



**HAL**  
open science

## Non-negative least-mean-square algorithm

Jie Chen, Cédric Richard, José C. M. Bermudez, Paul Honeine

► **To cite this version:**

Jie Chen, Cédric Richard, José C. M. Bermudez, Paul Honeine. Non-negative least-mean-square algorithm. *IEEE Transactions on Signal Processing*, 2011, 59 (11), pp.5225 - 5235. 10.1109/TSP.2011.2162508 . hal-01965584

**HAL Id: hal-01965584**

**<https://hal.science/hal-01965584v1>**

Submitted on 3 Jan 2019

**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.

# Non-negative least-mean-square algorithm

Jie Chen<sup>(1,2)</sup>, Cédric Richard<sup>(1)</sup>, *Senior Member, IEEE*,

Jose Carlos M. Bermudez<sup>(3)</sup>, *Senior Member, IEEE*, Paul Honeine<sup>(2)</sup>, *Member, IEEE*

<sup>(1)</sup> Université de Nice Sophia-Antipolis, UMR CNRS 6525, Observatoire de la Côte d'Azur

Laboratoire Hippolyte Fizeau, Parc Valrose, 06102 Nice cedex 2 - France

tel.: +33.4.92.07.63.94      fax.: +33.4.92.07.63.21

chenjieg@sina.com      cedric.richard@unice.fr

<sup>(2)</sup> Université de Technologie de Troyes, UMR CNRS 6279

Laboratoire LM2S, BP 2060, 10010 Troyes cedex - France

tel.: +33.3.25.71.56.25      fax.: +33.3.25.71.56.99

paul.honeine@utt.fr

<sup>(3)</sup> Federal University of Santa Catarina

Department of Electrical Engineering, 88040-900, Florianópolis, SC - Brazil

tel.: +55.48.3721.7719      fax.: +55.48.3721.9280

j.bermudez@ieee.org

## Abstract

Dynamic system modeling plays a crucial role in the development of techniques for stationary and non-stationary signal processing. Due to the inherent physical characteristics of systems under investigation, non-negativity is a desired constraint that can usually be imposed on the parameters to estimate. In this paper, we propose a general method for system identification under non-negativity constraints. We derive the so-called *non-negative least-mean-square algorithm* based on stochastic gradient descent, and we analyze its convergence. Experiments are conducted to illustrate the performance of this approach and consistency with the analysis.

## I. INTRODUCTION

In many real-life phenomena including biological and physiological ones, due to the inherent physical characteristics of systems under investigation, non-negativity is a desired constraint that can be imposed on the parameters to estimate in order to avoid physically absurd and uninterpretable results. For instance, in the study of a concentration field or a thermal radiation field, any observation is described with non-negative values (ppm, joule). Non-negativity as a physical constraint has received growing attention from the signal processing community during the last decade. For instance, consider the following non-negative least-square inverse problem

$$\begin{aligned} \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 \\ \text{subject to } [\mathbf{x}]_i \geq 0, \quad \forall i, \end{aligned} \quad (1)$$

with  $\mathbf{A}$  a real  $M \times N$  matrix of rank  $k \leq \min(M, N)$ ,  $\mathbf{b}$  an  $M$ -length real vector, and  $\mathbf{x}$  an  $N$ -length real vector.  $\|\cdot\|$  denotes the Euclidean 2-norm and  $[\cdot]_i$  the  $i$ -th entry of the vector. This problem has been addressed in various contexts, with applications ranging from image deblurring in astrophysics [1] to deconvolution of emission spectra in chemometrics [2]. Another similar problem is the non-negative matrix factorization (NMF), which is now a popular dimension reduction technique [3], [4], [5]. Given a matrix  $\mathbf{X}$  with non-negative entries, the squared error version of this problem can be stated as follows

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 \\ \text{subject to } [\mathbf{W}]_{ij} \geq 0, [\mathbf{H}]_{ij} \geq 0, \quad \forall i, j \end{aligned} \quad (2)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. This problem is closely related to the blind deconvolution one, and has found direct application in hyperspectral imaging [6]. Separation of non-negative mixture of non-negative sources has also been considered in [7], [8].

Over the last fifteen years, a variety of methods have been developed to tackle non-negative least-square problems (NNLS). Active set techniques for NNLS use the fact that if the set of variables which activate constraints is known, then the solution of the constrained least-square problem can be obtained by solving an unconstrained one that only includes inactive variables. The active set algorithm of Lawson and Hanson [9] is a batch resolution technique for NNLS problems. It has become a standard among the most frequently used methods. In [10], Bro and De Jong introduced a modification of the latter, called fast NNLS, which takes advantage of the special characteristics of iterative algorithms involving repeated use of non-negativity constraints. Another class of tools is the class of projected gradient algorithms [11], [12], [13], [14]. They are based on successive projections on the feasible region. In [15], Lin used this kind of algorithms for NMF

problems. Low memory requirements and simplicity make algorithms in this class attractive for large scale problems. Nevertheless, they are characterized by slow convergence rate if not combined with appropriate step size selection. The class of multiplicative algorithms is very popular for dealing with NMF problems [4], [16]. Particularly efficient updates were derived in this way for a large number of problems involving non-negativity constraints [17]. These algorithms however require batch processing, which is not suitable for online system identification problems.

In this paper, we consider the problem of system identification under non-negativity constraints on the parameters to estimate. The Karush-Kuhn-Tucker (KKT) conditions are established for any convex cost function, and a fixed-point iteration strategy is then applied in order to derive a gradient descent algorithm. Considering the square-error criterion as a particular case, a stochastic gradient scheme is presented. A convergence analysis of this algorithm is proposed. The resulting model accurately predicts the algorithm behavior for both transient and steady-state conditions. Finally, experiments are conducted to evaluate the algorithm performance and its consistency with the analysis.

## II. SYSTEM IDENTIFICATION WITH NON-NEGATIVITY CONSTRAINTS

Consider an unknown system, only characterized by a set of real-valued discrete-time responses to known stationary inputs. The problem is to derive a transversal filter model

$$y(n) = \boldsymbol{\alpha}^\top \mathbf{x}(n) + z_1(n), \quad (3)$$

with  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N]^\top$  the vector of the model parameters, and  $\mathbf{x}(n) = [x(n), x(n-1), \dots, x(n-N+1)]^\top$  the observed data vector. The input signal  $x(n)$  and the desired output signal  $y(n)$  are assumed stationary and zero-mean. The sequence  $z_1(n)$  accounts for measurement noise and modeling errors.

Due to the inherent physical characteristics of systems under investigation, in this paper, non-negativity is a desired constraint that is imposed on the coefficient vector  $\boldsymbol{\alpha}$ . Therefore, the problem of identifying the optimum model can be formalized as follows

$$\begin{aligned} \boldsymbol{\alpha}^o &= \arg \min_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}) \\ &\text{subject to } \alpha_i \geq 0, \quad \forall i, \end{aligned} \quad (4)$$

with  $J(\boldsymbol{\alpha})$  a continuously differentiable and strictly convex cost function in  $\mathbb{R}^N$ , and  $\boldsymbol{\alpha}^o$  the optimal solution to the constrained optimization problem.

### A. A fixed-point iteration scheme

In order to solve the problem (4), let us consider its Lagrangian function  $Q(\boldsymbol{\alpha}, \boldsymbol{\lambda})$  given by [18]

$$Q(\boldsymbol{\alpha}, \boldsymbol{\lambda}) = J(\boldsymbol{\alpha}) - \boldsymbol{\lambda}^\top \boldsymbol{\alpha},$$

where  $\boldsymbol{\lambda}$  is the vector of non-negative Lagrange multipliers. The Karush-Kuhn-Tucker conditions must necessarily be satisfied at the optimum defined by  $\boldsymbol{\alpha}^o$ ,  $\boldsymbol{\lambda}^o$ , namely,

$$\begin{aligned} \nabla_{\boldsymbol{\alpha}} Q(\boldsymbol{\alpha}^o, \boldsymbol{\lambda}^o) &= 0 \\ \alpha_i^o [\boldsymbol{\lambda}^o]_i &= 0, \quad \forall i \end{aligned}$$

where the symbol  $\nabla_{\boldsymbol{\alpha}}$  stands for the gradient operator with respect to  $\boldsymbol{\alpha}$ . Using  $\nabla_{\boldsymbol{\alpha}} Q(\boldsymbol{\alpha}, \boldsymbol{\lambda}) = \nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}) - \boldsymbol{\lambda}$ , these equations can be combined into the following expression

$$\alpha_i^o [-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}^o)]_i = 0, \quad (5)$$

where the extra minus sign is just used to make a gradient descent of  $J(\boldsymbol{\alpha})$  apparent. To solve Equation (5) iteratively, two important points have to be noticed. The first point is that  $\mathbf{D}(-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}))$  is also a gradient descent of  $J(\boldsymbol{\alpha})$  if  $\mathbf{D}$  is a symmetric positive definite matrix. The second point is that equations of the form  $\varphi(u) = 0$  can be solved with a fixed-point iteration algorithm, under some conditions on function  $\varphi$ , by considering the problem  $u = u + \varphi(u)$ . Implementing this strategy with Equation (5) leads us to the component-wise gradient descent algorithm

$$\alpha_i(n+1) = \alpha_i(n) + \eta_i(n) f_i(\boldsymbol{\alpha}(n)) \alpha_i(n) [-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}(n))]_i \quad (6)$$

with  $\eta_i(n)$  a positive step size required to get a contraction scheme and to control the convergence rate. Function  $f_i(\boldsymbol{\alpha}) > 0$  in (6) is the  $i$ -th entry of a diagonal matrix  $\mathbf{D}$ . It is an arbitrary positive function of  $\boldsymbol{\alpha}$ . Some criteria  $J(\boldsymbol{\alpha})$  are defined only for inputs  $\boldsymbol{\alpha}$  with positive entries, e.g., Itakura-Saito distance, Kullback-Leibler divergence. If necessary, this condition can be managed by an appropriate choice of the step size parameter. Let us assume that  $\alpha_i(n) \geq 0$ . Non-negativity of  $\alpha_i(n+1)$  is guaranteed if

$$1 + \eta_i(n) f_i(\boldsymbol{\alpha}(n)) [-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}(n))]_i \geq 0. \quad (7)$$

