\mathbb{E}\{\Phi_n \mathbf{Q} \tilde{\mathbf{w}}_{n-1}\} \approx \mathbb{E}\{\Phi_n\} \mathbf{Q} \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\}. \quad (39)$$

The proofs of (37)–(39) are not presented explicitly due to the limited space, but the simulation results are able to validate their effectiveness and rationality later. Substituting the approximations (37)–(39) into (36), it follows that:

$$\begin{aligned} \mathbb{E}\{\Phi_n\} \mathcal{A}^{-1} \mathbb{E}\{\tilde{\mathbf{w}}_n\} &= (\lambda \mathbb{E}\{\Phi_{n-1}\} - \gamma \mathbb{E}\{\Phi_n\} \mathbf{Q}) \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} \\ &\quad - \gamma \mathbb{E}\{\Phi_n\} \mathbf{Q} \mathbf{w}^*. \end{aligned} \quad (40)$$

Pre-multiplying both sides of (40) by $\mathcal{A} \mathbb{E}\{\Phi_n\}^{-1}$, it results that

$$\begin{aligned} \mathbb{E}\{\tilde{\mathbf{w}}_n\} &= \mathcal{A} (\lambda \mathbb{E}\{\Phi_n\}^{-1} \mathbb{E}\{\Phi_{n-1}\} - \gamma \mathbf{Q}) \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} \\ &\quad - \gamma \mathcal{A} \mathbf{Q} \mathbf{w}^* \end{aligned} \quad (41)$$

where relation (28) has been used.

3.2. Mean-Square Error Behavior Analysis

The network transient mean-square deviation (MSD) at time instant n is defined by [1, 2]:

$$\text{MSD}_n = \text{tr}\{\tilde{\mathbf{W}}_n\} / K \quad (42)$$

with the correlation matrix of block weight error vector over network $\tilde{\mathbf{W}}_n = \mathbb{E}\{\tilde{\mathbf{w}}_n \tilde{\mathbf{w}}_n^\top\}$. In order to investigate the mean-square error behavior of MT-DRLS algorithm, our next aim is to determine the update equation of $\tilde{\mathbf{W}}_n$. Post-multiplying (35) by its transpose, and taking the expectation of both sides, leads to:

$$\begin{aligned} \mathbf{T}_0 &= \lambda^2 \mathbf{T}_1 + \gamma^2 (\mathbf{T}_2 + \mathbf{T}_3) - \gamma \lambda (\mathbf{T}_4 + \mathbf{T}_4^\top) - \gamma \lambda (\mathbf{T}_5 + \mathbf{T}_5^\top) \\ &\quad + \gamma^2 (\mathbf{T}_6 + \mathbf{T}_6^\top) + \mathbf{S}_{xz} \end{aligned} \quad (43)$$

where

$$\mathbf{T}_0 = \mathbb{E}\{\Phi_n \mathcal{A}^{-1} \tilde{\mathbf{w}}_n \tilde{\mathbf{w}}_n^\top (\mathcal{A}^{-1})^\top \Phi_n\}, \quad (44)$$

$$\mathbf{T}_1 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \Phi_{n-1}\}, \quad (45)$$

$$\mathbf{T}_2 = \mathbb{E}\{\Phi_n \mathbf{Q} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \mathbf{Q}^\top \Phi_n\}, \quad (46)$$

$$\mathbf{T}_3 = \mathbb{E}\{\Phi_n \mathbf{Q} \mathbf{w}^* (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\}, \quad (47)$$

$$\mathbf{T}_4 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \mathbf{Q}^\top \Phi_n\}, \quad (48)$$

$$\mathbf{T}_5 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\}, \quad (49)$$

$$\mathbf{T}_6 = \mathbb{E}\{\Phi_n \mathbf{Q} \tilde{\mathbf{w}}_{n-1} (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\}, \quad (50)$$

$$\mathbf{S}_{xz} = \mathbb{E}\{\mathbf{s}_{xz,n} \mathbf{s}_{xz,n}^\top\}. \quad (51)$$

