

Blind adaptive SIMO acoustic system identification using a locally optimal step-size

Mathieu Hu¹, Dushyant Sharma², Simon Doclo³, Mike Brookes¹ and Patrick A. Naylor¹

¹*Department of Electrical and Electronic Engineering, Imperial College London, UK*

²*Voicemail-To-Text Research, Nuance Communications Inc., Marlow, UK*

³*Department of Medical Physics and Acoustics and Cluster of Excellence Hearing4All, University of Oldenburg, Oldenburg, Germany*

Correspondence should be addressed to Mathieu Hu (mathieu.hu12@imperial.ac.uk)

ABSTRACT

Blind adaptive identification of a Single-Input Multiple-Output (SIMO) acoustic system has useful applications including acoustic environment sensing, source localization and, in combination with multichannel equalization, dereverberation. An empirically chosen step-size is usually employed in blind system identification algorithms based on cross-relation error minimization. Although some adaptive step-size approaches have been proposed in the literature, the derivations rely, in some cases, on coarse approximations. In this paper, a locally optimal adaptive-step size exploiting the algebraic nature of the problem is derived. Experimental results using simulated room impulse responses show that the proposed algorithm has higher initial convergence rate.

1. INTRODUCTION

Sound signals captured within an enclosed environment with microphones placed at a distance from the sound source are affected by reverberation, such that the received signal includes not only the source component via direct path propagation from the sound source to each microphone but also source components due to propagation via reflections from surfaces in the environment including walls, ceilings and hard objects. While reverberation is desirable in music, it may degrade the quality and the intelligibility of speech [1].

A possible approach to dereverberation is to view the problem as a channel equalization problem. The recorded signal $y_i(n)$ at the i^{th} microphone is indeed modeled as the convolution of the dry signal $s(n)$ with the Room Impulse Response (RIR) h_i corrupted by some additive noise $v_i(n)$:

$$\forall i \in \{1, 2, \dots, M\}, y_i(n) = h_i * s(n) + v_i(n) \quad (1)$$

where M is the number of microphones and $*$ denotes the convolution product.

In that approach, the M RIRs are first blindly

estimated, then, given these estimates $\hat{\mathbf{h}}_i = [\hat{h}_i(0) \ \hat{h}_i(1) \ \dots \ \hat{h}_i(L-1)]^T$ a set of inversion filters $\mathbf{g}_i = [g_i(0) \ g_i(1) \ \dots \ g_i(L_{\text{inv}}-1)]^T$ is designed such that [2]:

$$\sum_{i=1}^M \hat{\mathbf{H}}_i \mathbf{g}_i = \mathbf{d} \quad (2)$$

where \mathbf{d} is the desired equalized impulse response, $\hat{\mathbf{H}}_i$ denotes the convolution matrix of size $(L + L_{\text{inv}} - 1) \times L_{\text{inv}}$ associated with $\hat{\mathbf{h}}_i$. The symbols L and L_{inv} respectively denote the length of $\hat{\mathbf{h}}_i$ and \mathbf{g}_i . In [2], a method to perfectly invert the room acoustics is proposed, i.e. \mathbf{d} in (2) has only one non-zero input. However, that inverse filter design for \mathbf{g} is very sensitive to estimation errors in the estimate of h_i [3]. By allowing early reflections in the desired equalized impulse response, inverse filter design algorithms robust to estimation errors have been proposed in the literature, e.g. [3, 4, 5].

Estimations of the RIRs can be obtained by using Blind System Identification (BSI) algorithms such as the Multi-Channel Least Mean Square (MCLMS) [6] or the Nor-

malized Multichannel Frequency Domain Least Mean Square (NMFCLMS) [7]. These algorithms minimize the Cross Relation [8] (CR) error via a Least Mean Squares (LMS) scheme to estimate the RIRs. Several trials are, however, required to find a proper fixed step-size that gives the best compromise between fast convergence and stability. The Variable Step-Size Unconstrained MCLMS (VSS-UMCLMS) algorithm, proposed in [9], uses an optimal adaptive step-size to overcome that drawback. The derivations, however, are made under the assumption that no noise is present in the recorded signal, which is unrealistic in practice.