If  $[-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}(n))]_i \leq 0$ , condition (7) is clearly satisfied and non-negativity does not impose any restriction on the step size. Conversely, if  $[-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}(n))]_i > 0$ , non-negativity of  $\alpha_i(n+1)$  holds if

$$0 \leq \eta_i(n) \leq \frac{1}{f_i(\boldsymbol{\alpha}(n)) [-\nabla_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}(n))]_i}. \quad (8)$$

Using a single step size  $\eta(n)$  in  $[0, \eta_{\max}(n)]$  for all entries of  $\boldsymbol{\alpha}$  so that

$$\eta_{\max}(n) = \min_i \frac{1}{f_i(\boldsymbol{\alpha}(n)) [\nabla_{\alpha} J(\boldsymbol{\alpha}(n))]_i}, \quad i = 1, \dots, N \quad (9)$$

the update equation can be written in vector form as

$$\boldsymbol{\alpha}(n+1) = \boldsymbol{\alpha}(n) + \eta(n) \mathbf{d}(n), \quad (10)$$

where the weight adjustment direction  $\mathbf{d}(n)$ , whose  $i$ -th entry is defined as follows

$$[\mathbf{d}(n)]_i = f_i(\boldsymbol{\alpha}(n)) \alpha_i(n) [-\nabla_{\alpha} J(\boldsymbol{\alpha}(n))]_i \quad (11)$$

is a gradient descent direction because  $f_i[\boldsymbol{\alpha}(n)] \alpha_i(n) \geq 0$ . It should be noted that condition (9) on the step size  $\eta(n)$  guarantees the non-negativity of  $\boldsymbol{\alpha}(n)$  for all  $n$ , but does not ensure the stability of the algorithm.

### B. The non-negative least-mean-square algorithm

Let us now consider the mean-square error criterion  $J_{mse}(\boldsymbol{\alpha})$  to be minimized with respect to  $\boldsymbol{\alpha}$ , that is,

$$\begin{aligned} \boldsymbol{\alpha}^o &= \arg \min_{\boldsymbol{\alpha}} E\{[y(n) - \boldsymbol{\alpha}^T \mathbf{x}(n)]^2\} \\ &\text{subject to } \alpha_i^o \geq 0, \quad \forall i, \end{aligned} \quad (12)$$

where we have included the non-negativity constraint only on the optimum solution because  $J_{mse}(\boldsymbol{\alpha})$  is defined for all  $\boldsymbol{\alpha}$ , that is, for all positive and negative entries  $\alpha_i$ . The gradient of  $J_{mse}(\boldsymbol{\alpha})$  can be easily computed as

$$\nabla_{\alpha} J(\boldsymbol{\alpha}) = 2(\mathbf{R}_x \boldsymbol{\alpha} - \mathbf{r}_{xy}) \quad (13)$$

with  $\mathbf{R}_x$  the autocorrelation matrix of  $\mathbf{x}(n)$  and  $\mathbf{r}_{xy}$  the correlation vector between  $\mathbf{x}(n)$  and  $y(n)$ . Using (10) and (11) with  $f_i(\boldsymbol{\alpha}) = \frac{1}{2}$  for all  $i$ , the update rule for minimizing the mean-square error under non-negativity constraints is given by

$$\boldsymbol{\alpha}(n+1) = \boldsymbol{\alpha}(n) + \eta(n) \mathbf{D}_{\alpha}(n) (\mathbf{r}_{xy} - \mathbf{R}_x \boldsymbol{\alpha}(n)) \quad (14)$$

where  $\mathbf{D}_{\alpha}(n)$  is the diagonal matrix with diagonal entries given by  $\boldsymbol{\alpha}(n)$ . Following a stochastic gradient approach, the second-order moments  $\mathbf{R}_x$  and  $\mathbf{r}_{xy}$  are replaced in (14) by the instantaneous estimates  $\mathbf{x}(n) \mathbf{x}^T(n)$  and  $y(n) \mathbf{x}(n)$ , respectively. This leads to the stochastic approximation of (14) given by<sup>1</sup>

$$\boldsymbol{\alpha}(n+1) = \boldsymbol{\alpha}(n) + \eta(n) e(n) \mathbf{D}_x(n) \boldsymbol{\alpha}(n), \quad \eta(n) > 0 \quad (15)$$

where  $\mathbf{D}_x(n)$  stands for the diagonal matrix with diagonal entries given by  $\mathbf{x}(n)$ , and  $e(n) = y(n) - \boldsymbol{\alpha}^T(n) \mathbf{x}(n)$ .

<sup>1</sup>Note that  $\mathbf{D}_{\alpha}(n) \mathbf{x}(n) = \mathbf{D}_x(n) \boldsymbol{\alpha}(n)$ .

It is interesting to notice how the term  $\alpha(n)$  in the update term on the right-hand side (r.h.s.) of (15) affects the dynamics of the coefficient update when compared with the well know LMS algorithm [19]. Note that the extra multiplying factor  $\alpha_i(n)$  in the update term of the  $i$ -th row of (15), which is not present in the LMS update, provides extra control of both the magnitude and the direction of the weight update, as compared to LMS. For a fixed step size  $\eta$ , the update term for the  $i$ -th component of  $\alpha(n)$  is proportional to  $-\alpha_i(n)e(n)x_i(n)$ , the stochastic gradient component. Thus, compared to the LMS stochastic gradient component  $-e(n)x_i(n)$ , the constrained algorithm includes the multiplying factor  $\alpha_i(n)$ . A negative  $\alpha_i(n)$  will then change the sign of the LMS adjustment, which on average tends to avoid convergence to negative coefficients of the unconstrained solution. Thus, coefficients that would normally converge, on average, to negative values using unconstrained LMS will tend to converge to zero using the constrained algorithm. In addition,  $\alpha_i(n)$  close to zero will tend to slow its own convergence unless the magnitude of  $e(n)x_i(n)$  is very large. Finally,  $\alpha_i(n) = 0$  is clearly a stationary point of Equation (15).

In the following, the adaptive weight behavior of the adaptive algorithm (15), called *non-negative LMS*, is studied in the mean and mean-square senses for a time-invariant step size  $\eta$ .

### III. MEAN BEHAVIOR ANALYSIS

We shall now propose a model to characterize the mean behavior of the non-negative LMS algorithm. Figure 1 shows a block diagram of the problem studied in this paper. The input signal  $x(n)$  and the desired output signal  $y(n)$  are assumed stationary and zero-mean. Let us denote by  $\alpha^*$  the solution of the unconstrained least-mean-square problem

$$\alpha^* = \arg \min_{\alpha} E\{[y(n) - \alpha^\top \mathbf{x}(n)]^2\}. \quad (16)$$

whose solution  $\alpha^*$  satisfies the Wiener-Hopf equations

$$\mathbf{R}_x \alpha^* = \mathbf{r}_{xy}. \quad (17)$$

The residual signal  $z(n) = y(n) - \mathbf{x}^\top(n) \alpha^*$  in Figure 1 accounts for measurement noise and modeling errors. It is assumed in the following that  $z(n)$  is stationary, zero-mean with variance  $\sigma_z^2$  and statistically independent of any other signal. Thus,  $E\{z(n) \mathbf{D}_x(n)\} = 0$ .

The adaptive algorithm (15) attempts to estimate the optimum  $\alpha^o$  for the constrained problem (12). The analytical determination of the optimal solution  $\alpha^o$  is not trivial in general. In the particular case of independent and identically distributed (i.i.d.) input samples, however,  $\mathbf{R}_x = \sigma_x^2 \mathbf{I}$  where  $\mathbf{I}$  is the identity matrix. In this case, the Karush-Kuhn-Tucker conditions imply that  $\alpha^o$  is obtained by turning the negative entries of  $\alpha^*$  to zero

$$\alpha^o = \{\alpha^*\}_+ \quad (18)$$

where  $\{u\}_+ = \max\{0, u\}$ . The minimum mean-square error produced by solution  $\boldsymbol{\alpha}^o$  is then

$$J_{msmin} = \sigma_y^2 - 2 \mathbf{r}_{xy} \{\boldsymbol{\alpha}^*\}_+ + \sigma_x^2 \{\boldsymbol{\alpha}^*\}_+^\top \{\boldsymbol{\alpha}^*\}_+ \quad (19)$$

with  $\sigma_y^2$  the variance of  $y(n)$ .

