

ENHANCED DIFFUSION LEARNING OVER NETWORKS

Ricardo Merched*, Stefan Vlaski† and Ali H. Sayed†

*Department of Electronics and Computer Engineering, Universidade Federal do Rio de Janeiro

† Institute of Electrical Engineering, École Polytechnique Fédérale de Lausanne

ABSTRACT

This work develops a variation of diffusion learning by incorporating an adaptive construction for the combination weights through local fusion steps. This leads to an implementation with enhanced convergence rate and mean-square-error performance while maintaining the same level of complexity as standard implementations. The approach is based on formulating optimal or close-to-optimal learning and fusion steps using a proximity function rationale within neighborhoods. The first version of the algorithm employs exact fusion in the least-squares sense using inverses of uncertainty matrices. The second version replaces these matrices by diagonal approximations with reduced complexity. The result is an LMS-complexity scheme with improved performance for distributed learning over networks.

Index Terms— diffusion networks, fusion, least-squares, adaptation, combination weights.

1. INTRODUCTION

In typical implementations of consensus and diffusion strategies for learning over networks, it is customary to combine estimates from neighborhoods by relying on convex combination weights [1]–[6]. In general, these weights are fixed and chosen as (scalar) entries of left-stochastic combination matrices. There have been works in the literature where the combination weights have also been learned as part of the adaptation process. For example, the *relative variance* combination rule from [7] was derived by optimizing the instantaneous mean-square-error measure as the learning algorithm evolves over time.

In this work, we take a different route to learning the combination weights, leading to enhanced performance. We attain this objective by formulating optimal or close-to-optimal fusion steps locally under a proximity rationale, and subsequently reduce the complexity of the iterations by replacing uncertainty matrices by scalar approximations. The main difference between our derivation and existing approaches is that uncertainties in the estimates by the agents are carried from

self-learning to the social learning phase, and back to self-learning again, in a continuous fashion. This results in adaptive combination weights and lead to lower mean-square-error and faster convergence compared to existing approaches.

2. DIFFUSION ALGORITHMS REVISITED

We consider a strongly-connected network of distributed agents, represented by a collection of N nodes in Fig. 1. Each agent k receives streaming data $\{d_k(i), \mathbf{u}_{k,i}\}$, assumed to be related via a linear regression model of the form

$$d_k(i) = \mathbf{u}_{k,i} \mathbf{w}_k^o + v_k(i) \quad (1)$$

where i is the time index, \mathbf{w}_k^o is an unknown local parameter of size $M \times 1$, $\mathbf{u}_{k,i}$ is a regression (row) vector of size $1 \times M$, and $v_k(i)$ is additive zero-mean white noise, which is temporally and spatially uncorrelated with other data.

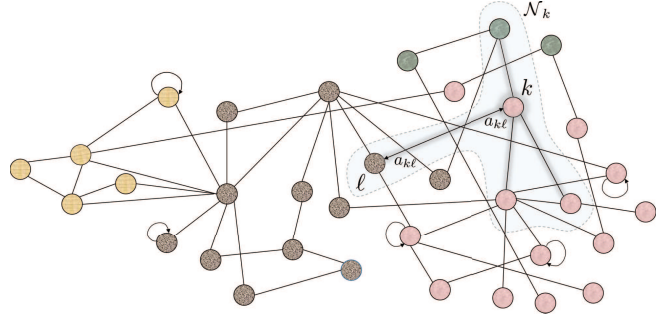


Fig. 1. Illustration of a network with $N = 34$ agents.

