

Accepted Manuscript

A Class of Stochastic Gradient Algorithms with Exponentiated Error Cost Functions

C. Boukis, D.P. Mandic, A.G. Constantinides

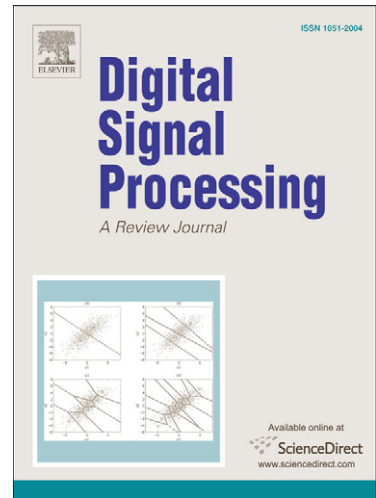
PII: S1051-2004(08)00181-4
DOI: [10.1016/j.dsp.2008.11.006](https://doi.org/10.1016/j.dsp.2008.11.006)
Reference: YDSPR 872

To appear in: *Digital Signal Processing*

Received date: 30 October 2007
Revised date: 16 June 2008
Accepted date: 23 November 2008

Please cite this article as: C. Boukis, D.P. Mandic, A.G. Constantinides, A Class of Stochastic Gradient Algorithms with Exponentiated Error Cost Functions, *Digital Signal Processing* (2008), doi: [10.1016/j.dsp.2008.11.006](https://doi.org/10.1016/j.dsp.2008.11.006)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



A Class of Stochastic Gradient Algorithms with Exponentiated Error Cost Functions

C. Boukis^{a,*}, D.P. Mandic^b and A.G. Constantinides^b

^a*Athens Information Technology, Autonomic and Grid Computing Group,
Markopoulo Ave, Peania/Athens 19002, Greece*

^b*Imperial College, Electrical and Electronic Engineering Dept, Communications
and Signal Processing Group, Exhibition Rd, London SW7 2BT, UK*

Abstract

A novel class of stochastic gradient descent algorithms is introduced based on the minimisation of convex cost functions with exponential dependence on the adaptation error, instead of the conventional linear combinations of even moments. The derivation is supported by rigorous analysis of the necessary conditions for convergence, the steady state mean square error is calculated and the optimal solutions in the least exponential sense are derived. The normalisation of the associated step size is also considered in order to fully exploit the dynamics of the input signal. Simulation results support the analysis.

Key words: Adaptive Filtering, Cost Functions, Stochastic Gradient Descent, Online Optimisation

PACS: 43.60.Mn, 45.10.Db

1 Introduction

Gradient descent (GD) algorithms aim at updating the coefficients of an adaptive tap-delay filter in a recursive manner, in order to minimise a chosen cost function. This is achieved in the following way

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu(n)f[e(n)]g[\mathbf{x}_n]$$

* Corresponding author

Email addresses: cbou@ait.edu.gr (C. Boukis), d.mandic@imperial.ac.uk (D.P. Mandic), agc@imperial.ac.uk (A.G. Constantinides).

24 where $\mathbf{w}_n = [w_0(n), w_1(n), \dots, w_{N-1}(n)]^t$ is the vector of filter coefficients, $e(n)$
 25 the adaptation error defined as the difference between the desired response
 26 and the output of the adaptive filter, $\mathbf{x}_n = [x(n), x(n-1), \dots, x(n-N+1)]^t$
 27 the input regressor vector, $\mu(n)$ the time-varying learning rate and $(\cdot)^t$ the
 28 vector transpose operator [1]. The (possibly nonlinear) functions $f[\cdot]$ and $g[\cdot]$,
 29 and consequently the performance of GD algorithms, depend critically on the
 30 choice of the cost function.

31 The vast majority of GD algorithms use quadratic cost functions, due to their
 32 mathematical tractability and convenience of analysis. In this case the func-
 33 tions $f[\cdot]$ and $g[\cdot]$ are linear. The so derived second order statistics (SOS)
 34 based algorithms have low computational complexity. Members of this class
 35 are the least mean square (LMS) [2] and the normalised least mean square
 36 (NLMS) [3].

37 Using high order even powers of the adaptation error as cost functions (non-
 38 linear $f[\cdot]$ and $g[\cdot]$) results in higher order statistics (HOS) adaptive algo-
 39 rithms [4]. These algorithms have potentially faster convergence than SOS
 40 based algorithms, due to their steeper error surfaces, that is they penalise
 41 heavier for deviations from the optimal solution. Moreover, unless the prob-
 42 ability distribution of the measurement noise is Gaussian, HOS algorithms
 43 exhibit reduced misadjustment as compared to SOS algorithms. Typical rep-
 44 resentatives of this class are the least mean fourth (LMF) [4] and the least
 45 mean kurtosis (LMK) [5] algorithm, which have been shown to outperform the
 46 LMS, especially in the presence of non-Gaussian additive measurement noise.

47 Mixed norm GD algorithms, that are robust under several noise conditions,
 48 can be derived when using finite sums of even error powers as cost functions.
 49 To that end Chambers, Tanrilulu and Constantinides introduced the Least
 50 Mean Mixed Norm (LMMN) algorithm [6] where second and fourth order
 51 moments were linearly combined in a convex manner. Later on, Chambers
 52 and Avlonitis presented the Robust Mixed Norm (RMN) algorithm [7] which
 53 is based on a convex mixing of the L1 and L2 norms. A normalised version
 54 of RMN was introduced in [8]. A generalisation of the mixed norm approach
 55 was introduced by Barros et al termed the weighted even moments (WEM)
 56 algorithm [9]. This algorithm is general enough to cater for as many even error
 57 powers as necessary, however, the weighting coefficients of the error powers
 58 need to be determined empirically.

59 In this paper, we propose stochastic gradient adaptation based on cost func-
 60 tions that have exponential dependence on the chosen error. Contrary to exist-
 61 ing approaches, this class of functions takes into account an infinite number of
 62 even moments of the error, resulting in nonlinear functions $f[\cdot]$ and $g[\cdot]$. Expo-
 63 nentiated error cost functions have much steeper surfaces than linear combina-
 64 tions of even moments, thus penalising heavily for deviation from the optimal

65 solution. Simulations in a system identification setting have shown that the
 66 proposed least exponentials class of algorithms (LE) outperform least mean
 67 square (LMS) algorithms in terms of convergence, together with increased
 68 robustness in the presence of impulsive noise.