#### A. Mean weight behavior model

Defining the weight-error vector  $\mathbf{v}(n) = \boldsymbol{\alpha}(n) - \boldsymbol{\alpha}^* = [v_1(n), v_2(n), \dots, v_N(n)]^\top$ , the update equation (15) can be written as

$$\mathbf{v}(n+1) = \mathbf{v}(n) + \eta e(n) \mathbf{D}_x(n) (\mathbf{v}(n) + \boldsymbol{\alpha}^*). \quad (20)$$

Using  $e(n) = y(n) - \boldsymbol{\alpha}^\top(n) \mathbf{x}(n) = z(n) - \mathbf{v}^\top(n) \mathbf{x}(n)$  leads us to the following expression

$$\begin{aligned} \mathbf{v}(n+1) = & \mathbf{v}(n) + \eta z(n) \mathbf{D}_x(n) \mathbf{v}(n) + \eta z(n) \mathbf{D}_x(n) \boldsymbol{\alpha}^* \\ & - \eta \mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) - \eta \mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n) \mathbf{v}(n). \end{aligned} \quad (21)$$

Taking the expectation of expression (21), neglecting the statistical dependence of  $\mathbf{x}(n)$  and  $\mathbf{v}(n)$ ,<sup>2</sup> and using that  $E\{z(n) \mathbf{D}_x(n)\} = 0$  yields

$$E\{\mathbf{v}(n+1)\} \approx (\mathbf{I} - \eta E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n)\}) E\{\mathbf{v}(n)\} - \eta E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n)\}. \quad (22)$$

The first expectation on the r.h.s. of (22) is given by

$$E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n)\} = E\{\mathbf{D}_{\alpha^*} \mathbf{x}(n) \mathbf{x}^\top(n)\} = \mathbf{D}_{\alpha^*} \mathbf{R}_x. \quad (23)$$

<sup>2</sup>This assumption is less restrictive than the well-known independence assumption [19, p. 247], as it does not require  $x(n)$  be Gaussian.

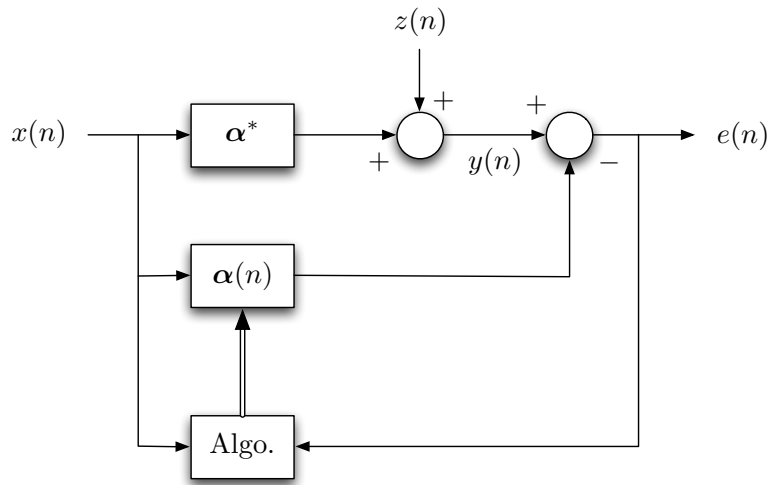


Fig. 1. Adaptive system under study



In order to evaluate the second expectation on the r.h.s. of (22), let us compute the  $i$ -th component of the vector  $\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n)$ . We have

$$[\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n)]_i = \sum_{j=1}^N x(n-i+1) v_i(n) v_j(n) x(n-j+1) \quad (24)$$

Taking the expectation of this expression, defining  $\mathbf{K}(n) = E\{\mathbf{v}(n) \mathbf{v}^\top(n)\}$ , and neglecting the statistical dependence of  $\mathbf{x}(n)$  and  $\mathbf{v}(n)$ , we obtain

$$\begin{aligned} [E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n)\}]_i &\approx \sum_{j=1}^N r_x(j-i) [\mathbf{K}(n)]_{ij} \\ &= [\mathbf{R}_x \mathbf{K}(n)]_{ii} \end{aligned} \quad (25)$$

This implies that  $E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n)\} \approx \text{diag}\{\mathbf{R}_x \mathbf{K}(n)\}$ , where  $\text{diag}\{\mathbf{A}\}$  denotes the vector whose  $i$ -th entry is defined by  $[\mathbf{A}]_{ii}$ . Using these results with Equation (22) yields the following mean weight-error vector update equation

$$E\{\mathbf{v}(n+1)\} = (\mathbf{I} - \eta \mathbf{D}_{\alpha^*} \mathbf{R}_x) E\{\mathbf{v}(n)\} - \eta \text{diag}\{\mathbf{R}_x \mathbf{K}(n)\}. \quad (26)$$

This equation requires second-order moments defined by  $\mathbf{K}(n)$  in order to update the first-order one  $E\{\mathbf{v}(n)\}$ . A recursive model will be derived for  $\mathbf{K}(n)$  in Section IV, see Equation (39). That model can be used along with (26) to predict the mean weight behavior of the algorithm. Nevertheless, we have found that a sufficiently accurate and more intuitive mean behavior model can be obtained using the following separation approximation

$$\mathbf{K}(n) \approx E\{\mathbf{v}(n)\} E\{\mathbf{v}^\top(n)\}. \quad (27)$$

Using (27) in (26) we obtain the following result

$$E\{\mathbf{v}(n+1)\} = (\mathbf{I} - \eta \mathbf{D}_{\alpha^*} \mathbf{R}_x) E\{\mathbf{v}(n)\} - \eta \text{diag}\{\mathbf{R}_x E\{\mathbf{v}(n)\} E\{\mathbf{v}^\top(n)\}\}. \quad (28)$$

Approximation (27) assumes that

$$\text{Cov}\{v_i(n), v_j(n)\} \ll E\{v_i(n)\} E\{v_j(n)\} \quad (29)$$

In general, (29) is valid when the adaptive weights are far from convergence, as the mean weight-error component tends to be much larger than the weight-error fluctuation about the mean. For correlated  $x(n)$ , the level of the weight-error fluctuations at convergence tends to be much less than the values of the nonzero optimal weights. For white input signals  $E\{v_i(n)\}$  tends to zero for those indexes corresponding to the positive coefficients of  $\alpha^o$ . In this case, approximation (29) will in fact tend to eliminate the weight estimation error at convergence. Extensive simulation results have shown that the simplified model in (28) yields a prediction of the mean weight behavior which is sufficient for design purposes. Furthermore, this simplification allows the more detailed analytical study of the mean weight behavior shown in the next section.

### B. Special case of a white input signal

In general, the behavior of (28) can become very complex to be studied analytically [20]. In order to obtain analytical results that allow some understanding of the mean weight behavior, we study here the particular case of  $x(n)$  i.i.d. and drawn from a zero-mean distribution. Unit variance  $\sigma_x^2$  is also assumed without loss of generality. Using  $\mathbf{R}_x = \mathbf{I}$  in (28) yields the component-wise expression

$$E\{v_i(n+1)\} = (1 - \eta \alpha_i^*) E\{v_i(n)\} - \eta E\{v_i(n)\}^2. \quad (30)$$

Function  $E\{v_i(n+1)\}$  in Equation (30) is a parabola in the variable  $E\{v_i(n)\}$  with roots at  $E\{v_i(n)\} = 0$  and  $E\{v_i(n)\} = \frac{1-\eta\alpha_i^*}{\eta}$ . It reaches its maximum value  $\frac{(1-\eta\alpha_i^*)^2}{4\eta}$  at  $E\{v_i(n)\} = \frac{1-\eta\alpha_i^*}{2\eta}$ . Fixed points are found by solving  $E\{v_i(n+1)\} = E\{v_i(n)\}$ , which yields  $E\{v_i(n)\} = 0$  or  $E\{v_i(n)\} = -\alpha_i^*$ . This result is consistent with solution (18) where

$$v_i^o = \begin{cases} 0 & \text{if } \alpha_i^* \geq 0 \\ -\alpha_i^* & \text{otherwise} \end{cases} \quad (31)$$

with  $v_i^o$  the  $i$ -th entry of  $\mathbf{v}^o = \boldsymbol{\alpha}^o - \boldsymbol{\alpha}^*$ .

Let us now derive conditions ensuring convergence of (30) to 0 and  $-\alpha_i^*$ . Writing  $u(n) = \frac{\eta}{1-\eta\alpha_i^*} E\{v_i(n)\}$ , where the index  $i$  has been dropped to simplify the notation, we obtain the following difference equation known as the *logistic map* [20], [21], [22]

$$u(n+1) = \rho u(n) (1 - u(n)) \quad (32)$$

with  $\rho = 1 - \eta\alpha_i^*$ , which is assumed nonzero. Fixed points defined in (31) now correspond to  $u = 0$  and  $u = \frac{\rho-1}{\rho}$ , respectively. Convergence of the logistic map to these values depends on the parameter  $\rho$  and on the initial condition  $u(0)$  as follows. See [20], [21], [22] for details and Figure 2 for illustration.

#### Case 1: $0 < \rho < 1$

An illustration of this case is shown in Figure 2 (left). The fixed point  $u = 0$  attracts all the trajectories originating in the interval  $]\frac{\rho-1}{\rho}; \frac{1}{\rho}[$ . The logistic map  $u(n)$  is identically equal to  $\frac{\rho-1}{\rho}$  for  $n \geq 1$  if  $u(0) = \frac{\rho-1}{\rho}$  or  $u(0) = \frac{1}{\rho}$ . Outside the interval,  $u(n)$  diverges to  $-\infty$ .