We collect the data measured at each agent k up to time i into the quantities:

$$\mathbf{H}_{k,i} = \text{col}\{\mathbf{u}_{k,1}, \mathbf{u}_{k,2}, \dots, \mathbf{u}_{k,i}\} \quad (2)$$

$$\mathbf{y}_{k,i} = \text{col}\{d_k(1), d_k(2), \dots, d_k(i)\} \quad (3)$$

In a non-cooperative setting, each agent k would estimate its parameter vector \mathbf{w}_k^o by solving a weighted least-squares problem of the form:

$$\min_{\mathbf{w}_k} \|\mathbf{y}_{k,i} - \mathbf{H}_{k,i} \mathbf{w}_k\|_{\Lambda_{k,i}}^2 \quad (4)$$

The work of A. H. Sayed was supported in part by NSF grant CCF-1524250. Emails: merched@lps.ufrj.br, {stefan.vlaski, ali.sayed}@epfl.ch.

for some weighting matrix $\mathbf{A}_{k,i} = \text{diag}\{\lambda^{i-1}, \dots, \lambda, 1\}$ where $0 \ll \lambda \leq 1$ is an exponential weighting factor. In many important situations, however, the individual models $\{\mathbf{w}_k^o\}$ are related, such as varying smoothly over a graph. In these cases, it is reasonable to encourage agents to seek estimates that are close to each other. One way to achieve this objective is to replace problem (4) by a regularized problem that enforces coupling among agents. To exploit this idea, let $\text{bdiag}(\cdot)$ denote a block diagonal operator and introduce the following extended quantities:

$$\mathbf{w}^o = \text{col}\{\mathbf{w}_1^o, \mathbf{w}_2^o, \dots, \mathbf{w}_N^o\} \quad (5)$$

$$\mathbf{y}_i = \text{col}\{\mathbf{y}_{1,i}, \mathbf{y}_{2,i}, \dots, \mathbf{y}_{N,i}\} \quad (6)$$

$$\mathbf{v}_i = \text{col}\{\mathbf{v}_{1,i}, \mathbf{v}_{2,i}, \dots, \mathbf{v}_{N,i}\} \quad (7)$$

$$\mathbf{H}_i = \text{bdiag}\{\mathbf{H}_{1,i}, \mathbf{H}_{2,i}, \dots, \mathbf{H}_{N,i}\} \quad (8)$$

$$\mathbf{A}_i = (1/N)\text{bdiag}\{\mathbf{A}_{1,i}, \dots, \mathbf{A}_{N,i}\} \quad (9)$$

Definitions (5)-(9) allow us to write a global linear model for the data collected across the network up to time i as follows:

$$\mathbf{y}_i = \mathbf{H}_i \mathbf{w}^o + \mathbf{v}_i \quad (10)$$

Introduce the extended parameter vector:

$$\mathbf{w} = \text{col}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\} \quad (11)$$

Now, one way to enforce coupling is via a proximity function. This is usually achieved by considering the following regularized cost

$$\min_{\mathbf{w}} \|\mathbf{y}_i - \mathbf{H}_i \mathbf{w}\|_{\mathbf{A}_i}^2 + \rho_i(\mathbf{w}) \quad (12)$$

for some regularizer $\rho_i(\mathbf{w})$ whose purpose is to encourage proximity to \mathbf{w}^o . However, since \mathbf{w}^o is unknown, we shall instead encourage proximity with respect to the best guess available at that moment. We will explain how to obtain this estimate in the sequel. For now, let us denote its entries by $\bar{w}_{k,i}$ at agent k ; i.e., this is the estimate for w_k^o at agent k at time i . We shall also explain how to associate with this estimate a local uncertainty matrix denoted by $\bar{\mathbf{P}}_{k,i}$: the better the quality of $\bar{w}_{k,i}$, the smaller $\bar{\mathbf{P}}_{k,i}$ will be so that the inverse matrix $\bar{\mathbf{P}}_{k,i}^{-1}$ serves as a measure of the uncertainty in the estimate $\bar{w}_{k,i}$.