In this paper, we propose to exploit the algebraic nature of the BSI problem to derive an optimal adaptive step-size which does not require the recorded signal to be noiseless.

In Section 2, the principles underlying BSI algorithms relying on the CR error are given. In Section 3, the proposed method is derived. The simulation setups and results are presented in Section 4. Conclusions are drawn in Section 5.

2. BACKGROUND

With respect to the identifiability conditions [8], we assume that the source has a full-rank covariance matrix and the RIRs have no common zeros. The RIRs are assumed all to be of known length L .

2.1. Cross-relation error

Cross-relation error based BSI algorithms exploit the Single-Input-Multiple-Output (SIMO) structure in (1) to estimate the RIRs. Let us consider two channels i and j and let us denote by $x_i(n) = h_i * s(n)$ the noiseless reverberant signal at the i^{th} microphone. We then have:

$$x_i(n) * h_j = h_i * (s(n) * h_j) = h_i * x_j(n). \quad (3)$$

Therefore, estimations of the RIRs are given by minimizing the energy of the error $e_{ij}(n)$ across channels:

$$e_{ij}(n) = \mathbf{y}_i^T(n) \hat{\mathbf{h}}_j(n) - \mathbf{y}_j^T(n) \hat{\mathbf{h}}_i(n) \quad (4)$$

$$\chi(n) = \sum_{i=1}^{M-1} \sum_{j=i+1}^M e_{ij}^2(n) \quad (5)$$

where $\chi(n)$ is the cross-relation error across channels and $\mathbf{y}_i(n) = [y_i(n) \ y_i(n-1) \ \dots \ y_i(n-L+1)]^T$.

Introducing $\mathbf{h}_i = [\mathbf{h}_i(0) \ \mathbf{h}_i(1) \ \dots \ \mathbf{h}_i(L-1)]^T$, estimates $\hat{\mathbf{h}}(n) = [\hat{\mathbf{h}}_1^T(n) \ \hat{\mathbf{h}}_2^T(n) \ \dots \ \hat{\mathbf{h}}_M^T(n)]^T$ of the stacked RIRs $\mathbf{h} = [\mathbf{h}_1^T \ \mathbf{h}_2^T \ \dots \ \mathbf{h}_M^T]^T$ are obtained by minimizing the expectation of $\chi(n)$, subject to the constraint $\hat{\mathbf{h}} \neq \mathbf{0}_{ML \times 1}$, where $\mathbf{0}_{ML \times 1}$ is the null vector of size ML .

This minimization problem is equivalent to computing an eigenvector corresponding the smallest eigenvalue of the expectation of the cross-relation matrix $\hat{\mathbf{R}}(n)$ [6] defined as

$$\hat{\mathbf{R}} = \begin{bmatrix} \sum_{i=2}^M \hat{\mathbf{R}}_{y_i y_i} & -\hat{\mathbf{R}}_{y_2 y_1} & \dots & -\hat{\mathbf{R}}_{y_M y_1} \\ -\hat{\mathbf{R}}_{y_1 y_2} & \sum_{i \neq 2}^M \hat{\mathbf{R}}_{y_i y_i} & \dots & -\hat{\mathbf{R}}_{y_M y_2} \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{\mathbf{R}}_{y_1 y_M} & -\hat{\mathbf{R}}_{y_1 y_M} & \dots & \sum_{i=1}^{M-1} \hat{\mathbf{R}}_{y_i y_i} \end{bmatrix} \quad (6)$$

where $\hat{\mathbf{R}}_{y_i y_j} = \mathbf{y}_i \mathbf{y}_j^T$ and the time index n has been dropped. We will refer to the expectation of $\hat{\mathbf{R}}(n)$ as \mathbf{R} .