69 The proposed LE class algorithms differ from the exponentiated gradient (EG)
 70 algorithms [10–12], since in our approach the coefficient adaptation formula
 71 is additive while the latter use multiplicative updating formulas. In addition,
 72 the cost function within EG algorithm attempts to minimise is the square of
 73 the error, while LE algorithms aim at the minimisation of exponentiated error
 74 cost functions.

75 Section 2 introduces the class of exponentiated error functions and Section
 76 3 presents the associated least exponential algorithms. The performance of
 77 these algorithms is then analysed in Section 4 within the energy conservation
 78 framework [14]. Simulation results are presented in Section 5 and Section 6
 79 concludes the paper.

80 2 Exponentiated Error Cost Functions

81 In order for GD algorithms to converge to global minima of error surfaces,
 82 they employ convex and unimodal cost functions. The most general choice of
 83 such cost functions is based on linear combination of even error powers [6,9],
 84 that is

$$J(n) = \sum_{i=1}^M \alpha_i e^{2i}(n) \quad (1)$$

85 where α_i is the weighting factor that represents the contribution of of the $(2i)$ -
 86 th power of the adaptation error $e(n)$. In the case of tap-delay (transversal)
 87 filters, the output error $e(n)$ is given by

$$e(n) = d(n) - \mathbf{w}_n^t \mathbf{x}_n \quad (2)$$

88 where $d(n)$ is the desired response. Functions of the form (1) are convex with
 89 a single minimum at $e(n) = 0$. Choosing $M = 1$ yields second order statistics
 90 (SOS) algorithms (e.g. the least mean square [2]), while for $k > 1$ we have
 91 higher order statistics (HOS) algorithms [4,6,9] (e.g. for $M = 2$ the least mean
 92 mixed norm algorithm is derived [6]). In most of the cases, the coefficients α_k
 93 are chosen empirically.

94 To reduce the dependence on an empirical choice of the parameters and to
 95 provide a closed form solution, we propose cost functions that have exponential
 96 dependence on the error. These are convex, unimodal and have steeper error
 97 surfaces than those given by (1). Moreover, they take into account all the even
 98 moments of the adaptation error. Two such functions are considered in this
 99 paper: the exponential of the squared error and the hyperbolic cosine. The
 100 cost function based on the exponential of the squared error (Fig. 1) is given
 101 by

$$J_{e2}(n) = \frac{1}{2} \exp[e^2(n)] \quad (3)$$

102 Evaluation of (3) as a Taylor Series Expansion (TSE) around zero gives

$$J_{e2}(n) = \frac{1}{2} \sum_{i=0}^{+\infty} \frac{1}{i!} e^{2i}(n) \quad (4)$$

103 indicating that the objective function $J_{e2}(n)$ takes into account all the even
 104 moments of the adaptation error. Also, as desired, since the weight of the
 105 $(2i)$ -th error moment is $1/(i!)$ more emphasis is given to lower order moments.

106 The second cost function considered is the hyperbolic cosine (sum of error
 107 exponentials) given by

$$J_{se}(n) = \frac{1}{2} \left(\exp[e(n)] + \exp[-e(n)] \right) \quad (5)$$

108 As illustrated in Fig. 1, $J_{se}(n)$ is less steep than $J_{e2}(n)$ and both are steeper
 109 than a quadratic function. This can also be observed from the TSE of $J_{se}(n)$
 110 around $e(n) = 0$

$$J_{se}(n) = \sum_{i=0}^{+\infty} \frac{1}{(2i)!} e^{2i}(n) \quad (6)$$

111 where the $2i$ -th power of the adaptation error is weighted by $1/(2i)!$. So $J_{se}(n)$
 112 emphasises less than $J_{e2}(n)$ on the high-order error moments. Notice also
 113 that for small error values, the exponentiated error cost functions become
 114 quadratic, that is

$$J_{e2}(n) \approx \frac{1}{2} \left(1 + e^2(n) \right) \quad (7)$$

115 and

$$J_{se}(n) = 1 + \frac{1}{2}e^2(n) \quad (8)$$

116 The LE algorithms therefore reduce to LMS for small errors; the term "1" in
117 (8) is irrelevant when taking gradients.

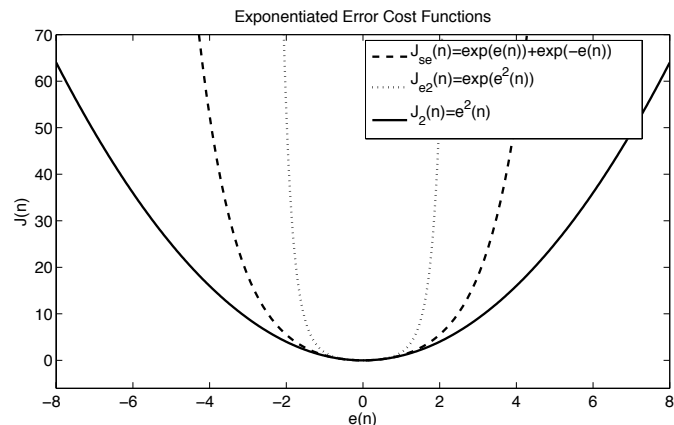


Fig. 1. Comparison of the standard and proposed cost functions.

118 3 The Class of Least Exponential Algorithms

119 Based on the gradient of $J(n)$ from (1) we have a general stochastic gradient
120 descent update formula [1]

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \mu(n) \frac{\partial J(n)}{\partial \mathbf{w}_n} \quad (9)$$

121 which yields

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu(n) \mathbf{x}_n \sum_{i=1}^M 2i\alpha_i e^{2i-1}(n) \quad (10)$$

122 For $M = 1$, $\alpha_1 = 1/2$ and a constant step size $\mu(n) = \mu$ this simplifies
123 into the least mean square (LMS) algorithm. For $M = 2$ the proposed algo-
124 rithms become similar to the least mean mixed norm (LMMN) algorithm [6].
125 Choosing $M > 2$ and heuristically finding the most appropriate values of α_k
126 ($k = 1, 2, \dots, M$) results in the weighted even moments (WEM) algorithm [9].