Using these intermediate estimates (to be constructed later), we replace the cost in (12) by one of the form

$$J_i(\mathbf{w}) = \|\mathbf{y}_i - \mathbf{H}_i \mathbf{w}\|_{\mathbf{A}_i}^2 + \sum_{m=1}^i \lambda^{i-m} \|\mathbf{w} - \bar{\mathbf{w}}_m\|_{\bar{\mathbf{P}}_m}^2 \quad (13)$$

where $\bar{\mathbf{w}}_m = \text{col}\{\bar{w}_{1,m}, \bar{w}_{2,m}, \dots, \bar{w}_{N,m}\}$, and

$$\bar{\mathbf{P}}_m^{-1} = \text{bdiag}\{\bar{\mathbf{P}}_{1,m}^{-1}, \bar{\mathbf{P}}_{2,m}^{-1}, \dots, \bar{\mathbf{P}}_{N,m}^{-1}\} \quad (14)$$

One possible choice for $\bar{w}_{k,i}$ is $w_{k,i-1}$, the estimate obtained at time $i-1$. Observe that (13) pushes \mathbf{w} towards the estimates $\bar{\mathbf{w}}_m$ by computing an exponentially weighted error

measure from time $m=1$ up to time $m=i$. Since this objective function is quadratic, it can be expressed in terms of its minimizer and Hessian matrix as

$$J_i(\mathbf{w}) = \|\mathbf{w} - \mathbf{w}_i\|_{\bar{\mathbf{P}}_i}^2 + c \quad (15)$$

for some constant c , which can be ignored, and where

$$\mathbf{w}_i = \arg \min_{\mathbf{w}} \|\mathbf{w} - \mathbf{w}_i\|_{\bar{\mathbf{P}}_i}^2 \quad (16)$$

$$\bar{\mathbf{P}}_i^{-1} = \nabla^2 J(\mathbf{w}) = \mathbf{H}_i^* \mathbf{A}_i \mathbf{H}_i + \sum_{m=1}^i \lambda^{i-m} \bar{\mathbf{P}}_m^{-1} \quad (17)$$

with $*$ denoting complex conjugate transposition. Of course, this reformulation is not useful in finding the minimizer of (13) since it involves \mathbf{w}_i itself. It does however allow us to obtain a recursive algorithm for the solution. To see this, let $\{\mathbf{d}_i, \mathbf{U}_i\}$ contain the most recent data at time i :

$$\mathbf{d}_i = \text{col}\{d_1(i), d_2(i), \dots, d_N(i)\} \quad (18)$$

$$\mathbf{U}_i = \text{bdiag}\{\mathbf{u}_{1,i}, \mathbf{u}_{2,i}, \dots, \mathbf{u}_{N,i}\} \quad (19)$$

Then, (13) can be expressed as

$$\begin{aligned} J_i(\mathbf{w}) &= \|\mathbf{d}_i - \mathbf{U}_i \mathbf{w}\|^2 + \|\mathbf{y}_{i-1} - \mathbf{H}_{i-1} \mathbf{w}\|_{\lambda \mathbf{A}_{i-1}}^2 \\ &\quad + \lambda \sum_{m=1}^{i-1} \lambda^{i-m-1} \|\mathbf{w} - \bar{\mathbf{w}}_m\|_{\bar{\mathbf{P}}_m}^2 + \|\mathbf{w} - \bar{\mathbf{w}}_i\|_{\bar{\mathbf{P}}_i}^2 \\ &= \|\mathbf{d}_i - \mathbf{U}_i \mathbf{w}\|^2 + \lambda \underbrace{\|\mathbf{w} - \mathbf{w}_{i-1}\|_{\bar{\mathbf{P}}_{i-1}}^2}_{J_{i-1}(\mathbf{w})} + \|\mathbf{w} - \bar{\mathbf{w}}_i\|_{\bar{\mathbf{P}}_i}^2 \end{aligned} \quad (20)$$

Note that $J_i(\mathbf{w})$ is composed of three components. The first term, $\|\mathbf{d}_i - \mathbf{U}_i \mathbf{w}\|^2$, fits \mathbf{w} to the most recent data; the second term, $\lambda \|\mathbf{w} - \mathbf{w}_{i-1}\|_{\bar{\mathbf{P}}_{i-1}}^2$, incorporates past information; and the last term $\|\mathbf{w} - \bar{\mathbf{w}}_i\|_{\bar{\mathbf{P}}_i}^2$ promotes smoothness and closeness to the intermediate estimate $\bar{\mathbf{w}}_i$ yet to be specified.