2.2. Baseline

The MCLMS algorithm [6] minimizes the Rayleigh Quotient

$$J(n) = \frac{\chi(n)}{\hat{\mathbf{h}}(n)^T \hat{\mathbf{h}}(n)} \quad (7)$$

to estimate the RIRs. The update equation is given by

$$\hat{\mathbf{h}}(n+1) = \frac{\hat{\mathbf{h}}(n) - \mu \nabla J(n)}{\|\hat{\mathbf{h}}(n) - \mu \nabla J(n)\|_2}, \quad (8)$$

$$\nabla J(n) = \frac{2}{\hat{\mathbf{h}}^T(n) \hat{\mathbf{h}}(n)} \{ \hat{\mathbf{R}}(n) \hat{\mathbf{h}}(n) - \chi(n) \hat{\mathbf{h}}(n) \} \quad (9)$$

where μ is a arbitrarily chosen constant step-size, $\|\cdot\|_2$ denotes the Euclidean norm and $\nabla J(n)$ is the gradient of $J(n)$.

In [9], an adaptive step-size $\mu_{\text{MCLMS}}(n)$ is derived such that the distance $\|\mathbf{h} - \alpha \hat{\mathbf{h}}(n)\|_2$ is minimized over α at each step in the absence of noise. The value of μ_{MCLMS} is given by

$$\mu_{\text{MCLMS}}(n) = \frac{\hat{\mathbf{h}}^T(n) \nabla J(n+1)}{\|\nabla J(n+1)\|_2^2}. \quad (10)$$

Neglecting the contribution of $\chi(n)$ in (9) as well as the normalization, the update of the VSS-UMCLMS is found:

$$\hat{\mathbf{h}}(n+1) = \hat{\mathbf{h}}(n) - 2\mu_{\text{MCLMS}}(n) \hat{\mathbf{R}}(n+1) \hat{\mathbf{h}}(n). \quad (11)$$

3. LOCALLY OPTIMAL STEP-SIZE

In the derivation of (10), the assumption that \mathbf{h} is in the null-space of $\hat{\mathbf{R}}(n)$ for all n is of key importance. In this section, the effect of additive noise on the solution given by the VSS-UMCLMS is analyzed. A locally optimal step-size, which does not require the input signal to be noiseless, is then derived. That locally optimal step-size exploits the algebraic nature of the problem as proposed in [10] for eigenvector estimation of fixed matrices.

3.1. Effect of additive noise

The effect of additive noise can be considered with reference to (6). Each the $\hat{\mathbf{R}}_{y_i y_j}$ blocks can be written as:

$$\begin{aligned}\hat{\mathbf{R}}_{y_i y_j} &= \mathbf{y}_i \mathbf{y}_j^T \\ &= \mathbf{x}_i \mathbf{x}_j^T + \mathbf{x}_i \mathbf{v}_j^T + \mathbf{x}_j \mathbf{v}_i^T + \mathbf{v}_i \mathbf{v}_j^T\end{aligned}\quad (12)$$

Taking the expectation of (12) under the assumptions that the source signal and the additive noise are uncorrelated and that the noise is zero-mean, we have

$$\mathbf{R}_{y_i y_j} = \mathbb{E}\{\mathbf{x}_i \mathbf{x}_j^T\} + \mathbb{E}\{\mathbf{v}_i \mathbf{v}_j^T\}$$

where \mathbb{E} denotes the expectation. In BSI, an eigenvector of the expected cross-relation matrix \mathbf{R} is sought. The expression of \mathbf{R} , in the presence of uncorrelated zero-mean additive noise, is given by

$$\mathbf{R} = \mathbf{R}_x + \mathbf{R}_v \quad (13)$$

where \mathbf{R}_x and \mathbf{R}_v are respectively the cross-relation matrices defined similarly to (6) for a reverberant noiseless input and a noise input.

In the absence of noise, \mathbf{R}_v is equal to $\mathbf{0}_{ML \times ML}$, the null matrix of size $ML \times ML$ and \mathbf{R} is a positive semi-definite matrix with a null-space of rank 1 spanned by \mathbf{h} [8].

In the presence of noise, however, \mathbf{R} is not necessarily positive semi-definite. If we assume that $\mathbf{v}_i(n)$ is a spatially uncorrelated white Gaussian noise, \mathbf{R}_v is an identity matrix of size $ML \times ML$ multiplied by the power of the noise. \mathbf{R} is then positive definite.

The update equation in (11) shows that the estimation procedure for $\hat{\mathbf{h}}(n+1)$ follows the gradient of a quadratic cost function given by $\hat{\mathbf{h}}^T(n) \hat{\mathbf{R}}(n+1) \hat{\mathbf{h}}(n)$. Therefore, in the presence of additive spatially white Gaussian noise, 0 is the minimum of the cost function and is achieved only at the trivial solution $\mathbf{0}_{ML \times 1}$.