For mathematical tractability of analysis, we introduce the following necessary approximations:

$$\begin{aligned} \mathbf{T}_0 &= \mathbb{E}\{\Phi_n \mathcal{A}^{-1} \tilde{\mathbf{w}}_n \tilde{\mathbf{w}}_n^\top (\mathcal{A}^{-1})^\top \Phi_n\} \\ &\approx \mathbb{E}\{\Phi_n\} \mathcal{A}^{-1} \tilde{\mathbf{W}}_n (\mathcal{A}^{-1})^\top \mathbb{E}\{\Phi_n\}, \end{aligned} \quad (52)$$

$$\mathbf{T}_1 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \Phi_{n-1}\} \quad (53)$$

$$\approx \mathbb{E}\{\Phi_{n-1}\} \tilde{\mathbf{W}}_{n-1} \mathbb{E}\{\Phi_{n-1}\}^\top,$$

$$\mathbf{T}_2 = \mathbb{E}\{\Phi_n \mathbf{Q} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \mathbf{Q}^\top \Phi_n\} \quad (54)$$

$$\approx \mathbb{E}\{\Phi_n\} \mathbf{Q} \tilde{\mathbf{W}}_{n-1} \mathbf{Q}^\top \mathbb{E}\{\Phi_n\},$$

$$\mathbf{T}_3 = \mathbb{E}\{\Phi_n \mathbf{Q} \mathbf{w}^* (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\} \quad (55)$$

$$\approx \mathbb{E}\{\Phi_n\} \mathbf{Q} \mathbf{w}^* (\mathbf{w}^*)^\top \mathbf{Q}^\top \mathbb{E}\{\Phi_n\},$$

$$\mathbf{T}_4 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} \tilde{\mathbf{w}}_{n-1}^\top \mathbf{Q}^\top \Phi_n\} \quad (56)$$

$$\approx \mathbb{E}\{\Phi_{n-1}\} \tilde{\mathbf{W}}_{n-1} \mathbf{Q}^\top \mathbb{E}\{\Phi_n\},$$

$$\mathbf{T}_5 = \mathbb{E}\{\Phi_{n-1} \tilde{\mathbf{w}}_{n-1} (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\} \quad (57)$$

$$\approx \mathbb{E}\{\Phi_{n-1}\} \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} (\mathbf{w}^*)^\top \mathbf{Q}^\top \mathbb{E}\{\Phi_n\},$$

$$\mathbf{T}_6 = \mathbb{E}\{\Phi_n \mathbf{Q} \tilde{\mathbf{w}}_{n-1} (\mathbf{w}^*)^\top \mathbf{Q}^\top \Phi_n\} \quad (58)$$

$$\approx \mathbb{E}\{\Phi_n\} \mathbf{Q} \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} (\mathbf{w}^*)^\top \mathbf{Q}^\top \mathbb{E}\{\Phi_n\}.$$

The corresponding proofs of (52)–(58) are omitted due to space constraints. The effectiveness and rationality can be testified by the simulation results in the next section. According to assumption 1 and the statistical property of measurement noise $z_{k,n}$, the matrix \mathbf{S}_{xz} can be determined as follows:

$$\mathbf{S}_{xz} = \text{bdiag}\{\sigma_{z,1}^2 \mathbf{R}_{x,1}, \dots, \sigma_{z,K}^2 \mathbf{R}_{x,K}\} = \Sigma_z \mathbf{R}_x \quad (59)$$

with block diagonal matrix $\Sigma_z = \text{bdiag}\{\sigma_{z,1}^2 \mathbf{I}_L, \dots, \sigma_{z,K}^2 \mathbf{I}_L\}$. Substituting (52)–(59) into (43), then multiplying from the left by $\mathcal{A} \mathbb{E}\{\Phi_n\}^{-1}$ and multiplying from the right by $\mathbb{E}\{\Phi_n\}^{-1} \mathcal{A}^\top$ simultaneously, we finally arrive at the recursion of $\tilde{\mathbf{W}}_n$ as follows:

$$\begin{aligned} \tilde{\mathbf{W}}_n = & \mathcal{A} \left[\lambda^2 \mathbb{E}\{\Phi_n\}^{-1} \mathbb{E}\{\Phi_{n-1}\} \tilde{\mathbf{W}}_{n-1} \mathbb{E}\{\Phi_{n-1}\} \mathbb{E}\{\Phi_n\}^{-1} \right. \\ & + \gamma^2 (\mathbf{Q} \tilde{\mathbf{W}}_{n-1} \mathbf{Q}^\top + \mathbf{Q} \mathbf{w}^* (\mathbf{w}^*)^\top \mathbf{Q}^\top) \\ & - \gamma \lambda (\mathbb{E}\{\Phi_n\}^{-1} \mathbb{E}\{\Phi_{n-1}\} \tilde{\mathbf{W}}_{n-1} \mathbf{Q}^\top \\ & + \mathbf{Q} \tilde{\mathbf{W}}_{n-1}^\top \mathbb{E}\{\Phi_{n-1}\} \mathbb{E}\{\Phi_n\}^{-1}) \\ & - \gamma \lambda (\mathbb{E}\{\Phi_n\}^{-1} \mathbb{E}\{\Phi_{n-1}\} \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} (\mathbf{w}^*)^\top \mathbf{Q}^\top \\ & + \mathbf{Q} \mathbf{w}^* \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\}^\top \mathbb{E}\{\Phi_{n-1}\} \mathbb{E}\{\Phi_n\}^{-1}) \\ & + \gamma^2 (\mathbf{Q} \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\} (\mathbf{w}^*)^\top \mathbf{Q}^\top + \mathbf{Q} \mathbf{w}^* \mathbb{E}\{\tilde{\mathbf{w}}_{n-1}\}^\top \mathbf{Q}^\top) \\ & \left. + \mathbb{E}\{\Phi_n\}^{-1} \Sigma_z \mathbf{R}_x \mathbb{E}\{\Phi_n\}^{-1} \right] \mathcal{A}^\top. \end{aligned} \quad (60)$$

It should be pointed out that the above recursive evaluation needs to employ relation (28) and (41). By (60), we can characterize the transient mean-square errors of the clustered MT-DRLS algorithm.

4. NUMERICAL TESTS

In this section, we provide an illustrative example to show the superior performance of clustered MT-DRLS algorithm, and to validate the obtained transient analytical models. We considered a connected network consisting of 14 nodes grouped into 3 clusters shown in Fig. 1(a). The optimal weight vectors to be estimated in each cluster were $\mathbf{w}_{\mathcal{C}_1}^* = [0.5196, -0.3667]^\top$, $\mathbf{w}_{\mathcal{C}_2}^* = [0.4952, -0.3783]^\top$, and $\mathbf{w}_{\mathcal{C}_3}^* = [0.4951, -0.4079]^\top$, respectively. The regression vectors $\mathbf{x}_{k,n}$ were zero-mean random vectors governed by a Gaussian distribution with covariance matrix $\mathbf{R}_{x,k} = \sigma_{x,k} \mathbf{I}_L$. The measurement noise $z_{k,n}$ was i.i.d. Gaussian with zero-mean and variances $\sigma_{z,k}^2$. The variances $\sigma_{z,k}^2$ and $\sigma_{x,k}^2$ are depicted in Fig. 1

(b), respectively. Each combination coefficient $a_{k\ell}$ was chosen as $|\mathcal{N}_k \cap \mathcal{C}(k)|^{-1}$ for all $\ell \in \mathcal{N}_k \cap \mathcal{C}(k)$, where $|\cdot|$ denotes the cardinality of its argument. The regularization weight $\rho_{k\ell}$ was uniformly chosen as $\rho_{k\ell} = |\mathcal{N}_k \setminus \mathcal{C}(k)|^{-1}$ for $\ell \in \mathcal{N}_k \setminus \mathcal{C}(k)$, and $\rho_{k\ell} = 0$ for any other ℓ . The step-sizes of MT-DLMS algorithm were all set to 0.03. The forgetting factor λ and the initialization parameter δ of clustered MT-DRLS algorithm were set to 0.995 and 0.05, respectively. All the empirical learning curves were obtained by averaging over 200 Monte-Carlo runs.