Using completion of squares in (20), we can combine the first two terms and write the minimization problem as:

$$\min_{\mathbf{w}} \|\mathbf{w} - \hat{\mathbf{w}}_i\|_{\hat{\mathbf{P}}_i}^2 + \|\mathbf{w} - \bar{\mathbf{w}}_i\|_{\bar{\mathbf{P}}_i}^2 \quad (21)$$

in terms of the updated local estimates:

$$\hat{\mathbf{w}}_i = \mathbf{w}_{i-1} + \hat{\mathbf{P}}_i \mathbf{U}_i^* (\mathbf{d}_i - \mathbf{U}_i \mathbf{w}_{i-1}) \quad (22)$$

$$\hat{\mathbf{P}}_i^{-1} = \lambda \bar{\mathbf{P}}_{i-1}^{-1} + \mathbf{U}_i^* \mathbf{U}_i \quad (23)$$

where $\hat{\mathbf{w}}_i = \text{col}\{\hat{w}_{1,i}, \hat{w}_{2,i}, \dots, \hat{w}_{N,i}\}$, and

$$\hat{\mathbf{P}}_{i-1}^{-1} = \text{bdiag}\{\hat{\mathbf{P}}_{1,i-1}^{-1}, \hat{\mathbf{P}}_{2,i-1}^{-1}, \dots, \hat{\mathbf{P}}_{N,i-1}^{-1}\} \quad (24)$$

We are now left to define $\{\bar{w}_{k,i}, \bar{\mathbf{P}}_{k,i}^{-1}\}$. While $\hat{w}_{k,i}$ is an improved estimate for w_k^o over $w_{k,i-1}$, since it includes the most recent data, we proceed one step further by allowing for an exchange of information within neighborhoods. This can

be obtained by designing $\bar{\mathbf{w}}_i$ as a fusion in the neighborhood \mathcal{N}_k of the agent's estimates $\{\hat{\mathbf{w}}_{k,i}, \hat{\mathbf{P}}_{k,i}\}$, defined in (22) and (23), in a weighted least-squares manner, however, without including agent k . This is because the contribution of agent k will be added by the first term of (21). Now, since (21) is equivalent to minimizing (15), instead of designing $\bar{\mathbf{w}}_i$, we can design the result of their fusion, say, $\{\mathbf{w}_i, \mathcal{P}_i\}$, directly as the solution to

$$\min_{\mathbf{w}'_k} \sum_{\ell \in \mathcal{N}_k} a_{k\ell} \|\mathbf{w}'_k - \hat{\mathbf{w}}_{\ell,i}\|_{\hat{\mathbf{P}}_{\ell,i}^{-1}}^2 \quad (25)$$

$$= \min_{\mathbf{w}'_k} \|\mathbf{w}'_k - \mathbf{w}'_{k,i}\|_{\mathbf{P}'_{k,i}^{-1}}, \quad k = 1, 2, \dots, N \quad (26)$$

$$= \min_{\mathbf{w}'} \|\mathbf{w}' - \mathcal{W}'_i\|_{\mathcal{P}'_i^{-1}} \quad (27)$$

for positive scalars $\{a_{k\ell}\}$. The solution to (25) is given by

$$\mathbf{w}'_{k,i} = \sum_{\ell \in \mathcal{N}_k} \mathbf{A}_{k\ell,i} \hat{\mathbf{w}}_{\ell,i}, \quad \mathbf{A}_{k\ell,i} \triangleq a_{k\ell} \mathbf{P}'_{k,i} \hat{\mathbf{P}}_{\ell,i}^{-1} \quad (28)$$

if $\ell \in \mathcal{N}_k$, while $\mathbf{A}_{k\ell,i} = \mathbf{0}$ if $\ell \notin \mathcal{N}_k$. The quantity $\mathbf{P}'_{k,i}^{-1}$ is the uncertainty that results from these estimates,

$$\mathbf{P}'_{k,i}^{-1} = \sum_{\ell \in \mathcal{N}_k} a_{k\ell} \hat{\mathbf{P}}_{\ell,i}^{-1} \quad (29)$$