3.2. Proposed method

Let us consider the standard update equation of a block LMS algorithm minimizing (7):

$$\hat{\mathbf{h}}(b+1) = \hat{\mathbf{h}}(b) - \mu(b) \nabla J(b) \quad (14)$$

where b is a block index and

$$\nabla J(b) = \frac{2}{\hat{\mathbf{h}}^T(b) \hat{\mathbf{h}}(b)} \{ \hat{\mathbf{R}}(b) \hat{\mathbf{h}}(b) - \chi(b) \hat{\mathbf{h}}(b) \}, \quad (15)$$

with $\hat{\mathbf{R}}(b) = \frac{1}{B} \sum_{l=0}^{B-1} \hat{\mathbf{R}}(bn_o + l)$, B and n_o respectively correspond to the block size and a sliding offset and $\chi(b) = \hat{\mathbf{h}}^T(b) \hat{\mathbf{R}}(b) \hat{\mathbf{h}}(b)$.

In (14), the updated estimate $\hat{\mathbf{h}}(b+1)$ is clearly a linear combination of the previous estimate $\hat{\mathbf{h}}(b)$ and the gradient $\nabla J(b)$ with a coefficient constrained to be 1 on $\hat{\mathbf{h}}$. If that constrain is removed, $\hat{\mathbf{h}}(b+1)$ can be written as:

$$\hat{\mathbf{h}}(b+1) = [\hat{\mathbf{h}}(b) \quad \nabla J(b)] \boldsymbol{\mu}_{\text{LOG}}(b) \quad (16)$$

where $\boldsymbol{\mu}_{\text{LOG}}(b) = \begin{bmatrix} \mu_1(b) \\ \mu_2(b) \end{bmatrix}$ is a 2×1 vector of weights.

Let us define the residual $r(b) = \hat{\mathbf{R}}(b) \hat{\mathbf{h}}(b+1) - \hat{\lambda}_{\min}(b+1) \hat{\mathbf{h}}(b+1)$ with $\hat{\lambda}_{\min}(b+1)$ an estimate of the smallest eigenvalue of \mathbf{R} . Since $\hat{\mathbf{h}}(b+1)$ belongs to the subspace spanned by $\hat{\mathbf{h}}(b)$ and $\nabla J(b)$, the new estimate is locally optimal when the residual is orthogonal to that subspace. The weight vector $\boldsymbol{\mu}_{\text{LOG}}$ is therefore such that

$$\begin{aligned} & \begin{bmatrix} \hat{\mathbf{h}}^T(b) \\ \nabla J(b)^T \end{bmatrix} \hat{\mathbf{R}}(b) [\hat{\mathbf{h}}(b) \quad \nabla J(b)] \boldsymbol{\mu}_{\text{LOG}} = \\ & \hat{\lambda}_{\min}(b+1) \begin{bmatrix} \hat{\mathbf{h}}^T(b) \\ \nabla J(b)^T \end{bmatrix} [\hat{\mathbf{h}}(b) \quad \nabla J(b)] \boldsymbol{\mu}_{\text{LOG}}. \end{aligned} \quad (17)$$

In other words, $\boldsymbol{\mu}_{\text{LOG}}(b)$ is a generalized eigenvector of the 2×2 matrices $\begin{bmatrix} \hat{\mathbf{h}}^T(b) \\ \nabla J^T(b) \end{bmatrix} \hat{\mathbf{R}}(b) [\hat{\mathbf{h}}(b) \quad \nabla J(b)]$ and $\begin{bmatrix} \hat{\mathbf{h}}^T(b) \\ \nabla J^T(b) \end{bmatrix} [\hat{\mathbf{h}}(b) \quad \nabla J(b)]$ for the generalized eigenvalue $\hat{\lambda}_{\min}(b+1)$.

