

Improved l_0 -RLS adaptive filter

B.K. Das and M. Chakraborty[✉]

In this Letter, the authors present an improved sparse recursive least squares (RLS) algorithm, which employs a novel approximation of the l_0 norm of the filter coefficient vector for regularising the RLS cost function. The proposed algorithm achieves improved performance over existing algorithms as demonstrated via numerical simulations.

Introduction: It is well known that the most ideal way of deriving sparsity aware adaptive filters is to add the l_0 norm penalty of the filter coefficient vector $\mathbf{w}(n)$ to the cost function of standard adaptive filters (e.g. least mean square (LMS), recursive least squares (RLS) etc.), and minimise the sum by differentiating it w.r.t. $\mathbf{w}(n)$ and setting it to zero. However, the l_0 norm, being just a count of non-zero entries of a vector, is not differentiable (neither it is a true norm) and thus, the above optimisation requires combinatorial search that is NP hard. A popular approach to overcome this difficulty has been to approximate the l_0 norm by other norms like the l_1 norm $\|\mathbf{w}(n)\|_1$ which are differentiable almost everywhere (a.e.). Recently, it has also been shown [1, 2] that the sparsity promoting function $\sum_{i=0}^{N-1} [1 - \exp(-\alpha|w_i(n)|)]$ is a better approximation of the l_0 norm of $\mathbf{w}(n)$ than the l_1 norm (where $w_i(n)$ denotes the i th coefficient of $\mathbf{w}(n)_{N \times 1}$ and α is a suitable constant). An RLS-based adaptive filter has been derived in [1] that adds the above function as a sparsity promoting penalty to the RLS cost function. Unfortunately, first-order differentiation of this function w.r.t. $\mathbf{w}(n)$ results in expressions given in terms of $\text{sgn}[\mathbf{w}(n)]$ (where the operator $\text{sgn}[\cdot]$ denotes element wise ‘sign’ or ‘signum’ function), meaning minimisation of the cost function leads to equations that are non-linear. Due to this, a closed form least squares (LS) solution of $\mathbf{w}(n)$ and its exact RLS update become very hard to determine. To circumvent this, the term $\text{sgn}[\mathbf{w}(n)]$ is replaced by $\text{sgn}[\mathbf{w}(n-1)]$ in [1], which, however, results in the introduction of substantial gradient noise especially for ‘inactive’ taps (i.e. the taps having zero or close to zero optimum values), as, in the steady state, estimates of such taps fluctuate around zero. In order to counter this, we follow the philosophy of [3] and propose a new approximation of the l_0 norm that results in exact LS solution of $\mathbf{w}(n)$ and its RLS update without giving rise to terms like $\text{sgn}[\mathbf{w}(n)]$.

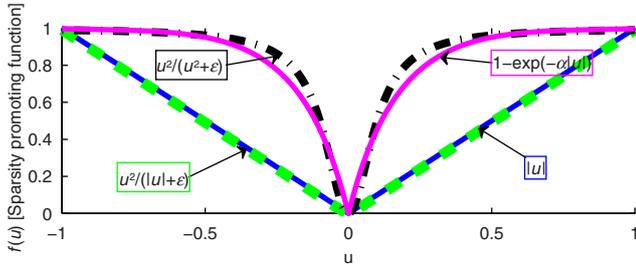


Fig. 1 Some sparsity promoting functions

New non-convex approximation of l_0 norm: Given a vector $\mathbf{a} \in \mathbb{R}^N$, we propose a new approximation of its l_0 norm by a non-convex, a.e. differentiable function $\sum_{i=0}^{N-1} (a_i^2)/(a_i^2 + \epsilon)$, where ϵ is a small, positive constant. In Fig. 1, we plot the univariate version of this with ϵ taken as 1×10^{-4} and also of the aforesaid l_0 norm approximation $\sum_{i=0}^{N-1} [1 - \exp(-\alpha|a_i|)]$ with $\alpha = 50$, by dashed ‘black’ and ‘magenta’, respectively, as a function of a single variable u . It is seen that both curves are very close to each other, and both are much better approximations of the l_0 norm than the l_1 norm (i.e. $|u|$, as shown in ‘blue’). Using the proposed approximation, the l_0 norm regularised LS cost function is given as