Hence, in extended vector form, this yields $\{\mathcal{W}'_i, \mathcal{P}'_i\}$. Now, observe that if we select $\mathbf{w}_i = \mathcal{W}'_i$ and $\mathcal{P}_i = \mathcal{P}'_i$, the costs (27) and (15) will have the same form. In other words, by selecting (27) and (15) to have the same minimizer with the same uncertainty, we are able to propagate these quantities from self-learning to social learning in a true recursion.

Note that while we could assume for simplicity that all agents in \mathcal{N}_k are equally important in the fusion process, say, $a_{k\ell} = \gamma_k \geq 0$, node k itself can be assigned a different weight a_{kk} , relative to its neighbors, implying that it can have more or less certainty of its own estimate. In order for these coefficients to add up to one, we select

$$a_{kk} = 1 - \gamma_k(n_k - 1) \quad (30)$$

where n_k is the degree of the neighborhood. Moreover, by associating the scalars $a_{k\ell}$ to entries of a $N \times N$ matrix \mathbf{A} , then \mathbf{A} is referred to a Laplacian matrix.

Finally, defining $\mathcal{A} = \mathbf{A} \otimes \mathbf{I}_M$, where \otimes denotes the Kronecker product, then, in extended matrix notation, it holds that

$$\mathbf{w}_i = \mathbf{A}_i \hat{\mathbf{w}}_i \quad (31)$$

where, using (24),

$$\mathbf{A}_i = \mathcal{P}_i \mathcal{A} \hat{\mathcal{P}}_i^{-1}, \quad \mathcal{P}_i^{-1} = \text{bdiag} \left[\mathcal{A} \hat{\mathcal{P}}_i^{-1} (\mathbf{1} \otimes \mathbf{I}_M) \right] \quad (32)$$

It is easily verified that $\mathbf{A}_i \mathbf{1} = (\mathbf{1} \otimes \mathbf{I}_M)$ is block right-stochastic by construction. The diffusion recursions derived

so far are listed in Table 1, and constitute what we shall refer to as the *Adapt-and-Fuse* (AAF) diffusion, which extends the usual description of diffusion strategy known as *Adapt-then-Combine* (ATC), where $\mathbf{A}_i = \mathcal{A}$. Figure 2 illustrates the equivalent global transmission scheme, where we have defined $\mathcal{G}_i = \hat{\mathcal{P}}_i \mathbf{U}_i^*$.

Initialization: $\mathbf{w}_0 = \mathbf{0}$, $\mathcal{P}_0^{-1} = \epsilon \mathbf{I}$ for small ϵ

$$\hat{\mathcal{P}}_i^{-1} = \lambda \mathcal{P}_{i-1}^{-1} + \mathbf{U}_i^* \mathbf{U}_i$$

$$\mathcal{P}_i^{-1} = \text{bdiag} \left[\mathcal{A} \hat{\mathcal{P}}_i^{-1} (\mathbf{1} \otimes \mathbf{I}_M) \right]$$

$$\mathbf{A}_i = \mathcal{P}_i \mathcal{A} \hat{\mathcal{P}}_i^{-1}$$

$$\hat{\mathbf{w}}_i = \mathcal{W}_{i-1} + \hat{\mathcal{P}}_i \mathbf{U}_i^* (d_i - \mathbf{U}_i \mathcal{W}_{i-1})$$

$$\mathbf{w}_i = \mathbf{A}_i \hat{\mathbf{w}}_i$$

Table 1. AAF Diffusion Adaptation.