Provided that $\hat{\mathbf{h}}(b)$ and $\nabla J(b)$ are made orthonormal before computing $\hat{\mathbf{h}}(b+1)$, $\boldsymbol{\mu}_{\text{LOG}}$ is an eigenvector corresponding to the smallest eigenvalue of

$$\hat{\mathbf{R}}'(b) = \begin{bmatrix} \hat{\mathbf{h}}^T(b) \\ \nabla J^T(b) \end{bmatrix} \hat{\mathbf{R}}(b) [\hat{\mathbf{h}}(b) \quad \nabla J(b)] \quad (18)$$

where $\hat{\mathbf{h}}(b) = \frac{\hat{\mathbf{h}}(b)}{\|\hat{\mathbf{h}}(b)\|_2}$ and $\nabla J(b)$ is the unit norm vector orthogonal to $\hat{\mathbf{h}}(b)$ obtained via the Gram-Schmidt orthonormalization process.

Since μ_{LOG} is defined up to a scaling factor, we chose the one that has a unit-norm so that $\hat{\mathbf{h}}(b+1)$ is also unit-norm. The update equation of the proposed Locally Optimal Gradient-MCLMS (LOG-MCLMS) is then given by

$$\hat{\mathbf{h}}(b+1) = \mu_1(b)\hat{\mathbf{h}}(b) + \mu_2(b)\nabla J(b) \quad (19)$$

where $\mu_1(b)$ and $\mu_2(b)$ are the components of the unit-norm eigenvector corresponding to the smallest eigenvalue of $\hat{\mathbf{R}}'(b)$.

4. SIMULATION

In this section, the proposed algorithm is evaluated against the baseline in various conditions.

4.1. Setup

The VSS-UMCLMS, NMCFLMS and LOG-MCLMS algorithms were evaluated for impulse responses of length $L = 128$, two different numbers of microphones, ($M = 3$ and $M = 5$), and two values of Signal-to-Noise Ratio (SNR)s ($\text{SNR} \in \{\infty, 30\text{dB}\}$). For each combination of L , M and SNR, the algorithms were evaluated for 50 realizations of the source and additive noise. The RIRs were generated using the image method [11] for a shoebox-shaped room of dimensions $5\text{ m} \times 6\text{ m} \times 3\text{ m}$, a reverberation time $T_{60} = 0.5\text{ s}$ and a sampling frequency $f_s = 8\text{ kHz}$. The length L was then fixed by truncating the obtained RIRs to the desired length. The input signals and the additive noise were white Gaussian noises, uncorrelated with each other. Although longer RIRs are present in practical scenarios, this setup allows us to compare the proposed algorithm against the baseline.

A sliding rectangular window of length $B = 2L$ overlapping by $n_o = L$ samples was used in the LOG-MCLMS. The constant step-size shown in [7] was used in the NMCFLMS and was equal to 0.8.

The accuracy of the different algorithms was evaluated using the scale-independent Normalized Projection Misalignment (NPM) [1] defined as follows:

$$\text{NPM}(\hat{\mathbf{h}}, \mathbf{h}) = 10 \log_{10} \left(\frac{\|\mathbf{h} - \frac{\hat{\mathbf{h}}^T \mathbf{h}}{\hat{\mathbf{h}}^T \hat{\mathbf{h}}} \hat{\mathbf{h}}\|_2^2}{\mathbf{h}^T \mathbf{h}} \right). \quad (20)$$

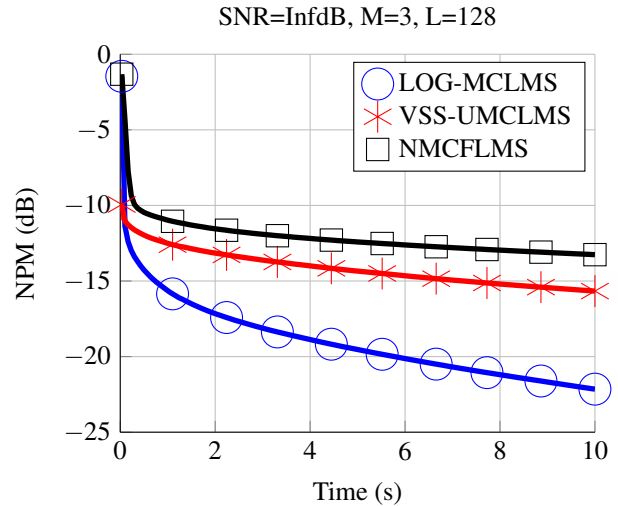


Fig. 1: Average NPM against time when no noise is present with $M = 3$

The norms of $\hat{\mathbf{h}}$ across time are given as well to verify that the algorithms are not converging to the trivial solution.