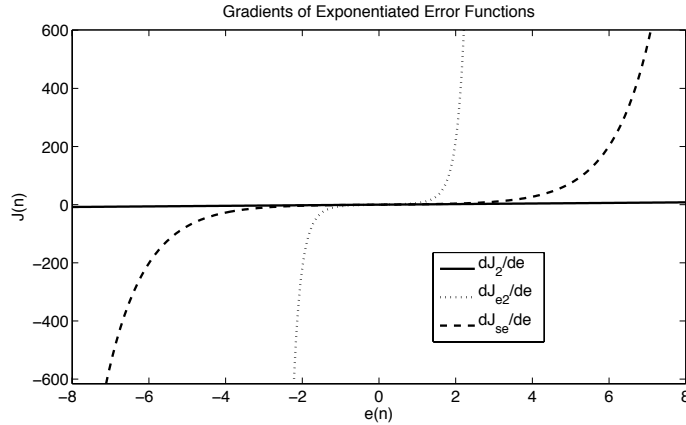


Fig. 2. The gradients of the square, exponential of the square, and sum of exponentials cost functions.

127 3.1 The Least Exponentiated Square (LE2) Algorithm

128 Taking the gradient of $J_{e2}(n)$ with respect to the vector of the filter coefficients
129 yields

$$\frac{\partial J_{e2}(n)}{\partial \mathbf{w}_n} = -e(n)\mathbf{x}(n) \exp[e^2(n)] \quad (11)$$

130 Substituting (11) into the general SGD update from (9), and assuming a time
131 invariant step size results in the least exponentiated square (LE2) algorithm,
132 which updates its coefficient estimates according to

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu e(n)\mathbf{x}_n \exp[e^2(n)] \quad (12)$$

133 Since the cost function $J_{e2}(n)$ is much steeper than the square of the error
134 (Fig. 1), the value of the gradient $\partial J_{e2}(n)/\partial \mathbf{w}_n$ is significantly larger than that
135 of the gradient of the squared error with respect to the coefficients (Fig. 2).
136 Hence, the LE2 converges faster than the least mean square (LMS) algorithm
137 for a given constant learning rate, provided stability conditions. The nonlinear-
138 earity of this algorithm with respect to the adaptation error is obvious from

$$f_{e2}[e(n)] = e(n)\exp[e^2(n)] = \sum_{i=0}^{+\infty} \frac{1}{i!} e^{2i+1}(n) \quad (13)$$

139 Using (13), the weight update (12) can be re-written as

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu \mathbf{x}_n \sum_{i=0}^{+\infty} \frac{1}{i!} e^{2i+1}(n) \quad (14)$$

140 illustrating that LE2 algorithm comprises the odd moments of the adaptation
141 error.

142 3.2 The Least Sum of Exponentials (LSE) Algorithm

143 The Least Sum of Exponentials (LSE) algorithm is derived by substituting
144 $\partial J(n)/\partial \mathbf{w}_n$ in the general SGD formula given by (9), with the partial gradient
145 of $J_{se}(n)$ w.r.t. the coefficients vector \mathbf{w}_n . Its update is given by

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \frac{\mu}{2} \mathbf{x}_n n (\exp[e(n)] - \exp[-e(n)]) \quad (15)$$

146 since

$$\frac{\partial J_{se}(n)}{\partial \mathbf{w}_n} = -\frac{1}{2} \mathbf{x}_n n (\exp[e(n)] - \exp[-e(n)]) \quad (16)$$

147 The error nonlinearity in the recursion of LSE is therefore

$$f_{se}[e(n)] = \frac{1}{2} [\exp[e(n)] - \exp[-e(n)]] = \sum_{i=0}^{+\infty} \frac{1}{(2i+1)!} e^{2i+1}(n) \quad (17)$$

148 Combining (15) and (17) yields the final weight update in the form

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \mu \mathbf{x}_n \sum_{i=0}^{+\infty} \frac{1}{(2i+1)!} e^{2i+1}(n) \quad (18)$$

149 The only difference between LSE and LE2 is the fact that the $(2i+1)$ -th
150 power is weighted by $1/(2i+1)!$ in the LSE update, whereas in the LE2 the
151 coefficient associated with the same error power is $1/i!$.

152 4 Convergence Analysis

153 In this section the performance of the least exponential algorithms is exam-
154 ined in terms of the optimal solution, mean-square stability and steady state

155 behaviour. This is achieved mainly based on the energy conservation frame-
 156 work [14,15], which relies on the observation that [16]

$$\|\tilde{\mathbf{w}}_{n+1}\|^2 + \frac{1}{\|\mathbf{x}_n\|^2} |e_a(n)|^2 = \|\tilde{\mathbf{w}}_n\|^2 + \frac{1}{\|\mathbf{x}_n\|^2} |e_p(n)|^2 \quad (19)$$

157 where $e_a(n)$ is the a priori error given by

$$e_a(n) = [\mathbf{w}_o - \mathbf{w}_n]^t \mathbf{x}_n = \tilde{\mathbf{w}}_n^t \mathbf{x}_n \quad (20)$$

158 and

$$e_p(n) = [\mathbf{w}_o - \mathbf{w}_{n+1}]^t \mathbf{x}_n = \tilde{\mathbf{w}}_{n+1}^t \mathbf{x}_n \quad (21)$$

159 is the a posteriori error.

160 In this analysis we also make the following standard assumptions [14]

- The desired response is produced by a linear model given by

$$d(n) = \mathbf{w}_o^t \mathbf{x}_n + v(n) \quad (22)$$

161 where \mathbf{w}_o is a vector containing the optimal coefficient values, and $v(n)$ is
 162 an additive noise component;

- 163 • The noise sequence $\{v(n)\}$ is independent, identically distributed and inde-
 164 pendent of the input sequence $\{\mathbf{x}_n\}$;
- 165 • The filter is long enough such that the a priori error is Gaussian. This
 166 assumption implies that the systematic component of the unknown signal
 167 is adequately modelled and hence the modelling error is not biased;
- The random variables $\|\mathbf{x}_n\|^2$ and $f^2[e(n)]$ are asymptotically uncorrelated.
 This assumption can be mathematically expressed as

$$\lim_{n \rightarrow \infty} E \left[\|\mathbf{x}_n\|^2 f^2[e(n)] \right] = E \left[\|\mathbf{x}_n\|^2 \right] \lim_{n \rightarrow \infty} E \left[f^2[e(n)] \right]$$

168 4.1 Optimal Solution

169 In order for the behaviour of the algorithm to be controllable the optimal
 170 solution should be unique (unimodal error surfaces). This is also the case with
 171 the exponentiated square and the sum of exponentials cost functions.