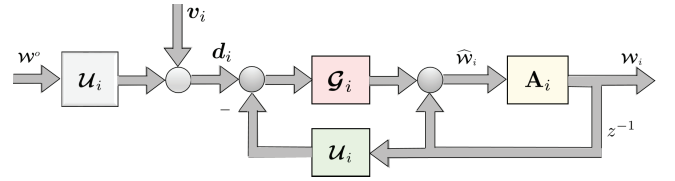


Fig. 2. Global Description of the AAF scheme.

3. SIMPLIFIED AAF RECURSIONS

Because the optimal coefficients of \mathbf{A}_i and all its defining covariances w.r.t. node k are complex and matrix-valued, they require complex matrix \times matrix operations and inversions at each time i . These computations can be cumbersome considering that agents should perform elementary operations when fusing their estimates within \mathcal{N}_k . By assuming uncorrelated input regressors, we can simplify the AAF recursions by restricting the covariances to diagonal matrices, i.e., we set $\mathbf{P}_{k,i} \approx \sigma_k^2(i) \mathbf{I}$, and $\hat{\mathbf{P}}_{k,i} \approx \hat{\sigma}_k^2(i) \mathbf{I}$. With these approximations, the $N \times N$ block entry of \mathbf{A}_i in (32) simplifies to

$$[\mathbf{A}_i]_{k\ell} = a_{k\ell}(i) \mathbf{I} = a_{k\ell} \sigma_k^2(i) \hat{\sigma}_\ell^{-2}(i) \mathbf{I}$$

The resulting algorithm is shown in Table 2, which we refer to as the Simplified-AAF (SAAF) algorithm.

Initialization: $\mathbf{w}_{\ell,0} = \mathbf{0}$, $\sigma_k^{-2}(0) = \epsilon$ for small ϵ

for $k = 1$ to N :

$$\hat{\sigma}_k^{-2}(i) = \lambda \sigma_k^{-2}(i-1) + |\mathbf{u}_k(i)|^2$$

$$\sigma_k^{-2}(i) = \sum_{\ell \in \mathcal{N}_k} a_{k\ell} \hat{\sigma}_\ell^{-2}(i)$$

$$a_{\ell k}(i) = a_{k\ell} \sigma_k^2(i) \hat{\sigma}_\ell^{-2}(i)$$

$$\hat{\mathbf{w}}_{k,i} = \mathbf{w}_{k,i-1} + \hat{\sigma}_k^2(i) \mathbf{u}_{k,i}^* [d_k(i) - \mathbf{u}_{k,i} \mathbf{w}_{k,i-1}]$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k}(i) \hat{\mathbf{w}}_{\ell,i}$$

Table 2. SAAF Algorithm.

Note that the computational complexity for updating the combination matrix in the SAAF recursions is of $\mathcal{O}(NM)$.

This is in contrast to the relative-variance rule of [7] which requires $\mathcal{O}(NNM)$ computations per iteration in order to accomplish a similar task, where $\mathcal{N} = \sum_{k=1}^N \mathcal{N}_k$. That is, the uncertainties of each node in [7] are assumed to vary over the neighborhoods, while in the proposed construction, the nodes uncertainties are absolute; they are updated independently, and combined through the fixed Laplacian matrix \mathbf{A} .

4. SIMULATIONS

In order to illustrate the performance of the proposed AAF based algorithms in comparison with existing cooperative-based schemes, we consider a topology with $N = 20$ agents with unknown vectors of size $M = 10$, and compare: (i) the power normalized LMS-based algorithms employing a fixed combination policy; (ii) a normalized LMS version of the adaptive relative variance diffusion algorithm of [7], and (iii) the RLS-based diffusion of, e.g., [11].