#### Case 2: $\rho = 1$

The fixed point  $u = 0$  attracts all the trajectories originating in the interval  $[0; 1]$ . The logistic map  $u(n)$  is identically equal to 0 for  $n \geq 1$  if  $u(0) = 0$  or 1. It diverges to  $-\infty$  if  $u(0) \notin [0; 1]$ .

#### Case 3: $1 < \rho \leq 3$

An illustration of this case is shown in Figure 2 (right). The fixed point  $u = \frac{\rho-1}{\rho}$  attracts all the trajectories

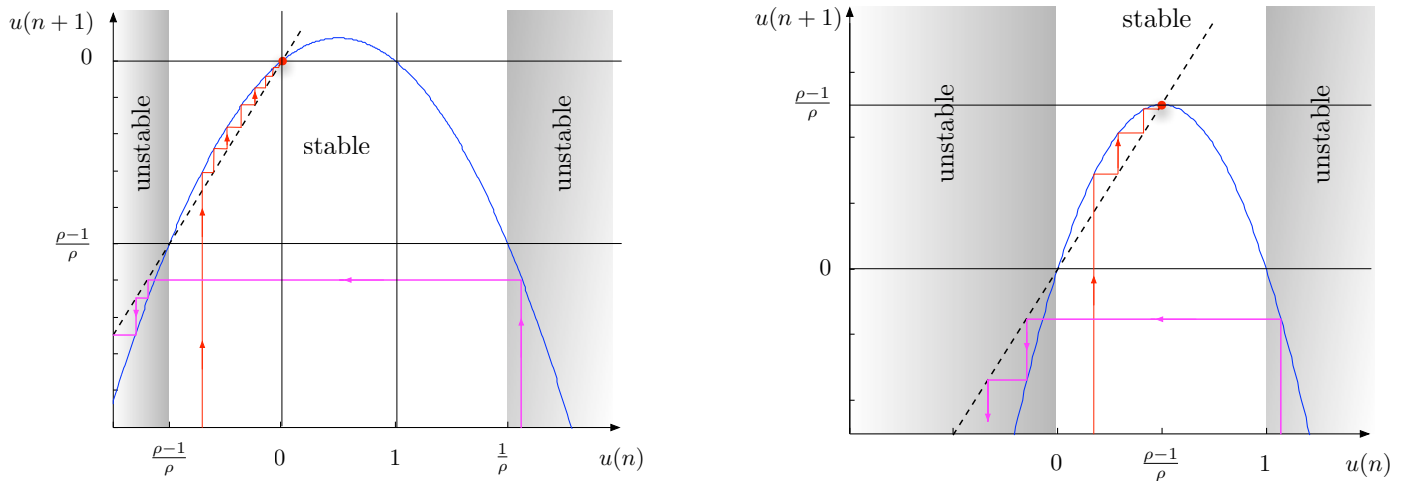


Fig. 2. Convergence of the logistic map, in the Case 1 (left) and in the Case 3 (right). The dashed line is the line of equation  $u(n+1) = u(n)$ .

originating in  $]0; 1[$ . With the initial conditions  $u(0) = 0$  or  $u(0) = 1$ , we have  $u(n) = 0$  for all  $n > 0$ . It can be shown that the logistic map diverges to  $-\infty$  if  $u(0) \notin [0; 1]$ .

#### Case 4: $\rho > 3$

Fixed points become unstable. New fixed points appear between which the system alternates in stable cycles of period  $2^k$ , with  $k$  tending to infinity as  $\rho$  increases. This case may even lead to a chaotic behavior, and falls out of the scope of our study.

To derive conditions for convergence of the difference equation (30) to 0 or  $-\alpha_i^*$ , we must consider separately components of  $E\{v_i(n)\}$  associated with positive or negative unconstrained optimum  $\alpha_i^*$ , respectively. On the one hand, based on the analysis of the logistic map (32), convergence of (30) to 0 corresponds to the conditions on  $\rho$  and  $u(0)$  satisfying Case 1 and Case 2 above. This yields

$$0 < \eta < \frac{1}{\alpha_i^*} \quad -\alpha_i^* < v_i(0) < \frac{1}{\eta} \quad (33)$$

in the case where  $\alpha_i^* > 0$ . If  $\alpha_i^* = 0$ , these two conditions become  $\eta > 0$  and  $0 < v_i(0) < \frac{1}{\eta}$ . On the other hand,  $\rho$  and  $u(0)$  must obey the conditions presented in Case 3 for convergence of Equation (30) to  $-\alpha_i^*$ . This leads to

$$0 < \eta \leq -\frac{2}{\alpha_i^*} \quad 0 < v_i(0) < \frac{1}{\eta} - \alpha_i^* \quad (34)$$

in the case where  $\alpha_i^* < 0$ . Finally, combining these inequalities leads to the following theoretical conditions for convergence of  $E\{v(n)\}$ :

$$0 < \eta \leq \min_i \frac{1}{|\alpha_i^*|} \quad \text{and} \quad 0 < v_i(0) < \frac{1}{\eta} \quad \text{for all } i \quad (35)$$

or, using also Equations (33) and (34), for convergence of  $E\{\alpha(n)\}$ :

$$0 < \eta \leq \min_i \frac{1}{|\alpha_i^*|} \quad \text{and} \quad 0 < \alpha_i(0) < \frac{1}{\eta} \quad \text{for all } i. \quad (36)$$

Conditions (35) and (36) on  $v_i(0)$  and  $\alpha_i(0)$  show that there is more freedom in choosing  $\alpha_i(0)$  for small values of  $\eta$ . They guarantee the convergence of the difference equation (30).

### C. Simulation examples for the first-order moment analysis

This section presents simulation examples to verify the validity of the first-order moment analysis of the non-negative LMS algorithm. We illustrate the accuracy of the model (30) through a first example where inputs  $x(n)$  and noise  $z(n)$  are i.i.d. and drawn from a zero-mean Gaussian distribution with variance  $\sigma_x^2 = 1$  and  $\sigma_z^2 = 10^{-2}$ , respectively. The impulse response  $\alpha^*$  is given by

$$\alpha^* = [0.8 \quad 0.6 \quad 0.5 \quad 0.4 \quad 0.3 \quad 0.2 \quad 0.1 \quad -0.1 \quad -0.3 \quad -0.6]^\top \quad (37)$$

The initial impulse response  $\alpha(0)$  is drawn from the uniform distribution  $\mathcal{U}([0; 1])$ , and kept unchanged for all the simulations. The algorithm's stability limit was determined experimentally to be  $\eta_{\max} \approx 5 \times 10^{-3}$ . As usually happens with adaptive algorithms, this limit is more restrictive than the mean weight convergence limit given by (36), as stability is determined by the weight fluctuations [19]. The mean value  $E\{\alpha_i(n)\}$  of each coefficient is shown in Figure 3 for  $\eta = 10^{-3} = \eta_{\max}/5$  and  $\eta = 5 \times 10^{-5} = \eta_{\max}/10$ . The simulation curves (blue line) were obtained from Monte Carlo simulation averaged over 100 realizations. The theoretical curves (red line) were obtained from model (30). One can notice that all the curves are perfectly superimposed and, as predicted by the result (18), each coefficient  $\alpha_i(n)$  tends to  $\{\alpha_i^*\}_+$  as  $n$  goes to infinity.

It is interesting to note that convergence towards the solution  $\{\alpha^*\}_+$  agrees with the theoretically predicted behavior of (32). For each positive entry  $\alpha_i^*$  of  $\alpha^*$ , the corresponding value of  $\rho = 1 - \eta \alpha_i^*$  is in  $]0; 1[$ . This corresponds to Case 1 in Section III-B, where the fixed point  $u = 0$  attracts all the trajectories and  $v_i(n)$  converges to zero. It can also be verified that each  $\rho$  associated with a negative entry  $\alpha_i^*$  is in  $]1; 3]$ . This corresponds to Case 3 where  $u = (\rho - 1)/\rho$  attracts all the trajectories and  $\lim_{n \rightarrow \infty} v_i(n) = -\alpha_i^*$ .

The second simulation example illustrates the accuracy of the model (30) for inputs  $x(n)$  correlated in time. We consider a first-order AR model given by  $x(n) = a x(n-1) + w(n)$ , with  $a = \frac{1}{2}$ . The noise  $w(n)$  is i.i.d. and drawn from a zero-mean Gaussian distribution with variance  $\sigma_w^2 = 1 - \frac{1}{4}$ , so that  $\sigma_x^2 = 1$  as in the first example. The other parameters of the initial experimental setup remain unchanged, except for the step size values. In order to verify the model's accuracy also for large step sizes we performed the simulations for  $\eta = 2.5 \times 10^{-3} = \eta_{\max}/2$  and  $\eta = 5 \times 10^{-5} = \eta_{\max}/10$ . The mean value  $E\{\alpha_i(n)\}$  of each coefficient is shown in Figure 4. As before, the simulation curves (blue line) and the theoretical curves (red line) are superimposed.