172 The optimal solution of the LE2 algorithm can be analysed based on

$$-\frac{1}{2} \left(E \left\{ d(n) \mathbf{x}_n \exp[e^2(n)] \right\} - E \left\{ \mathbf{x}_n \mathbf{x}_n^t \exp[e^2(n)] \right\} \mathbf{w}_{o,e2} \right) = 0 \quad (23)$$

173 where \mathbf{w}_o is the optimal solution in the least exponentiated square sense.
 174 Substituting the TSE of $\exp[e(n)]$ into (23) results in

$$E \left\{ d(n) \mathbf{x}_n \left[\sum_{i=0}^{+\infty} \frac{1}{i!} e^{2i}(n) \right] \right\} = E \left\{ \mathbf{x}_n \mathbf{x}_n^t \left[\sum_{i=0}^{+\infty} \frac{1}{i!} e^{2i}(n) \right] \right\} \mathbf{w}_{o,e2} \quad (24)$$

175 Consequently, the optimal solution in the least exponential squares sense is
 176 given by

$$\mathbf{w}_{o,e2} = \left[\mathbf{R} + \sum_{i=1}^{+\infty} \frac{1}{i!} E \left\{ \mathbf{x}_n \mathbf{x}_n^t e^{2i}(n) \right\} \right]^{-1} \cdot \left[\mathbf{p} + \sum_{i=1}^{+\infty} \frac{1}{i!} E \left\{ d(n) \mathbf{x}_n e^{2i}(n) \right\} \right] \quad (25)$$

177 where $\mathbf{R} = E \left\{ \mathbf{x}_n \mathbf{x}_n^t \right\}$ is the autocorrelation matrix of the input signal and
 178 $\mathbf{p} = E \left\{ d(n) \mathbf{x}_n \right\}$ the cross-correlation between the input \mathbf{x}_n and the desired
 179 signal $d(n)$. Moreover, assuming that $\{\mathbf{x}_n\}$ and $\{e(n)\}$ are asymptotically
 180 uncorrelated yields

$$\mathbf{w}_{o,e2} = \left[\mathbf{R} \sum_{i=0}^{+\infty} \frac{1}{i!} E \left\{ e^{2i}(n) \right\} \right]^{-1} \cdot \left[\mathbf{p} \sum_{i=0}^{+\infty} \frac{1}{i!} E \left\{ e^{2i}(n) \right\} \right] \quad (26)$$

181 Unless the measurements of the desired signal $d(n)$ are extremely noisy, the
 182 adaptation error $e(n)$ is small at the steady state and the terms that include
 183 high powers of $e(n)$ can be neglected resulting in

$$\mathbf{w}_{o,e2} = \mathbf{R}^{-1} \mathbf{p}$$

184 hence conforming with the Wiener solution. This is due to the fact that when
 185 $e(n)$ is very small, which is the case in the steady state, $J_{e2}(n)$ becomes ap-
 186 proximately quadratic. Eqn. (25) can be further analysed by expressing the
 187 desired response $d(n)$ and the adaptation error $e(n)$ as functions of the a priori
 188 error $e_a(n)$ as dictated by eqn. (22) and by

$$e(n) = e_a(n) + v(n) = \tilde{\mathbf{w}}_n^t \mathbf{x}_n + v(n) \quad (27)$$

189 Similarly to (25) the optimal coefficient vector in the least sum of exponentials
 190 sense $\mathbf{w}_{o,se}$ is derived by minimising the expectation of $J_{se}(n)$ as

$$\mathbf{w}_{o,se} = \left[\mathbf{R} \sum_{i=0}^{+\infty} \frac{1}{(2i+1)!} E \{ e^{2i}(n) \} \right]^{-1} \cdot \left[\mathbf{p} \sum_{i=0}^{+\infty} \frac{1}{(2i+1)!} E \{ e^{2i}(n) \} \right] \quad (28)$$

191 The least sum of exponentials solution from (28) differs from the least expo-
 192 nentiated square solution (eqn (26)) only in the weighting of the terms that
 193 contain high order powers of $e(n)$. Assuming that the measurement noise is
 194 negligible, the terms containing high order powers of the adaptation error can
 195 be ignored, resulting in

$$\mathbf{w}_{o,se} = \mathbf{R}^{-1} \mathbf{p}$$

196 that is the Wiener solution.

197 4.2 Step-size bounds for stability

198 Similar to all gradient descent algorithms, the choice of the step size in least
 199 exponential algorithms is crucial. To guarantee stability, the step size should
 200 satisfy

$$E \{ \|\tilde{\mathbf{w}}_{n+1}\|^2 \} \leq E \{ \|\tilde{\mathbf{w}}_n\|^2 \} \quad (29)$$

201 Embarking upon (29), and using (19) in order to preserve stability the bound
 202 on the step size can be calculated as

$$\mu \leq \frac{2}{[\|\mathbf{x}_n\|^4]^{1/2}} \left(\inf_{E\{e_a^2\} \in \Omega} \frac{E \{ e_a^2 \} h_G [E \{ e_a^2 \}]}{\sqrt{h_C [E \{ e_a^2 \}]} } \right) \quad (30)$$

203 where

$$h_G [E \{ e_a^2 \}] \triangleq \frac{E \{ e_a(i) f[e(i)] \}}{E \{ e_a^2 \}}, \quad (31)$$

$$h_C [E \{ e_a^2 \}] \triangleq E \{ f^4[e(i)] \} \quad (32)$$

204 and the set Ω'' is defined as

$$\Omega'' = \left\{ E \{ e_a^2 \} : \lambda \leq E \{ e_a^2 \} \leq \frac{1}{4} \text{Tr}(\mathbf{R}) E \{ \|\tilde{\mathbf{w}}_o\|^2 \} \right\} \quad (33)$$

205 with λ the Cramer-Rao bound [17]. Using (30) the set of values of the step
 206 size that guarantee stability can be computed by finding the values of the
 207 functions $h_G[E \{ e_a^2 \}]$ and $h_C[E \{ e_a^2 \}]$ for the LE2 and the LSE algorithms.
 208 Indeed, substituting the nonlinearity of the LE2 algorithm $f_{e2}[e(n)]$ from (13)
 209 into (31) yields

$$h_G[E \{ e_a^2(n) \}] = \frac{E \{ e_a(n) e(n) \exp[e^2(n)] \}}{E \{ e_a^2(n) \}}. \quad (34)$$

210 Replacing $\exp[e^2(n)]$ with its TSE, expressing $e(n)$ as a function of the a
 211 priori error according to (27), and assuming that $e_a(n)$ and $v(n)$ are mutually
 212 independent, results in

$$h_G[E \{ e_a^2(n) \}] = \sum_{i=0}^{+\infty} \frac{1}{i!} \frac{E \{ e_a^{2i+2}(n) \}}{E \{ e_a^2(n) \}} \quad (35)$$