4.2. Result

Figures 1 and 2 show that when no noise is present in the recorded signal, the algorithms converge and that an increase in the number of microphone leads to estimates with lower NPMs as near common zeros between the channels are less likely [12]. In these conditions, the proposed algorithm achieves significantly lower NPMs on average than the VSS-UMCLMS and the NMCFLMS. For $M = 5$, for example, the VSS-UMCLMS and the NMCFLMS respectively achieve average NPMs of -13 dB and -20 dB after 10s of input data while the proposed algorithm achieves -32 dB .

As shown in Figures 3 and 4, in the presence of additive noise at $\text{SNR}=30\text{ dB}$, similar observations can be made although the improvement over the baseline is less significant. As is well known with cross-correlation error-based algorithms, misconvergence occurs when in the input contains significant noise [13]. As expected therefore, we observed such misconvergence in the proposed as well as baseline algorithms for noise levels greater than around $\text{SNR} = 30\text{ dB}$. In the presence of noise, the norm of $\hat{\mathbf{h}}(n)$ was observed and found to be close to zero after 10s of input data in the case of VSS-UMCLMS.

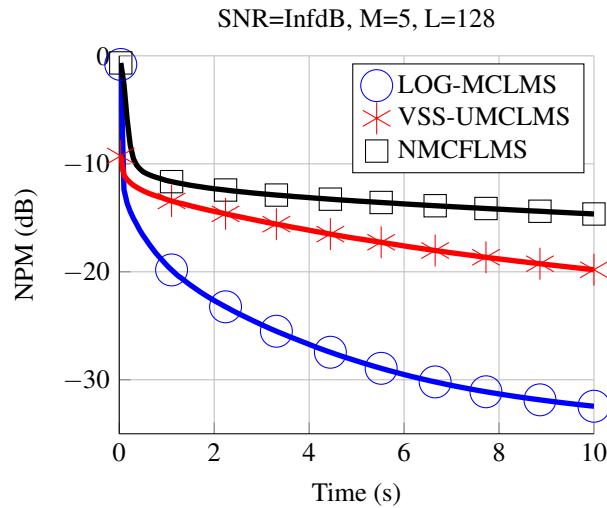


Fig. 2: Average NPM against time when no noise is present with $M = 5$

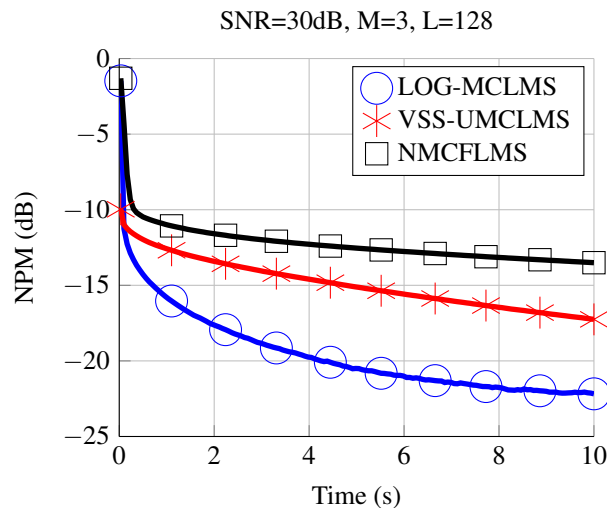


Fig. 3: Average NPM against time with SNR = 30dB and $M = 3$

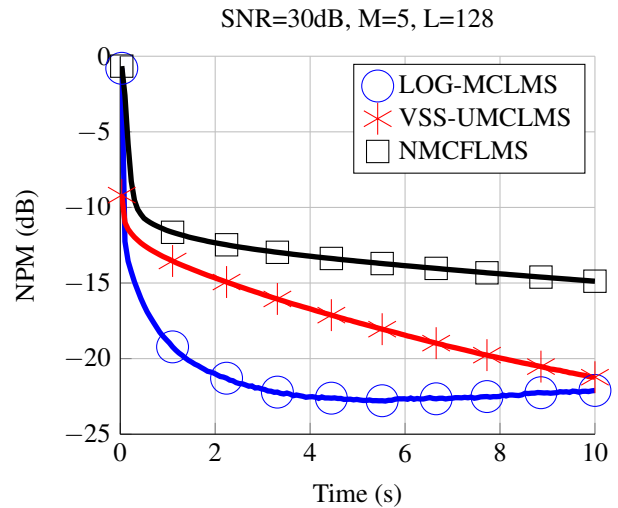


Fig. 4: Average NPM against time with SNR = 30dB and $M = 5$

5. CONCLUSION

In this paper, we have proposed a locally optimal adaptive-step size, for blind SIMO acoustic system identification. The formulation of the adaptive step-size exploits the eigenvector estimation framework of cross-correlation error-based BSI. Although the proposed adaptive step-size algorithm does not solve the well known misconvergence problem, the adaptive step-size shown in this paper leads to improved initial convergence rates in our test, measured in terms of NPM.

6. ACKNOWLEDGEMENT

The research leading to these results has received funding from the European Unions Seventh Framework Programme (FP7/2007-2013) under grant agreement n ITN-GA-2012-316969

7. REFERENCES

- [1] P. A. Naylor and N. D. Gaubitch, Eds., *Speech Dereverberation*, Springer, 2010.