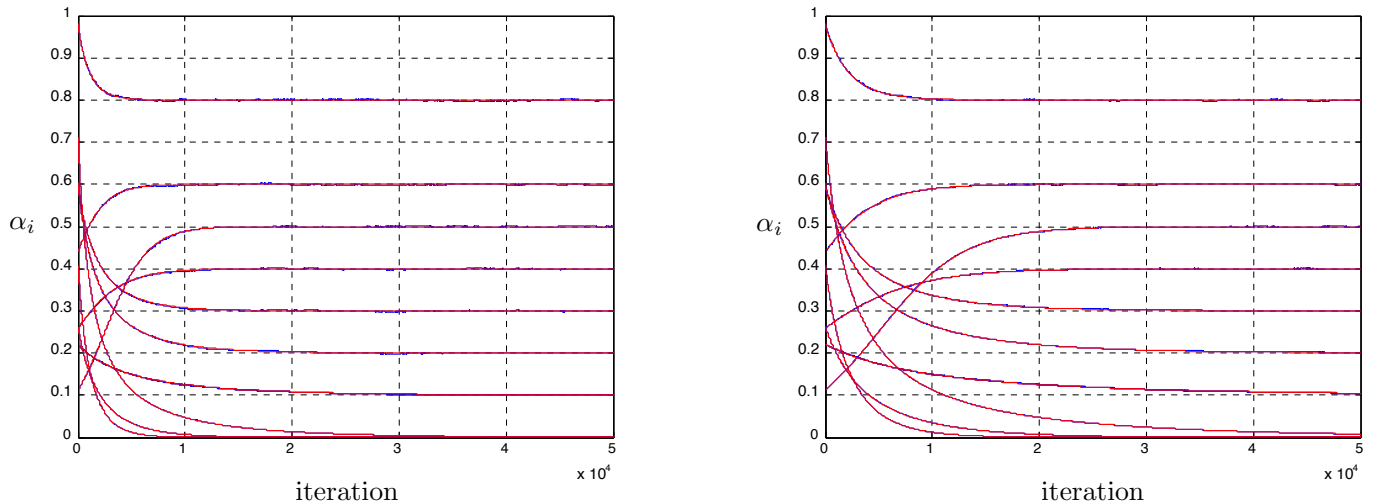


Fig. 3. Convergence of the coefficients  $\alpha_i(n)$  in the case where input  $x(n)$  and noise  $z(n)$  are i.i.d. Two different step sizes are considered:  $\eta = 10^{-3}$  on the left figure, and  $\eta = 5 \times 10^{-4}$  on the right figure. The theoretical curves (red line) obtained from the model (30) and simulation curves (blue line) are perfectly superimposed.

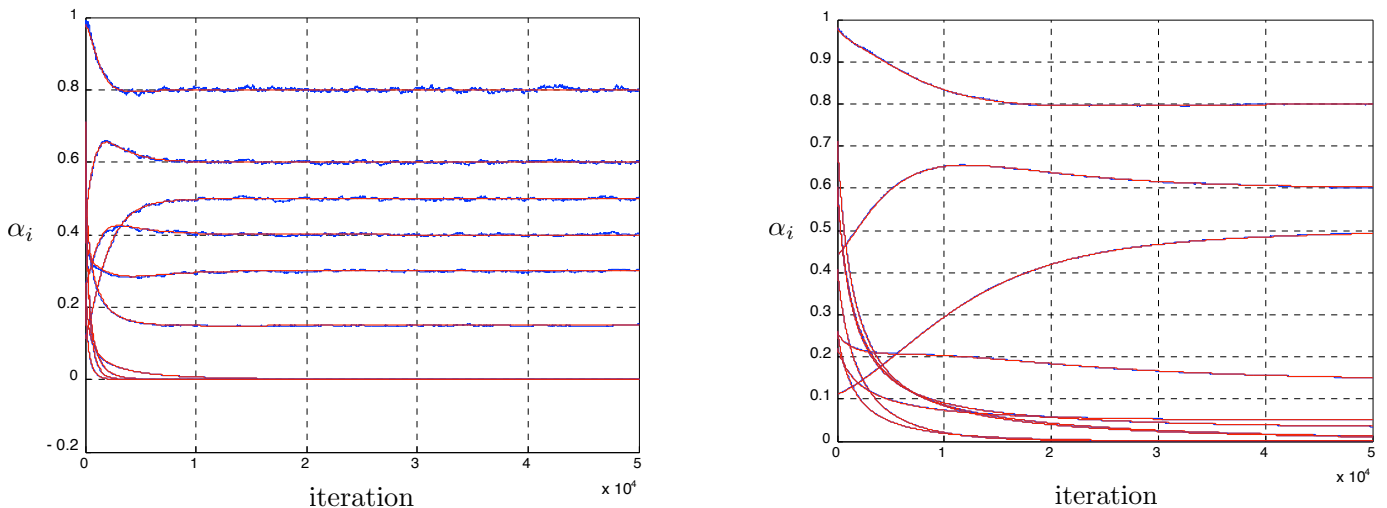


Fig. 4. Same experiment as in Figure 3 except that input sequence  $x(n)$  is generated by a first-order AR process. Two different step sizes are considered:  $\eta = 2.5 \times 10^{-3}$  on the left figure, and  $\eta = 5 \times 10^{-4}$  on the right figure.

It can be noticed that  $\alpha(n)$  no longer converges to  $\{\alpha^*\}_+$  since the input samples  $x(n)$  are now correlated. We can easily verify that  $E\{e^2(n)\} = 4.97$  dB using  $\{\alpha^*\}_+$ , and  $E\{e^2(n)\} = 3.82$  dB at convergence of the non-negative LMS algorithm.

#### IV. SECOND-ORDER MOMENT ANALYSIS

We now present a model for the behavior of the second-order moments of the adaptive weights. To allow further analysis progress, we assume in this section that the input  $x(n)$  is Gaussian.

### A. Second moment behavior model

Using  $e(n) = z(n) - \mathbf{v}^\top(n) \mathbf{x}(n)$ , neglecting the statistical dependence of  $\mathbf{x}(n)$  and  $\mathbf{v}(n)$ , and using the properties assumed for  $z(n)$  yields an expression for the mean-square estimation error (MSE):

$$\begin{aligned} E\{e^2(n)\} &= E\{(z(n) - \mathbf{v}^\top(n) \mathbf{x}(n))(z(n) - \mathbf{v}^\top(n) \mathbf{x}(n))\} \\ &= \sigma_z^2 + E\{\mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{v}(n)\} \\ &\approx \sigma_z^2 + \text{trace}\{\mathbf{R}_x \mathbf{K}(n)\}. \end{aligned} \quad (38)$$

Equation (26) clearly shows that the mean behavior of each coefficient is a function of a single diagonal entry of matrix  $\mathbf{R}_x \mathbf{K}(n)$ . In this case, approximation (28) could be used without compromising the accuracy of the resulting mean behavior model. This accuracy has been verified through Monte Carlo simulations in Section III-C. The MSE in (38), however, is a function of the trace of  $\mathbf{R}_x \mathbf{K}(n)$ . Thus, the effect of the second order moments of the weight-error vector entries on the MSE behavior becomes more significant than in (26), and in general cannot be neglected. Thus, we determine a recursion for  $\mathbf{K}(n)$  starting again from the weight error update equation (21).

Premultiplying Equation (21) by its transpose, taking the expected value, and using the statistical properties of  $z(n)$ ,<sup>3</sup> yields

$$\begin{aligned} \mathbf{K}(n+1) &= \mathbf{K}(n) - \eta \mathbf{P}_1(n) \mathbf{K}(n) - \eta \mathbf{K}(n) \mathbf{P}_1^\top(n) + \eta^2 \sigma_z^2 \mathbf{P}_2(n) + \eta^2 \sigma_z^2 [\mathbf{P}_3(n) + \mathbf{P}_3^\top(n)] \\ &\quad + \eta^2 \sigma_z^2 \mathbf{P}_4(n) - \eta [\mathbf{P}_5(n) + \mathbf{P}_5^\top(n)] + \eta^2 \mathbf{P}_6(n) + \eta^2 \mathbf{P}_7(n) + \eta^2 \mathbf{P}_8(n) + \eta^2 \mathbf{P}_9(n) \end{aligned} \quad (39)$$

where matrices  $\mathbf{P}_1$  to  $\mathbf{P}_9$  are defined by

$$\mathbf{P}_1 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n)\} \quad (40)$$

$$\mathbf{P}_2 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \quad (41)$$

$$\mathbf{P}_3 = E\{\mathbf{D}_x(n) \mathbf{v}(n) \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \quad (42)$$

$$\mathbf{P}_4 = E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \quad (43)$$

$$\mathbf{P}_5 = E\{\mathbf{v}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \quad (44)$$

$$\mathbf{P}_6 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \quad (45)$$

$$\mathbf{P}_7 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \quad (46)$$

$$\mathbf{P}_8 = E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \boldsymbol{\alpha}_*^\top \mathbf{D}_x(n)\} \quad (47)$$

$$\mathbf{P}_9 = E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \quad (48)$$

<sup>3</sup>The two important properties of  $z(n)$  used in evaluating (39) are its independence of any other signal and its zero-mean.

The expected values in (40)–(48) are calculated in the following. In order to keep the calculations mathematically tractable, the following statistical assumptions are employed:

**A1:** The input vector  $\mathbf{x}(n)$  is zero-mean Gaussian.

**A2:** The weight-error vector  $\mathbf{v}(n)$  is statistically independent of  $\mathbf{x}(n) \mathbf{x}^\top(n)$ . The reasoning for this approximation has been discussed in detail in [23].