213 In a similar manner the value of the function $h_C[e_a^2(n)]$ can be found for the
 214 LE2 algorithm by combining (13) and (32) that is

$$h_C[E \{ e_a^2(n) \}] = \sum_{i=0}^{+\infty} 2^{2i} \frac{1}{i!} \left[E \{ e_a^{2i+4}(n) \} + E \{ v^{2i+4}(n) \} \right] \quad (36)$$

215 Hence, in order to guarantee stability, the step size of the LE2 algorithm
 216 should be upper bounded by

$$0 \leq \mu \leq \frac{2}{[\|\mathbf{x}_n\|^4]^{1/2}} \left(\inf_{E \{ e_a^2 \} \in \Omega''} \frac{\sum_{i=0}^{+\infty} \frac{1}{i!} \frac{E \{ e_a^{2i+2}(n) \}}{E \{ e_a^2(n) \}}}{\sqrt{\sum_{i=0}^{+\infty} 2^{2i} \frac{1}{i!} [E \{ e_a^{2i+4}(n) \} + E \{ v^{2i+4}(n) \}]}} \right) \quad (37)$$

From (37), it is apparent that the step size depends strongly on the even moments of the measurement noise. It appears that the larger the amount of the injected noise the smaller the maximum step size is. Similar conditions for

the LSE algorithm can be derived as

$$h_G[E\{e_a^2(n)\}] = \sum_{i=1}^{+\infty} \frac{1}{(2i+1)!} \frac{E\{e_a^{2i+2}(n)\}}{E\{e_a^2(n)\}} \quad (38)$$

and

$$h_C[E\{e_a^2(n)\}] = 2 \sum_{i=1}^{+\infty} (4^{2i} - 4^{i+1}) \frac{1}{2i!} [E\{e_a^{2i+4}(n)\} + E\{v^{2i+4}(n)\}] \quad (39)$$

217 Therefore the step size of the LSE algorithm should satisfy

$$0 \leq \mu \leq \frac{2}{[\|\mathbf{x}_n\|^4]^{1/2}} \left(\inf_{E\{e_a^2\} \in \Omega} \frac{\sum_{i=1}^{+\infty} \frac{1}{(2i+1)!} \frac{E\{e_a^{2i+2}(n)\}}{E\{e_a^2(n)\}}}{\sqrt{\sum_{i=1}^{+\infty} (4^{2i} - 4^{i+1}) \frac{1}{(2i)!} [E\{e_a^{2i}(n)\} + E\{v^{2i}(n)\}]}} \right) \quad (40)$$

218 4.3 Step size normalisation

219 Using a constant step size in LE gradient descent algorithms is very restrictive;
 220 it should be very small in order for the algorithm to converge – especially in
 221 the presence of large modelling or measurement error $v(n)$ – according to (37)
 222 and (40). This does not allow for the full exploitation of the benefits of the
 223 exponential cost function.

224 To circumvent this problem, recall that minimisation of the *a posteriori* error
 225 during every iteration results in time varying normalised step sizes [18], given
 226 by (Appendix A)

$$\mu_{e2}(n) = \frac{\mu}{\mathbf{x}_n^t \mathbf{x}_n \exp[e^2(n)]} \quad (41)$$

227 and

$$\mu_{se}(n) = \frac{\mu e(n)}{\mathbf{x}_n^t \mathbf{x}_n (\exp[e(n)] - \exp[-e(n)])} \quad (42)$$

228 These normalised step sizes completely remove the exponential factor in the
 229 update of the LE algorithms, resulting in the standard normalised least mean

230 square (NLMS) algorithm. In order to control the steepness of the error sur-
 231 face, we shall introduce a positive multiplicative factor α in the exponential
 232 terms (41) and (42), which results in partially normalised step sizes given by

$$\mu_{e2}(n) = \frac{\mu}{\mathbf{x}_n^t \mathbf{x}_n \exp[\alpha e^2(n)]} \quad (43)$$

233 and

$$\mu_{se}(n) = \frac{\mu e(n)}{\mathbf{x}_n^t \mathbf{x}_n (\exp[\alpha e(n)] - \exp[-\alpha e(n)])} \quad (44)$$

234 The closer α to unity, the less pronounced the effect of the exponential term
 235 and the greater the similarity with the NLMS algorithm is. For $\alpha < 1$, algo-
 236 rithms that are faster, but less robust than the standard NLMS algorithm are
 237 derived. Having values of α greater than unity results in algorithms with slow
 238 response but increased robustness to impulsive noise.

239 4.4 Excess Mean Squared Error

240 The excess mean square error (EMSE) is defined as the expectation of the
 241 square value of the a priori error in the steady state, that is

$$S \triangleq \lim_{n \rightarrow \infty} E \{ |e_a(n)|^2 \} \quad (45)$$

242 This quantity measures the ability of the algorithm to model the desired signal;
 243 the lower the value of S the more accurate the modelling is. According to the
 244 energy preservation framework [13], EMSE is the fixed point of the equation

$$S = Tr(\mathbf{R}) \frac{h_U[S]}{h_G[S]} \quad (46)$$

245 where function $h_G[\cdot]$ is given by eqn (31) and $h_U[\cdot]$ is defined as the expectation
 246 of the square of the nonlinear error function,

$$h_U[E\{e_a^2(n)\}] = E \{ f^4[e(i)] \} \quad (47)$$

247 Taking into account the nonlinearity within the LE2 algorithm (13), and using
 248 the standard assumptions of the energy conservation framework yields

$$h_U[E\{e_a^2(n)\}] = \sum_{i=0}^{+\infty} \frac{1}{i!} 2^i [E\{e_a^{2i+2}(n)\} + E\{v^{2i+2}(n)\}] \quad (48)$$

249 Similarly, for the LSE algorithm, we have

$$h_U[E\{e_a^2(n)\}] = 2 \sum_{i=1}^{+\infty} \frac{1}{(2i)!} 2^{2i} [E\{e_a^{2i}(n)\} + E\{v^{2i}(n)\}] \quad (49)$$

250 Hence, the EMSE of the LE2 algorithm is the positive root of the nonlinear
251 equation

$$S = \frac{\mu}{2} Tr(\mathbf{R}) \frac{S + \sigma_U^2 + \sum_{i=1}^{+\infty} \frac{2^i}{i!} [E\{e_a^{2i+2}\} + E\{v^{2i+2}\}]}{1 + \sum_{i=1}^{+\infty} \frac{1}{i!} \frac{E\{e_a^{2i+2}\}}{S}} \quad (50)$$