- [2] M. Miyoshi and Y. Kaneda, "Inverse filtering of room acoustics," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 36, no. 2, pp. 145–152, Feb. 1988.
- [3] W. Zhang, E. A. P. Habets, and P. A. Naylor, "On the use of channel shortening in multichannel

- acoustic system equalization,” in *Proc. Intl. Workshop Acoust. Echo and Noise Control (IWAENC)*, Tel Aviv, Israel, Aug. 2010.
- [4] I. Kodrasi, S. Goetze, and S. Doclo, “Regularization for partial multichannel equalization for speech dereverberation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 21, no. 9, pp. 1879–1890, Sept. 2013.
- [5] Felicia Lim, Wancheng Zhang, Emanuël A. P. Habets, and Patrick A. Naylor, “Robust multichannel dereverberation using relaxed multichannel least squares,” *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, vol. 22, no. 9, pp. 1379–1390, Sept. 2014.
- [6] Y. Huang and J. Benesty, “Adaptive multi-channel least mean square and Newton algorithms for blind channel identification,” *Signal Processing*, vol. 82, pp. 1127–1138, Aug. 2002.
- [7] Y. Huang and J. Benesty, “A class of frequency-domain adaptive approaches to blind multichannel identification,” *IEEE Trans. Signal Process.*, vol. 51, no. 1, pp. 11–24, Jan. 2003.
- [8] G. Xu, H. Liu, L. Tong, and T. Kailath, “A least-squares approach to blind channel identification,” *IEEE Trans. Signal Process.*, vol. 43, no. 12, pp. 2982–2993, Dec. 1995.
- [9] J. Huang, J. Benesty, and J. Chen, “Optimal step size of the adaptive multichannel lms algorithm for blind simo identification,” *Signal Processing Letters, IEEE*, vol. 12, no. 3, pp. 173–176, March 2005.
- [10] Andrew V. Knyazev, “Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method,” *SIAM Journal on Scientific Computing*, vol. 23, no. 2, pp. 517–541, 2001.
- [11] J. B. Allen and D. A. Berkley, “Image method for efficiently simulating small-room acoustics,” *J. Acoust. Soc. Am.*, vol. 65, no. 4, pp. 943–950, Apr. 1979.
- [12] Xiang (Shawn) Lin, Andy W. H. Khong, and Patrick A Naylor, “A forced spectral diversity algorithm for speech dereverberation in the presence of near-common zeros,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 3, pp. 888–899, Mar. 2012.
- [13] M. K. Hasan and P. A. Naylor, “Analyzing effect of noise on LMS-type approaches to blind estimation of SIMO channels: robustness issue,” in *Proc. European Signal Processing Conf. (EUSIPCO)*, Florence, Italy, Sept. 2006.