**A3:** The statistical dependence of  $\mathbf{v}(n) \mathbf{v}^\top(n)$  and  $\mathbf{v}(n)$  is neglected. This assumption follows the same reasoning valid for assumption **A2**, see [23].

**A4:**  $\mathbf{v}(n)$  and  $(\mathbf{x}^\top(n) \mathbf{v}(n))^2$  are statistically independent given **A2**. This is because  $(\mathbf{x}^\top(n) \mathbf{v}(n))^2$  is a linear combination of the entries of  $\mathbf{v}(n) \mathbf{v}^\top(n)$ . Thus, this approximation follows basically the same reasoning discussed in [23] to justify **A2**.

### Expected value $P_1$

This expected value has been already calculated in (23). Remember that

$$\mathbf{P}_1 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^*\} = E\{\mathbf{D}_{\alpha^*} \mathbf{x}(n) \mathbf{x}^\top(n)\} = \mathbf{D}_{\alpha^*} \mathbf{R}_x. \quad (49)$$

### Expected value $P_2$

Basic linear algebra gives

$$\mathbf{P}_2 = E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} = E\{\mathbf{D}_{\alpha^*} \mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{D}_{\alpha^*}\} = \mathbf{D}_{\alpha^*} \mathbf{R}_x \mathbf{D}_{\alpha^*}. \quad (50)$$

### Expected value $P_3$

Neglecting the statistical dependence of  $\mathbf{x}(n)$  and  $\mathbf{v}(n)$  yields

$$\mathbf{P}_3 = E\{\mathbf{D}_x(n) \mathbf{v}(n) \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \approx E\{\mathbf{D}_v(n)\} \mathbf{R}_x \mathbf{D}_{\alpha^*}. \quad (51)$$

### Expected value $P_4$

The  $(i, j)$ -th entry of the matrix within the expectation in  $\mathbf{P}_4$  is given by

$$[\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)]_{ij} = x(n-i+1) v_i(n) v_j(n) x(n-j+1). \quad (52)$$

Using **A2**,  $E\{x(n-i+1) v_i(n) v_j(n) x(n-j+1)\} \approx E\{x(n-i+1) x(n-j+1)\} E\{v_i(n) v_j(n)\}$  and

$$\mathbf{P}_4 \approx \mathbf{R}_x \circ \mathbf{K}(n) \quad (53)$$

where  $\circ$  denotes the so-called Hadamard entry-wise product.

### Expected value $P_5$

Defining  $\mathbf{D}_v(n)$  as the diagonal matrix with diagonal entries given by  $\mathbf{v}(n)$ , we first note that

$$E\{\mathbf{v}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} = E\{\mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{D}_v(n)\}. \quad (54)$$

Now, using **A2** and **A3**, the expectation can be approximated as

$$E\{\mathbf{v}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \approx E\{\mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n)\} E\{\mathbf{D}_v(n)\}. \quad (55)$$

Finally, using again **A2** we obtain

$$\mathbf{P}_5 \approx \mathbf{K}(n) \mathbf{R}_x E\{\mathbf{D}_v(n)\}. \quad (56)$$

### Expected value $\mathbf{P}_6$

Basic linear algebra gives

$$\begin{aligned} \mathbf{P}_6 &= E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \\ &= \mathbf{D}_{\alpha^*} E\{\mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n)\} \mathbf{D}_{\alpha^*}. \end{aligned} \quad (57)$$

Under **A1** and applying the same methodology used to derive [24, Equation (29)],

$$\begin{aligned} \mathbf{P}_6 &\approx \mathbf{D}_{\alpha^*} \left( 2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x + E\{\mathbf{v}^\top(n) \mathbf{R}_x \mathbf{v}(n)\} \mathbf{R}_x \right) \mathbf{D}_{\alpha^*} \\ &= \mathbf{D}_{\alpha^*} \left( 2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x + E\{\text{trace}\{\mathbf{v}^\top(n) \mathbf{R}_x \mathbf{v}(n)\}\} \mathbf{R}_x \right) \mathbf{D}_{\alpha^*} \\ &= \mathbf{D}_{\alpha^*} (2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x + \text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x) \mathbf{D}_{\alpha^*}. \end{aligned} \quad (58)$$

### Expected value $\mathbf{P}_7$

Using basic algebra, **A2** and **A3** as done to obtain (55), we have

$$\begin{aligned} \mathbf{P}_7 &= E\{\mathbf{D}_x(n) \boldsymbol{\alpha}^* \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{D}_x(n)\} \\ &\approx \mathbf{D}_{\alpha^*} E\{\mathbf{x}(n) \mathbf{x}^\top(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{x}^\top(n)\} E\{\mathbf{D}_v(n)\}. \end{aligned} \quad (59)$$

Finally, under **A1** and applying the same methodology as in [24, Equation (29)], yields

$$\mathbf{P}_7 \approx \mathbf{D}_{\alpha^*} (2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x + \text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x) E\{\mathbf{D}_v(n)\}. \quad (60)$$

### Expected value $\mathbf{P}_8$

Using basic algebra we obtain

$$\begin{aligned} \mathbf{P}_8 &= E\{\mathbf{D}_x(n) \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \mathbf{v}^\top(n) \mathbf{x}(n) \boldsymbol{\alpha}^{*\top} \mathbf{D}_x(n)\} \\ &= E\{\mathbf{D}_v(n) (\mathbf{x}^\top(n) \mathbf{v}(n))^2 \mathbf{x}(n) \mathbf{x}^\top(n)\} \mathbf{D}_{\alpha^*}. \end{aligned} \quad (61)$$

Using **A4**,  $\mathbf{P}_8$  becomes

$$\mathbf{P}_8 \approx E\{\mathbf{D}_v(n)\} E\{(\mathbf{x}^\top(n) \mathbf{v}(n))^2 \mathbf{x}(n) \mathbf{x}^\top(n)\} \mathbf{D}_{\alpha^*}. \quad (62)$$

The expected value  $E\{(\mathbf{x}^\top(n) \mathbf{v}(n))^2 \mathbf{x}(n) \mathbf{x}^\top(n)\}$  for zero-mean Gaussian signal  $\mathbf{x}(n)$  has already been evaluated in [24, equations (7)–(9)], using results from [25]. Following the same procedure as in [24] yields

$$\begin{aligned} E\{(\mathbf{x}^\top(n) \mathbf{v}(n))^2 \mathbf{x}(n) \mathbf{x}^\top(n) | \mathbf{v}(n)\} &\approx \mathbf{v}^\top(n) \mathbf{R}_x \mathbf{v}(n) \mathbf{R}_x + 2 \mathbf{R}_x \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{R}_x \\ &= \text{trace}\{\mathbf{R}_x \mathbf{v}(n) \mathbf{v}^\top(n)\} \mathbf{R}_x + 2 \mathbf{R}_x \mathbf{v}(n) \mathbf{v}^\top(n) \mathbf{R}_x. \end{aligned} \quad (63)$$



Now, taking the expected value with respect to  $\mathbf{v}(n)$ ,

$$E\{(\mathbf{x}^\top(n) \mathbf{v}(n))^2 \mathbf{x}(n) \mathbf{x}^\top(n)\} \approx \text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x + 2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x. \quad (64)$$

Then we obtain the final result

$$\mathbf{P}_8 \approx E\{\mathbf{D}_v(n)\} (\text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x + 2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x) \mathbf{D}_{\alpha^*}. \quad (65)$$

### Expected value $\mathbf{P}_9$

Computing the  $(i, j)$ -th entry of matrix  $\mathbf{P}_9$  within the expectation, and using **A2**, yields

$$\begin{aligned} [\mathbf{P}_9]_{ij} &= \sum_{\ell} \sum_k E\{x(n-i+1) [\mathbf{v}(n) \mathbf{v}^\top(n)]_{ik} [\mathbf{x}(n) \mathbf{x}^\top(n)]_{k\ell} [\mathbf{v}(n) \mathbf{v}^\top(n)]_{\ell j} x(n-j+1)\} \\ &= \sum_{\ell} \sum_k E\{x(n-i+1) x(n-k+1) x(n-\ell+1) x(n-j+1)\} E\{v_i(n) v_j(n) v_k(n) v_\ell(n)\}. \end{aligned} \quad (66)$$

For  $\mathbf{x}(n)$  zero-mean Gaussian (**A1**), we know that [26]

$$E\{x(n-i+1) x(n-k+1) x(n-\ell+1) x(n-j+1)\} = r_x(k-i) r_x(j-\ell) + r_x(\ell-i) r_x(j-k) + r_x(j-i) r_x(\ell-k). \quad (67)$$

The expectation  $E\{v_i(n) v_j(n) v_k(n) v_\ell(n)\}$  cannot be evaluated directly, as the statistics of  $\mathbf{v}(n)$  are unknown. Approximate expressions can be obtained using numerous different approaches. We have chosen to use the following approximation which preserves relevant information about the second moment behavior of the adaptive weights while keeping the mathematical problem tractable. We first note that

$$E\{v_i(n) v_k(n) v_\ell(n) v_j(n)\} = E\{v_i(n) v_j(n)\} E\{v_k(n) v_\ell(n)\} + \text{Cov}\{v_i(n) v_j(n), v_k(n) v_\ell(n)\}. \quad (68)$$