252 Solving (50) for S yields

$$S = \frac{\beta \sigma_U^2 + \sum_{i=1}^{+\infty} \left[\frac{2^i \beta}{i!} E\{v^{2i+2}\} + \frac{2^{i\beta-1}}{i!} E\{e_a^{2i+2}\} \right]}{1 - \beta} \quad (51)$$

253 where $\beta = \frac{\mu}{2} Tr(\mathbf{R})$. The EMSE of the LSE algorithm is the positive root of
254 the nonlinear equation

$$S = \mu Tr(\mathbf{R}) \frac{S + \sigma_U^2 + \sum_{i=2}^{+\infty} \frac{2^i}{(2i)!} [E\{e_a^{2i}\} + E\{v^{2i}\}]}{1 + \sum_{i=1}^{+\infty} \frac{1}{(2i+1)!} \frac{E\{e_a^{2i+2}\}}{S}} \quad (52)$$

255 As a consequence the steady state error of the LSE algorithm is found to be

$$S = \frac{2\beta \sigma_U^2 + \sum_{i=1}^{+\infty} \left[\frac{2^{i+1}\beta}{(2i)!} E\{v^{2i+2}\} + \frac{2^{i+1}(2i+1)\beta-1}{(2i+1)!} E\{e_a^{2i+2}\} \right]}{1 - 2\beta} \quad (53)$$

256 This completes the analysis of the proposed class of algorithms. Obviously, the
257 study of the steady state behaviour of the proposed LE algorithms through
258 eqn. (51) and (53) is a demanding task that requires a thorough analysis.
259 However, the following straightforward conclusions can be drawn

- 260 • The steady state MSE of LE algorithms depends on the even high order
- 261 moments of the *a priori* error $e_a(n)$
- 262 • The EMSE is a function of the even high order moments of the additive
- 263 noise $v(n)$

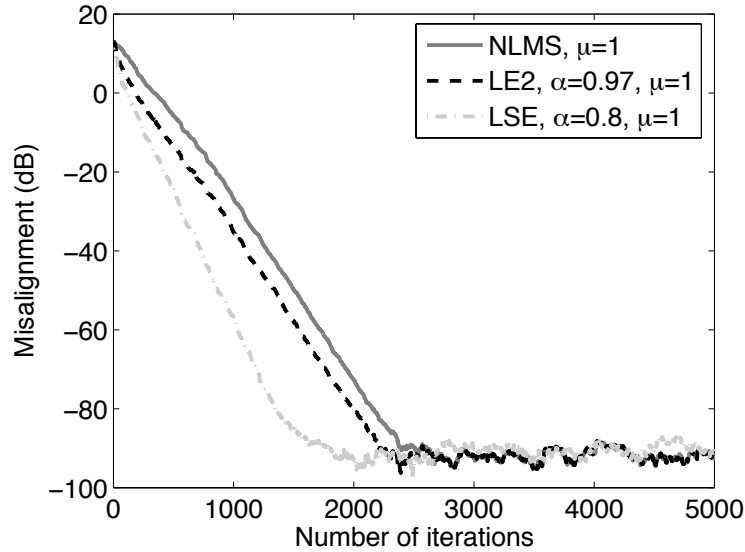


Fig. 3. Misalignment of the normalised LMS, the LE2 and the LSE algorithms for additive white measurement noise in a system identification context.

- The steady state MSE of the LSE is slightly lower than that of the LE2.

5 Simulation Results

The performance of the proposed LSE and LE2 algorithms was evaluated in a system identification setting. Both the unknown channel and the identifying filter were FIR filters of the same order. The input signal was coloured noise, that was produced by passing a white noise signal with Gaussian distribution $\mathcal{N} \sim (0, 1)$ through an autoregressive model with transfer function

$$A(z) = \frac{1}{1 - 1.79z^{-1} + 1.85z^{-2} - 1.27z^{-3} + 0.41z^{-4}} \quad (54)$$

The quantitative performance measure was the misalignment defined as

$$\|\tilde{\mathbf{w}}_n\|^2 = (\mathbf{w}_o - \mathbf{w}_n)^t (\mathbf{w}_o - \mathbf{w}_n) \quad (55)$$

where $\mathbf{w}_n = [w_0(n), w_1(n), \dots, w_{N-1}(n)]^t$ the values of the digital filter coefficients at time instant n and $\mathbf{w}_o = [w_{0o}, w_{1o}, \dots, w_{(N-1)o}]^t$ the samples of the impulse response of the linear time invariant (LTI) unknown channel. The performance of the NLMS algorithm whose step size is given by

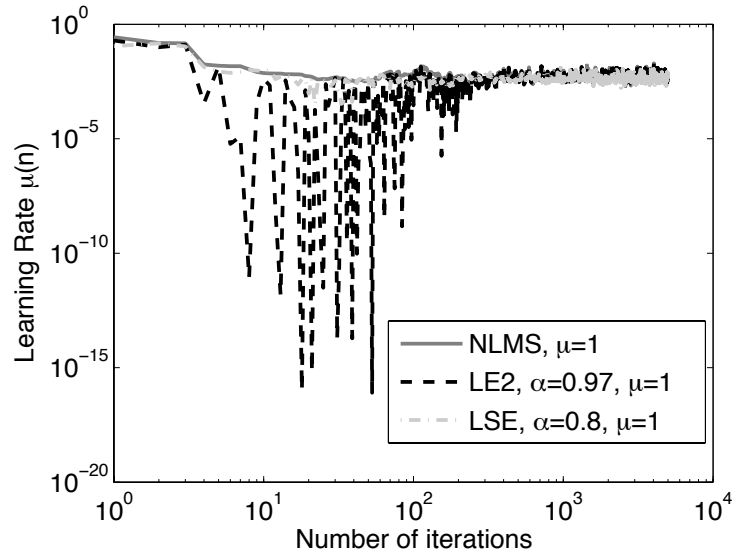


Fig. 4. The time varying step sizes for the normalised LMS and the normalised least exponential squared error algorithms, for the misalignment curves from Fig. 3

$$\mu_{NLMS}(n) = \frac{\mu}{\|\mathbf{x}_n\|^2} \quad (56)$$

276 was used as benchmark. A set of simulations for partially normalised learning
277 rates within the LE2 and the LSE algorithms is also provided.