In order to set one possible theoretical benchmark for comparison on the minimum mean-square-deviation (MSD), we consider the one corresponding to the well established diffusion LMS algorithm. The network MSD in this case, for sufficiently small μ is given by (see [8], pp. 606)

$$\text{MSD}_{\text{dist,av}} = \frac{M}{2} \left(\sum_{k=1}^N \mu_k^2 p_k^2 \sigma_{u,k}^2 (\sigma_{v,k}^2 + \sigma_{u,k}^2 \|\mathbf{w}^o - \mathbf{w}_k^o\|^2) \right) \cdot \left(\sum_{k=1}^N \mu_k p_k \sigma_{u,k}^2 \right)^{-1} \quad (33)$$

in terms of the input and noise variances, $\sigma_{u,k}^2$ and $\sigma_{v,k}^2$ respectively, and the entries p_k of the Perron vector associated to the combination matrix, here chosen as the optimal relative-variance rule (for the cases when $\mathbf{w}_k^o = \mathbf{w}^o$). This is illustrated by a thick straight line.

◆ **Scenario 1 (Performance of the SAAF algorithm)**: Figure 3 shows typical ensemble average learning curves for uncorrelated inputs. We set the step-size as $\mu = 0.002$ for the LMS-based algorithms, $\lambda = 1$ for the RLS algorithms, and $\gamma_k = 0.0024$ in (30) for the proposed recursions.

We see that the proposed simplified algorithm outperforms LMS-based algorithms, and with reduced complexity compared to the relative-variance policy. We clearly see a difference in terms of the MSD attained. Moreover, it exhibits approximately the same performance of existing RLS-based recursions. It is worth noting that for the latter, the use of a Metropolis combination rule \mathbf{C} does not yield improvement against the case when $\mathbf{C} = \mathbf{I}$.

◆ **Scenario 2 (Performance of the full AAF algorithm for colored input and $\lambda < 1$ for all RLS algorithms)**: Figure 4 shows the curves for the RLS algorithms when $\lambda = 0.99$, and for a slightly colored AR process, with pole at 0.5.

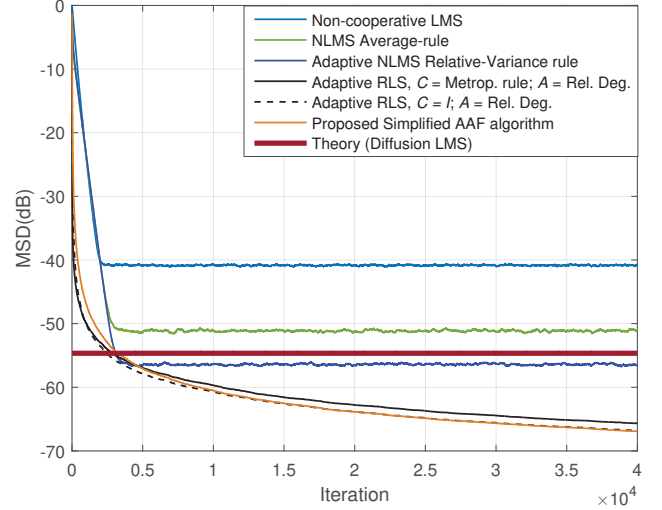


Fig. 3. Comparison among algorithms, uncorrelated input.

We verify that existing diffusion RLS algorithms have their performance degraded. The exact AAF-RLS, despite the computational complexity, outperforms in speed and MSE level, which continues to decrease beyond $8 \cdot 10^4$ iterations.

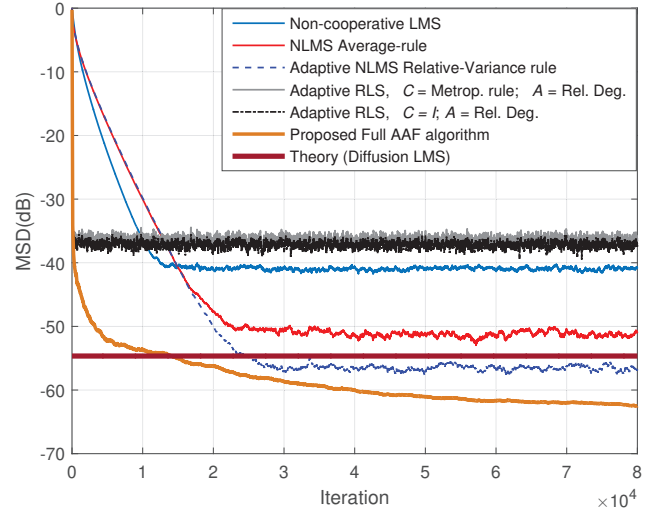


Fig. 4. AAF performance with colored input and for $\lambda < 1$.