Now, writing  $v_i(n+1) v_j(n+1) = (v_i(n) + \eta \Delta v_i(n)) (v_j(n) + \eta \Delta v_j(n))$ , we see that the fluctuations in  $v_i(n+1) v_j(n+1)$  are proportional to  $\eta$ . Using the same reasoning for  $v_k(n) v_\ell(n)$  we finally note that the covariance in (68) is proportional to  $\eta^2$ . The higher order moments of the entries of  $\mathbf{v}(n)$  in (68) will then be proportional to  $\eta^p$  with  $p \geq 2$ . Thus, for sufficiently small values of  $\eta$ , neglecting these terms yields the approximation

$$E\{v_i(n) v_k(n) v_\ell(n) v_j(n)\} \approx E\{v_i(n) v_j(n)\} E\{v_k(n) v_\ell(n)\}. \quad (69)$$

This approximation is supported by the simulation results presented in Section IV-B. Substituting the two equations above into the expression of  $[\mathbf{P}_9]_{ij}$  leads to

$$\begin{aligned} [\mathbf{P}_9]_{ij} &\approx r_x(j-i) \sum_{\ell} \sum_k r_x(\ell-k) [\mathbf{K}(n)]_{k\ell} [\mathbf{K}(n)]_{ij} \\ &\quad + \sum_{\ell} \sum_k r_x(k-i) r_x(j-\ell) [\mathbf{K}(n)]_{k\ell} [\mathbf{K}(n)]_{ij} \\ &\quad + \sum_{\ell} \sum_k r_x(\ell-i) r_x(j-k) [\mathbf{K}(n)]_{k\ell} [\mathbf{K}(n)]_{ij}. \end{aligned} \quad (70)$$

The first right-hand term of Equation (70) can be expressed as follows

$$\begin{aligned} r_x(j-i) \sum_{\ell} \sum_k r_x(\ell-k) [\mathbf{K}(n)]_{k\ell} [\mathbf{K}(n)]_{ij} &= [\mathbf{R}_x]_{ij} \sum_k \left( \sum_{\ell} [\mathbf{R}_x]_{k\ell} [\mathbf{K}(n)]_{k\ell} \right) [\mathbf{K}(n)]_{ij} \\ &= [\text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x]_{ij} [\mathbf{K}(n)]_{ij}. \end{aligned} \quad (71)$$

The second and third right-hand terms write

$$\begin{aligned} \sum_{\ell} \sum_k r_x(k-i) r_x(j-\ell) [\mathbf{K}(n)]_{k\ell} [\mathbf{K}(n)]_{ij} &= \left( \sum_{\ell} \sum_k [\mathbf{R}_x]_{ik} [\mathbf{K}(n)]_{k\ell} [\mathbf{R}_x]_{\ell j} \right) [\mathbf{K}(n)]_{ij} \\ &= [\mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x]_{ij} [\mathbf{K}(n)]_{ij}. \end{aligned} \quad (72)$$

This leads to the following close-form expression

$$[\mathbf{P}_9] = (\text{trace}\{\mathbf{R}_x \mathbf{K}(n)\} \mathbf{R}_x + 2 \mathbf{R}_x \mathbf{K}(n) \mathbf{R}_x) \circ \mathbf{K}(n). \quad (73)$$

Using the expected values  $\mathbf{P}_1$  to  $\mathbf{P}_9$  in Equation (39), we finally obtain a recursive analytical model for the behavior of  $\mathbf{K}(n)$ . This result can be used to study the convergence properties of  $E\{e^2(n)\}$ , and can be applied to design.<sup>4</sup> The next section illustrates the model accuracy in predicting the non-negative LMS algorithm behavior.

### B. Simulation examples for the second-order moment analysis

This section presents simulation examples to verify the accuracy of the model (39). Figures 5 and 6 show the behavior of the excess MSE  $J_{emse}(n) = \text{trace}\{\mathbf{R}_x \mathbf{K}(n)\}$  corresponding to the experiments conducted in Section III-C. The simulation curves (blue line) were obtained from Monte Carlo simulation averaged over 100 realizations. The theoretical curves (red line) were obtained from model (39). Note the model's accuracy even for step sizes as large as  $\eta_{\max}/2$  (left side of Figure 6). Also note that the theoretical value of the minimum excess mean-square error  $J_{emse_{min}}$  is represented in Figure 5.<sup>5</sup> It can be observed that  $J_{emse}(n)$  tends to  $J_{emse_{min}}$  as  $n$  goes to infinity. Figure 7 highlights the performance of the model for uncorrelated and correlated input signals  $x(n)$  through the same experimental setup as described before, except that the noise variance  $\sigma_z^2$  is now set to 1. All these experiments illustrate the accuracy of the model, which can provide important guidelines for the use of the non-negative LMS algorithm in practical applications.

<sup>4</sup>This model can also be used in (26) for the mean weight behavior if needed. However, our experience has been that the simplified model (28) suffices for predicting the mean weight behavior for most practical needs. It also makes the analytical study presented in Section III-B tractable.

<sup>5</sup>It can be easily shown, from Equations (17)–(19), that  $J_{emse_{min}} = \|\boldsymbol{\alpha}^* - (\boldsymbol{\alpha}^*)_+\|^2$  in the case where  $\mathbf{R}_x = \mathbf{I}$ .

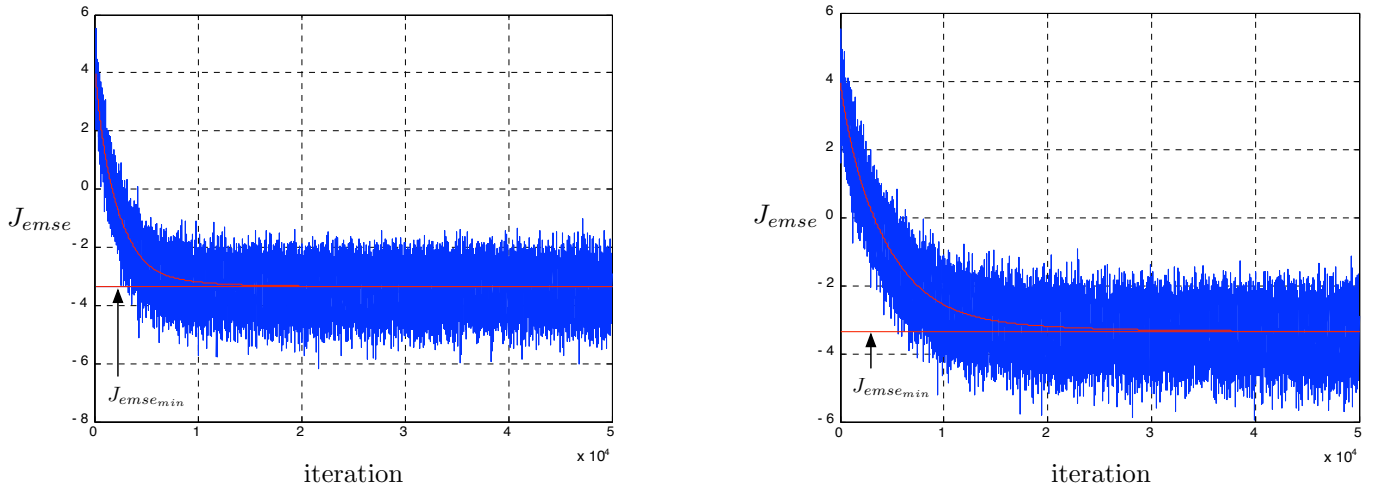


Fig. 5. Convergence of  $J_{emse}(n)$  in the case where input  $x(n)$  and noise  $z(n)$  are i.i.d. Two different step sizes are considered:  $\eta = 10^{-3}$  on the left figure, and  $\eta = 5 \times 10^{-4}$  on the right figure. The theoretical curves (red line) obtained from (38) and the model (39) and simulation curves (blue line) are perfectly superimposed.

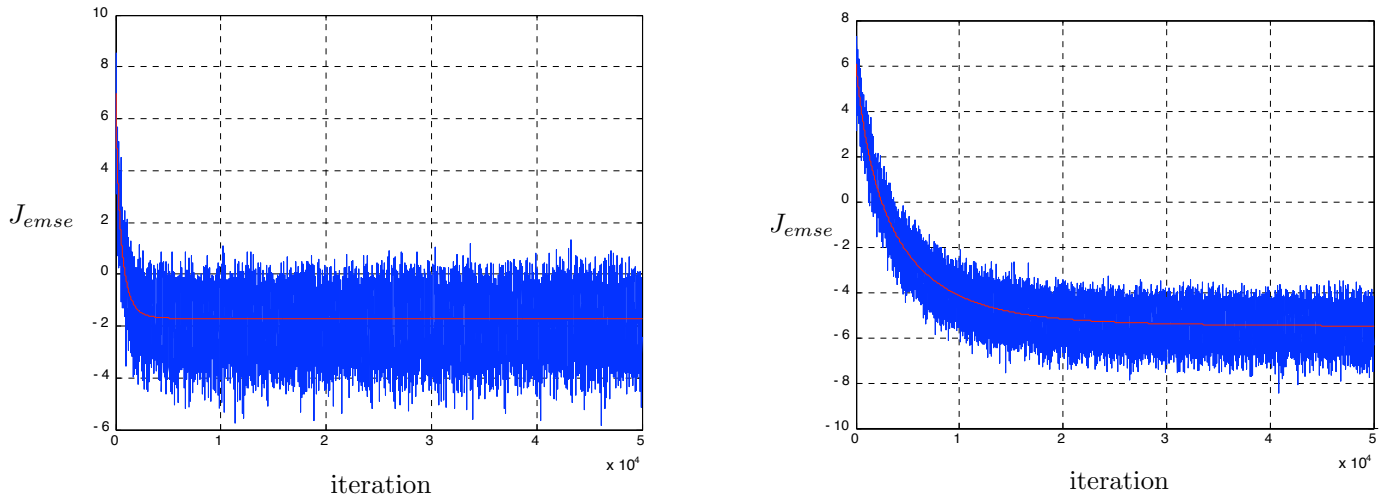


Fig. 6. Same experiment as in Figure 5 except that input sequence  $x(n)$  is generated by a first-order AR process. Two different step sizes are considered:  $\eta = 2.5 \times 10^{-3}$  on the left figure, and  $\eta = 5 \times 10^{-4}$  on the right figure.