278 In Fig. 3 the misalignment curves of the LE2, the LSE and the NLMS algo-
279 rithms are shown when the system output is contaminated with additive white
280 Gaussian noise $v(n)$ with 80dB SNR. The step size of the LE2 was varying
281 according to (43) with $\mu = 1$ and $\alpha = 0.97$, while that of the LSE was calcu-
282 lated from (44) with $\mu = 1$ and $\alpha = 0.8$. Both LE algorithms outperformed
283 the NLMS, since they reached the same steady state error level within fewer
284 iterations.

285 Assuming that the error is the steady state is negligible ($\lim_{n \rightarrow \infty} e(n) = 0$)
286 results in

$$\exp[\alpha e^2(n)] \approx 1 \quad (57)$$

287 and

$$(\exp[\alpha e(n)] - \exp[-\alpha e(n)]) \approx \alpha. \quad (58)$$

288 Substituting (57) and (58) into (43) and (44) respectively yields

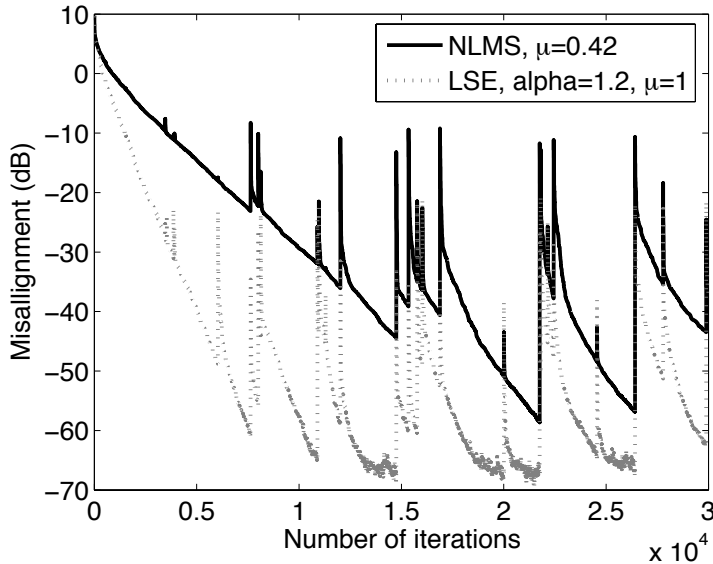


Fig. 5. Misalignment curves of the normalised LMS and the normalised LSE algorithms for impulsive noise disturbances.

$$\mu_{e2}(n) \approx \frac{\mu}{\|\mathbf{x}_n\|^2} = \mu_{NLMS}(n) \quad (59)$$

289 and

$$\mu_{se}(n) \approx \frac{\mu}{\alpha \|\mathbf{x}_n\|^2} = \frac{\mu_{NLMS}(n)}{\alpha} \quad (60)$$

290 This is also observed from the step size trajectories, depicted in Fig. 4, where
291 all the evaluated SGD algorithms have the same step size at the steady state.

292 The performance of the LSE algorithm when the output of the unknown system
293 is contaminated with impulsive noise along with 80 dB SNR of white
294 Gaussian noise, is presented in Fig. 5. The step size was varied according to
295 (44), with $\alpha = 1.2$ and $\mu = 1$. The misalignment curve for the NLMS algorithm
296 with $\mu = 0.42$, is also provided. Observe that the LSE algorithm is more
297 immune to impulsive noise than the NLMS, since it has faster convergence,
298 lower steady state error, and less pronounced overshoot every time an impulse
299 occurs. The learning rates of the LSE and the NLMS algorithms were chosen
300 so as to have similar values at the steady state. This is presented in Fig. 6,
301 where it is shown that both the LSE and the NLMS algorithms, have similar
302 learning rate values. As desired, when an impulse occurs, the learning rate of
303 the LSE algorithm becomes very small, preventing erroneous updating of the
304 values of the filter coefficients.

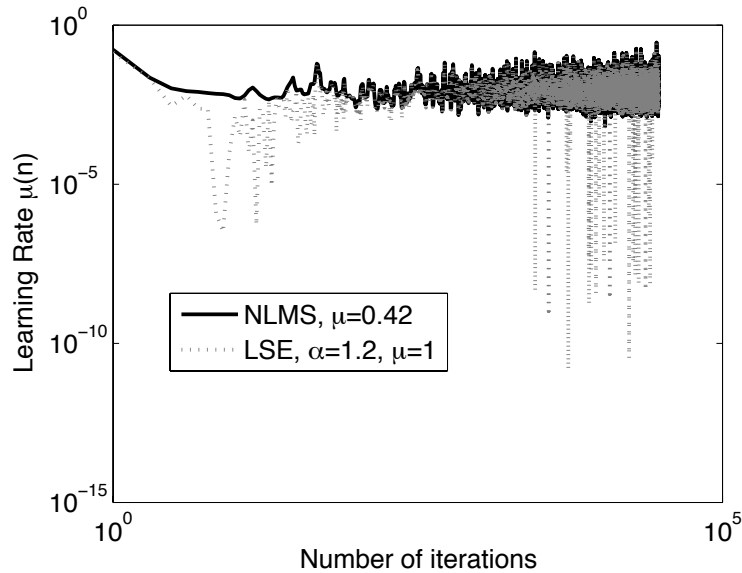


Fig. 6. The corresponding step sizes for the misalignment curves of Fig. 5.

305 Notice that the parameter α of the partially normalised learning rate is very
 306 important, since it greatly affects the performance of the LSE (or the LE2)
 307 algorithm for the same μ . For $\alpha \in (0, 1)$ algorithms that converge faster than
 308 the conventional NLMS are obtained, but are very sensitive. Having $\alpha > 1$,
 309 on the other hand, results in algorithms that are more robust under impulsive
 310 noise than the NLMS but have slower convergence.

311 6 Conclusions

312 A novel class of least exponential (LE) algorithms has been presented . These
 313 have been derived by minimising cost functions that have exponential depen-
 314 dence on the adaptation error. It has been shown that LE algorithms can be
 315 considered as a generalisation of the mixed norm stochastic gradient descent
 316 algorithms since they take into account an infinite number of error norms. A
 317 rigorous mathematical analysis has been provided resulting in closed form
 318 expressions for the optimal solutions and an upper bounds for the learning
 319 rate. For robustness, normalisation of the step size of the proposed algorithms
 320 has been addressed. Simulation results in a system identification setting and
 321 under various noise conditions support the analysis.