5. CONCLUSIONS

We have proposed a construction for optimized diffusion networks which fuses estimates and uncertainties at every node in the LS sense. The proposed AAF recursions outperform existing diffusion algorithms, especially when $\lambda < 1$, and for colored inputs. For uncorrelated data, we verified that its simplified version outperforms all other algorithms, and exhibiting the same convergence performance of existing RLS diffusion schemes, however, under LMS complexity.

6. REFERENCES

- [1] A. H. Sayed, "Adaptation, learning, and optimization over networks," *Foundations and Trends in Machine Learning*, vol. 7, no. 4-5, pp. 311-801, July 2014.
- [2] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460-497, April 2014.
- [3] J. Tsitsiklis and M. Athans, "Convergence and asymptotic agreement in distributed decision problems," *IEEE Trans. Autom. Control*, vol. 29, no. 1, pp. 42-50, Jan. 1984.
- [4] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multiagent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48- 61, Jan. 2009.
- [5] A. G. Dimakis, S. Kar, J. M. F. Moura, M. G. Rabbat, and A. Scaglione, "Gossip algorithms for distributed signal processing," *Proc. IEEE*, vol. 98, no. 11, pp. 1847-1864, Nov. 2010.
- [6] S. Kar and J. M. F. Moura, "Convergence rate analysis of distributed gossip (linear parameter) estimation: Fundamental limits and tradeoffs," *IEEE Journal on Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 674-690, Aug. 2011.
- [7] S-Y Tu and A. H. Sayed, "Optimal combination rules for adaptation and learning over networks," *Proc. IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, San Juan, Puerto Rico, pp. 317-320, December 2011.
- [8] A. H. Sayed, *Adaptation, Learning, and Optimization over Networks*, *Foundations and Trends in Machine Learning*, vol. 7, issue 4-5, NOW Publishers, Boston-Delft, 518pp, 2014. ISBN 978-1-60198-850-8, DOI 10.1561/22000000051.
- [9] J. Chen and A. H. Sayed, "On the limiting behavior of distributed optimization strategies," *Proc. 50th Annual Allerton Conference on Communication, Control, and Computing*, pp. 1535-1542, Monticello, IL, October 2012.
- [10] R. G. P. Pinto and R. Merched "A compressed sensing approach to block-iterative equalizers", *IEEE Trans. on Signal Processing*, vol. 66, no. 4, pp. 1007-1022, Dec. 2017.
- [11] A. H. Sayed and C. Lopes, "Distributed recursive least-squares strategies over adaptive networks," *Proc. 40th Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, CA, pp. 233-237, October-November, 2006.
- [12] H. Salami, B. Ying, and A. H. Sayed, "Diffusion social learning over weakly-connected graphs," *Proc. IEEE ICASSP*, pp. 4119-4123, Shanghai, China, March 2016.
- [13] S. Vlaski, L. Vandenberghe, and A. H. Sayed, "Diffusion stochastic optimization with non-smooth regularizers," *Proc. IEEE ICASSP*, pp. 4149-4153, Shanghai, China, March 2016.
- [14] K. Yuan, B. Ying, X. Zhao, and A. H. Sayed, "Exact diffusion for distributed optimization and learning – Part I: Algorithm development," *IEEE Trans. Signal Processing*, vol. 67, no. 3, pp. 708-723, Feb. 2019.
- [15] R. Nassif, S. Vlaski, and A. H. Sayed, "Learning over multitask graphs — Part I: Stability analysis," available as arXiv:1805.08535v1, May 2018.