## V. CONCLUSION

In many real-life phenomena, due to the inherent physical characteristics of systems under investigation, non-negativity is a desired constraint that can be imposed on the parameters to estimate in order to avoid physically absurd and uninterpretable results. In this paper, we proposed a general method for system identification under non-negativity constraints, and we derived the so-called non-negative LMS based on stochastic gradient descent. This algorithm switches automatically between a gradient descent mechanism and a gradient ascent one depending whether the non-negativity constraint is violated or not. Finally, we analyzed the algorithm convergence in the mean sense and in the mean-square sense. In future research

efforts, we intend to explore these models in practical applications since they provide important guidelines to algorithm designers. We also plan to derive variants of this approach, e.g., in the spirit of the normalized-LMS and the sign-LMS algorithms.

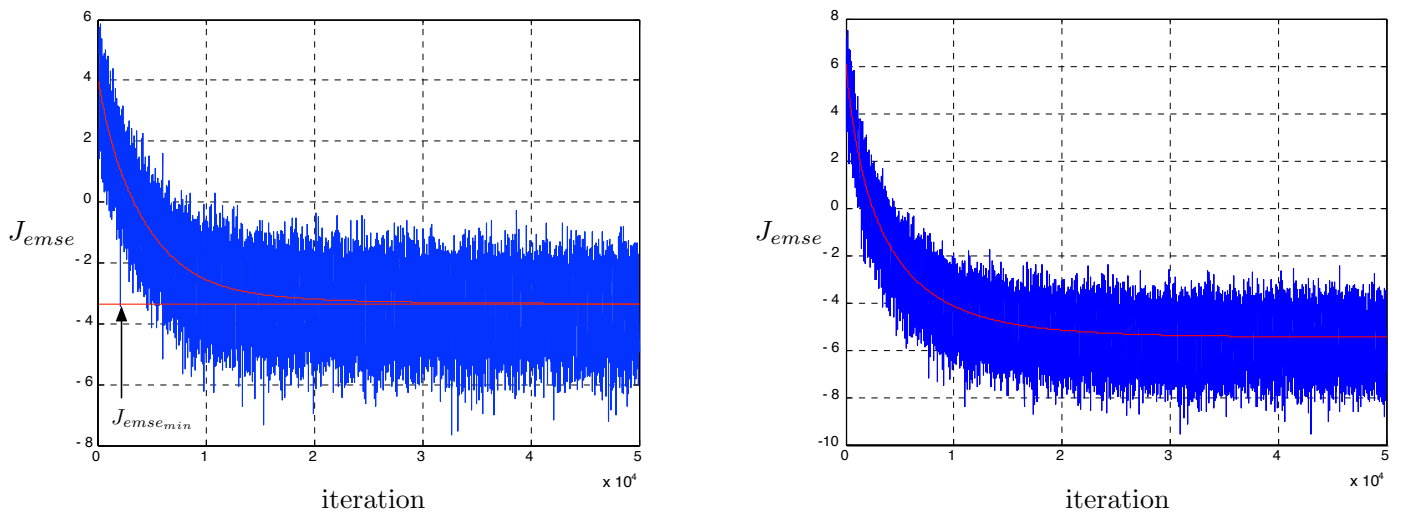


Fig. 7. Convergence of  $J_{emse}(n)$  with step size  $\eta = 5 \times 10^{-4}$ , in the case where input  $x(n)$  is i.i.d. on the left figure, and generated by a first-order AR process on the right figure. Compared to Figure 5 (right) and 6 (right), the variance of the noise  $z(n)$  has been increased from  $10^{-2}$  to 1.

## REFERENCES

- [1] F. Benvenuto, R. Zanella, L. Zanni, and M. Bertero, “Nonnegative least-squares image deblurring: improved gradient projection approaches,” *Inverse Problems*, vol. 26, no. 1, 2010.
- [2] M. H. Van Benthem and M. R. Keenan, “Fast algorithm for the solution of large-scale non-negativity-constrained least squares problems,” *Journal of Chemometrics*, vol. 18, pp. 441–450, 2004.
- [3] D.D. Lee and H.S. Seung, “Learning the parts of objects by non-negative matrix factorization,” *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.
- [4] D.D. Lee and H.S. Seung, “Algorithms for non-negative matrix factorization,” *Advances in neural information processing systems, NIPS*, pp. 556–562, 2001.
- [5] A. Cichocki, R. Zdunek, A. H. Phan, and S. Amari, *Nonnegative Matrix and Tensor Factorizations*, Wiley, 2009.
- [6] M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, and R. J. Plemmons, “Algorithms and applications for approximate nonnegative matrix factorization,” *Computational Statistics and Data Analysis*, vol. 52, no. 1, pp. 155–173, 2007.
- [7] M. D. Plumbley, “Algorithms for nonnegative independent component analysis,” *IEEE Transactions on Neural Networks*, vol. 14, no. 3, pp. 534–543, 2003.
- [8] S. Moussaoui, D. Brie, A. Mohammad-Djafari, and C. Carteret, “Separation of non-negative mixture of non-negative sources using a bayesian approach and MCMC sampling,” *IEEE Transactions on Signal Processing*, vol. 54, no. 11, pp. 4133–4145, 2006.
- [9] C. L. Lawson and R. J. Hanson, *Solving Least Squares Problems*, Society for Industrial and Applied Mathematics, 1995.
- [10] R. Bro and S. De Jong, “A fast non-negativity-constrained least squares algorithm,” *Journal of Chemometrics*, vol. 11, no. 5, pp. 393–401, 1997.
- [11] J. B. Rosen, “The gradient projection method for nonlinear programming. part 1: Linear constraints,” *Journal of the Society for Industrial and Applied Mathematics*, vol. 8, no. 1, pp. 181–217, 1960.
- [12] P. H. Calamai and J. J. Moré, “Projected gradient methods for linearly constrained problems,” *Mathematical Programming*, vol. 39, no. 1, pp. 93–116, 1987.
- [13] J. Barzilai and J. M. Borwein, “Two-point step size gradient methods,” *IMA Journal of Numerical Analysis*, vol. 8, no. 1, pp. 141–148, 1988.
- [14] S. Theodoridis, K. Slavakis, and I. Yamada, “Adaptive learning in a world of projections: a unifying framework for linear and nonlinear classification and regression tasks,” *IEEE Signal Processing Magazine*, vol. 28, no. 1, pp. 97–123, January 2011.
- [15] C. J. Lin, “Projected gradient methods for nonnegative matrix factorization,” *Neural Computation*, vol. 19, no. 10, pp. 2756–2779, 2007.
- [16] C. J. Lin, “On the convergence of multiplicative update algorithms for nonnegative matrix factorization,” *IEEE Transactions on Neural Networks*, vol. 18, no. 6, pp. 1589–1596, 2007.
- [17] H. Lantéri, M. Roche, O. Cuevas, and C. Aime, “A general method to devise maximum-likelihood signal restoration multiplicative algorithms with non-negativity constraints,” *Signal Processing*, vol. 81, no. 5, pp. 945–974, 2001.
- [18] S. Boyd and L. Vandenberghe, *Convex Optimization*, University Press, Cambridge, 2004.
- [19] Ali Sayed, *Adaptive Filters*, Wiley-Interscience, New York, 2008.
- [20] R. M. May, “Simple mathematical models with very complicated dynamics,” *Nature*, vol. 261, no. 10, pp. 459–467, 1976.
- [21] K. Alligood, T. Sauer, and J. A. Yorke, *Chaos: an introduction to dynamical systems*, Springer-Verlag, 1997.

- [22] Daniel Perrin, “La suite logistique et le chaos,” Tech. Rep., Département de Mathématiques d’Orsay, Université de Paris-Sud, France, 2008.
- [23] J. Minkoff, “Comment: On the unnecessary assumption of statistical independence between reference signal and filter weights in feedforward adaptive systems,” *IEEE Trans. Signal Process.*, vol. 49, no. 5, pp. 1109, May 2001.
- [24] P.I. Hubscher and J.C.M. Bermudez, “An improved statistical analysis of the least mean fourth (LMF) adaptive algorithm,” *IEEE Transactions on Signal Processing*, vol. 51, no. 3, pp. 664 – 671, March 2003.
- [25] N. J. Bershad, P. Celka, and J.-M. Vesin, “Stochastic analysis of gradient adaptive identification of nonlinear systems with memory for gaussian data and noisy input and output measurements,” *IEEE Transactions on Signal Processing*, vol. 47, no. 3, pp. 675 –689, March 1999.
- [26] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*, McGraw-Hill, New York, 3rd edition, 1991.