322 **A Normalisation of the step size**

323 Normalisation of the step size can be achieved through the minimisation of
 324 the magnitude of the *a posteriori* error given by

$$e_p(n) = d(n) - \mathbf{w}_{n+1}^t \mathbf{x}_n \quad (\text{A.1})$$

325 as was presented in [18] for the case of the LMS algorithm. Applying similar
 326 considerations, normalised step sizes for the LE algorithms can be derived.
 327 Indeed, substituting (12) in (A.1), yields

$$e_p(n) = d(n) - \left[\mathbf{w}_n + \mu_{e2} e(n) \mathbf{x}_n \exp[e^2(n)] \right]^t \mathbf{x}_n \quad (\text{A.2})$$

328 Using (22), (A.2) can be re-written as

$$\varepsilon_p(n) = \left[1 - \mu_{e2} \mathbf{x}_n^t \mathbf{x}_n \exp[e^2(n)] \right] e(n) \quad (\text{A.3})$$

329 The magnitude of the *a posteriori* error is minimised when a time varying step
 330 size is employed, that is

$$\mu_{e2}(n) = \frac{1}{\mathbf{x}_n^t \mathbf{x}_n \exp[e^2(n)]} \quad (\text{A.4})$$

331 Thus an the optimal learning rate of the LE2 algorithm becomes

$$\mu_{e2}(n) = \frac{\mu_{e2}}{\mathbf{x}_n^t \mathbf{x}_n \exp[e^2(n)]} \quad (\text{A.5})$$

332 where $0 \leq \mu_{e2} \leq 2$. Similarly, an appropriate choice for the step size for the
 333 LSE algorithm is

$$\mu_{se}(n) = \frac{\mu_{se} e(n)}{\mathbf{x}_n^t \mathbf{x}_n (\exp[e(n)] - \exp[-e(n)])} \quad (\text{A.6})$$

334 These normalised step sizes completely remove the exponential terms from
 335 the recursive equations of the LE2 and the LSE algorithms, given by (12)
 336 and (15) respectively, and reduce the derived LE algorithms to the standard
 337 NLMS algorithm. Introducing a positive factor α such that

$$\mu_{e2}(n) = \frac{\mu_{e2}}{\mathbf{x}_n^t \mathbf{x}_n \exp[\alpha e^2(n)]} \quad (\text{A.7})$$

338 and

$$\mu_{se}(n) = \frac{\mu_{se} e(n)}{\mathbf{x}_n^t \mathbf{x}_n (\exp[\alpha e(n)] - \exp[-\alpha e(n)])} \quad (\text{A.8})$$

339 the effect of these exponential term can be controlled.

340 References

- 341 [1] A. Benveniste, M. Metivier, P. Priouret, Adaptive Algorithms and Stochastic
342 Approximations, New York:Springer-Verlang, 1990.
- 343 [2] B. Widrow, S. Stearns, Adaptive Signal Processing, Prentice Hall, 1985.
- 344 [3] N. Bershad, Analysis of the Normalized LMS with Gaussian Inputs, IEEE
345 Trans. Acoust., Speech, Signal Processing 34 (4) (1986) 793–806.
- 346 [4] E. Walach, B. Widrow, The least mean fourth (LMF) adaptive algorithm and its
347 family, IEEE Trans. Inform. Theory 30 (2) (2003) 275–283.
- 348 [5] O. Tanrikulu, A. Constantinides, Least mean kurtosis: A novel higher-order
349 statistics based adaptive filtering algorithm, Elec. Lett 30 (3) (1994) 189–190.
- 350 [6] J. Chambers, O. Tanrikulu, A. Constantinides, Least mean mixed-norm
351 adaptive filtering, Elec. Lett., vol. 30, no 3, pp. 1–190, 1994.
- 352 [7] J. Chambers, A. Avlonitis, A Robust Mixed-Norm Adaptive Filter Algorithm,
353 IEEE Signal Processing Lett., vol. 4, no 2, pp. 46–48, 1997.
- 354 [8] D. Mandic, E. Papoulis, C. Boukis, A Normalized Mixed-Norm Adaptive
355 Filtering Algorithm Robust Under Impulsive Noise Interference, in: IEEE
356 International Conference on Acoustics, Speech, and Signal Processing, vol. VI,
357 pp. 333–336, 2003.
- 358 [9] A. Barros, J. Principe, Y. Takeuchi, C. Sales, N. Ohnishi, An algorithm based
359 in the even moments of the error, in: IEEE XIII Workshop on Neural Networks
360 for Signal Processing, pp. 879–885, 2003.
- 361 [10] J. Kivinen, M. Warmuth, Exponentiated gradient versus gradient descent for
362 linear predictors, Inform. Comp., vol. 132, no 1, pp. 1–64, 1997.
- 363 [11] R. Martin, W. Sethares, R. Williamson, C. Johnson Jr., Exploiting sparsity in
364 adaptive filters, IEEE Trans. Signal Processing, vol. 50, no 8, pp. 1883–1893,
365 2002.

- 366 [12] J. Benesty, Y. Huang, J. Chen, An exponentiated gradient adaptive algorithm
367 for blind identification of sparse simo systems, in: IEEE International
368 Conference on Acoustics, Speech, and Signal Processing, vol. 2, pp. 829–832,
369 2002.
- 370 [13] M. Rupp, A.H. Sayed, A Time-Domain Feedback Analysis of Filtered-Error
371 Adaptive Gradient Algorithms, IEEE Trans. Signal Processing, vol. 44, no 6,
372 pp. 1428–1429, 1996.
- 373 [14] T.Y. Al-Naffouri, A.H. Sayed, Adaptive Filters with Error Nonlinearities: Mean-
374 Square Analysis and Optimum Design, EURASIP Journal of Applied Signal
375 Processing, vol. 4, pp. 192–205, 2001.
- 376 [15] T.Y. Al-Naffouri, A.H. Sayed, Transient Analysis of Adaptive Filters with Error
377 Nonlinearities, IEEE Trans. Signal Processing, vol. 51, no 3, pp. 653–661, 2003.
- 378 [16] N.R. Yousef, A.H. Sayed, A Unified Approach to the Steady-State and Tracking
379 Analyses of Adaptive Filters, IEEE Trans. Signal Processing, vol. 49, no 2,
380 pp. 314–324, 2001.
- 381 [17] L. Van Trees, Detection, Estimation and Modulation Theory: Part I, John Wiley
382 and Sons, New York, 1968.
- 383 [18] E. Soria-Olivas, J. Calpe-Maravilla, J. Guerrero-Martinez, M. Martinez-
384 Sober, J. Espi-Lopez, An Easy Demonstration of the Optimum Value of the
385 Adaptation Constant in the LMS Algorithm, IEEE Transactions on Education,
386 vol. 41, no 1, pp. 81, 1998.