Fig. 1(c) shows that the clustered MT-DRLS algorithms significantly outperforms the counterpart clustered MT-DLMS algorithms in terms of convergence rate, steady-state errors, and parameter estimation accuracy. As shown in Fig. 1(c), the MT-DRLS algorithm with parameter $\gamma = 0.1$ gains about 3 dB over the MT-DRLS algorithm with parameter $\gamma = 0$ in the steady-state MSD, and needs about 300 iterations less before attaining the steady-state phase of MT-DRLS algorithm with parameter $\gamma = 0$. More importantly, we can also see that the consistent agreement between empirical and theoretical MSD curves validates the accuracy and effectiveness of the transient theoretical analysis. Last, this consistency also validates all the necessary approximations introduced in the analysis.

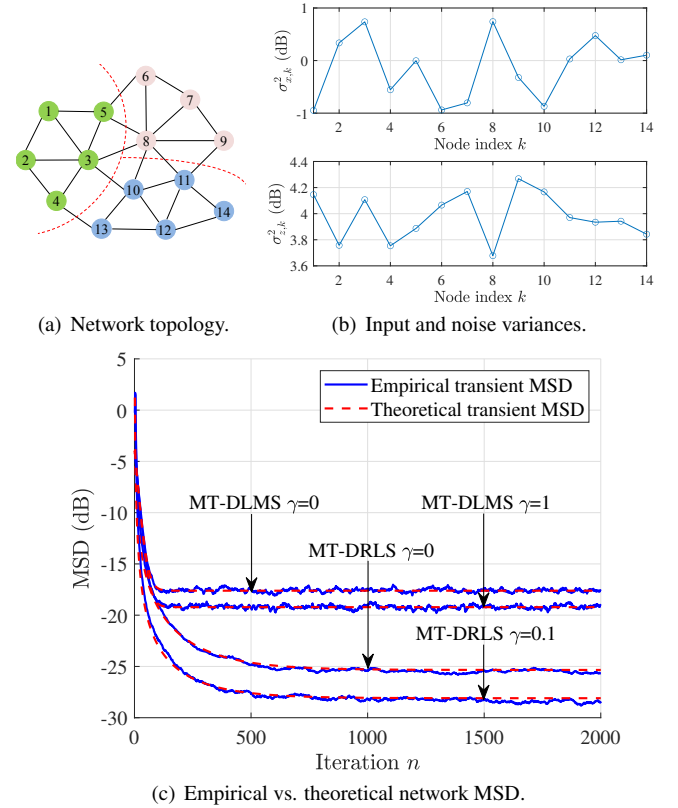


Fig. 1. Network setup and simulation results.

5. CONCLUSION

In this paper, we presented the DRLS algorithm with ATC diffusion strategy over clustered multitask networks to improve the performance of clustered MT-DLMS algorithm. We also provided a transient analysis of the algorithm in the mean and mean-square error sense. In future works, we will study its steady-state behavior.

6. REFERENCES

- [1] A. H. Sayed, "Adaptive networks," *Proc. IEEE*, vol. 102, no. 4, pp. 460–497, Apr. 2014.
- [2] A. H. Sayed, "Diffusion adaptation over networks," in *Academic Press Library in Signal Processing*, vol. 3, pp. 323–453. Elsevier, 2014.
- [3] P. M. Djuric and C. Richard, *Cooperative and Graph Signal Processing: Principles and Applications*, Academic Press, New York, USA, 2018.
- [4] J. Chen, C. Richard, and A. H. Sayed, "Diffusion LMS for clustered multitask networks," in *Proc. IEEE ICASSP*, May 2014, pp. 5487–5491.
- [5] 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.
- [6] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Multitask diffusion adaptation over asynchronous networks," *IEEE Trans. on Signal Process.*, vol. 64, no. 11, pp. 2835–2850, Jun. 2016.
- [7] J. Chen, C. Richard, and A. H. Sayed, "Multitask diffusion adaptation over networks with common latent representations," *IEEE J. Sel. Topics Signal Process.*, vol. 11, no. 3, pp. 563–579, Apr. 2017.
- [8] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Diffusion LMS for multitask problems with local linear equality constraints," *IEEE Trans. Signal Process.*, vol. 65, no. 19, pp. 4979–4993, Oct. 2017.
- [9] R. Nassif, C. Richard, A. Ferrari, and A. H. Sayed, "Proximal multitask learning over networks with sparsity-inducing coregularization," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6329–6344, Dec. 2016.
- [10] D. Jin, J. Chen, C. Richard, and J. Chen, "Online proximal learning over jointly sparse multitask networks with $\ell_{\infty,1}$ regularization," *IEEE Trans. on Signal Process.*, vol. 68, pp. 6319–6335, 2020.
- [11] V. C. Gogineni, S. P. Talebi, and S. Werner, "Performance of clustered multitask diffusion LMS suffering from inter-node communication delays," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 68, no. 7, pp. 2695–2699, Jul. 2021.
- [12] R. Nassif, C. Richard, J. Chen, and A. H. Sayed, "Multitask learning over graphs: An approach for distributed, streaming machine learning," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 14–25, May 2020.
- [13] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion recursive least-squares for distributed estimation over adaptive networks," *IEEE Trans. on Signal Process.*, vol. 56, no. 5, pp. 1865–1877, May 2008.
- [14] A. Bertrand, M. Moonen, and A. H. Sayed, "Diffusion bias-compensated RLS estimation over adaptive networks," *IEEE Trans. on Signal Process.*, vol. 59, no. 11, pp. 5212–5224, Nov. 2011.
- [15] R. Arablouei, K. Dogancay, S. Werner, and Y. Huang, "Adaptive distributed estimation based on recursive least-squares and partial diffusion," *IEEE Trans. on Signal Process.*, vol. 62, no. 14, pp. 3510–3522, Jul. 2014.
- [16] Z. Liu, Y. Liu, and C. Li, "Distributed sparse recursive least-squares over networks," *IEEE Trans. on Signal Process.*, vol. 62, no. 6, pp. 1386–1395, Mar. 2014.
- [17] V. Vahidpour, A. Rastegarnia, A. Khalili, and S. Sane'i, "Analysis of partial diffusion recursive least squares adaptation over noisy links," *IET Signal Processing*, vol. 11, no. 6, pp. 749–757, 2017.
- [18] Y. Yu, H. Zhao, R. C. de Lamare, Y. Zakharov, and L. Lu, "Robust distributed diffusion recursive least squares algorithms with side information for adaptive networks," *IEEE Trans. on Signal Process.*, vol. 67, no. 6, pp. 1566–1581, Mar. 2019.
- [19] A. Rastegarnia, "Reduced-communication diffusion RLS for distributed estimation over multi-agent networks," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 67, no. 1, pp. 177–181, Jan. 2020.
- [20] W. Gao, J. Chen, and C. Richard, "Transient theoretical analysis of diffusion RLS algorithm for cyclostationary colored inputs," *IEEE Signal Process. Lett.*, vol. 28, pp. 1160–1164, 2021.
- [21] S. Haykin, *Adaptive Filter Theory*, Prentice-Hall, New Jersey, 2nd edition, 1991.
- [22] A. H. Sayed, *Fundamentals of Adaptive Filtering*, Wiley, New York, 2003.
- [23] X. Cao and K. J. R. Liu, "Decentralized sparse multitask RLS over networks," *IEEE Trans. on Signal Process.*, vol. 65, no. 23, pp. 6217–6232, Dec. 2017.
- [24] E. Eweda, N. J. Bershad, and J. C. M. Bermudez, "Stochastic analysis of the recursive least squares algorithm for cyclostationary colored inputs," *IEEE Trans. on Signal Process.*, vol. 68, pp. 676–686, 2020.