$$J(\mathbf{w}(n)) = \sum_{k=1}^n \lambda^{n-k} [d(k) - \mathbf{x}^T(k)\mathbf{w}(n)]^2 + \frac{\rho}{2} \sum_{i=0}^{N-1} \frac{w_i^2(n)}{w_i^2(n) + \epsilon} \quad (1)$$

$$\simeq \sum_{k=1}^n \lambda^{n-k} [d(k) - \mathbf{x}^T(k)\mathbf{w}(n)]^2 + \frac{\rho}{2} \sum_{i=0}^{N-1} \frac{w_i^2(n)}{w_i^2(n-1) + \epsilon},$$

where $d(n) = \mathbf{x}^T(n)\mathbf{w}_0 + v(n)$ is the observable system output, $\mathbf{x}(n) = [x(n)x(n-1) \dots x(n-N+1)]^T$, $x(n)$ is the system input, $v(n)$

is an additive observation noise, and \mathbf{w}_0 is the system impulse response vector. Note that in the denominator of the second term in (1), $w_i(n)$ is approximated by $w_i(n-1)$ which is treated as a known constant. It is easy to check that as long as $w_i(n)$ and $w_i(n-1)$ remain close to each other, $\sum_{i=0}^{N-1} w_i^2(n)/(w_i^2(n-1) + \epsilon)$ still remains a good approximation of $\|\mathbf{w}(n)\|_0$. The above approximation, however, ensures that on derivation by $w_i(n)$, one obtains a first-order expression in $w_i(n)$, similar to the first term on the RHS of (1), which is the standard LS cost function and, whose derivative w.r.t. $\mathbf{w}(n)$ is given by

$$(\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{X}(n))\mathbf{w}(n) - \mathbf{X}^T(n)\mathbf{A}(n)\mathbf{d}(n),$$

where $\mathbf{X}(n) = [\mathbf{X}^T(n-1) \quad \mathbf{x}(n)]^T$, $\mathbf{X}(1) = \mathbf{x}^T(1) = [x(1), 0, 0, \dots, 0]$, $\mathbf{d}(n) = [d(1), d(2), \dots, d(n)]^T$ and $\mathbf{A}(n) = \text{diag}\{\lambda^{n-1}, \dots, \lambda^2, \lambda, 1\}$ with λ ($0 < \lambda < 1$) being a forgetting factor.

The second term on the RHS of (1), which is the sparsity promoting term, on the other hand, produces $\rho\mathbf{D}_0(n)\mathbf{w}(n)$ when differentiated w.r.t. $\mathbf{w}(n)$, where $\mathbf{D}_0(n) = \text{diag}\{d_{0,0}(n), d_{0,1}(n), \dots, d_{0,N-1}(n)\}$, and $d_{0,i}(n) = 1/(w_i^2(n-1) + \epsilon)$. Combining, the optimal LS solution $\mathbf{w}(n)$ for (1) is given by

$$\mathbf{w}(n) = \arg \min_{\mathbf{w}(n)} J(\mathbf{w}(n)) = \mathbf{P}(n)\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{d}(n), \quad (2)$$

where $\mathbf{P}(n) = (\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{X}(n) + \rho\mathbf{D}_0(n))^{-1}$.

To obtain the RLS update relation from (2), we follow the approach of [3]. In [3], an l_1 norm regularised RLS algorithm (termed as ZA-RLS-II) was derived which approximates the l_1 norm of a given $\mathbf{a} \in \mathbb{R}^N$ by $\sum_{i=0}^{N-1} (a_i^2)/(a_i^2 + \epsilon)$ (shown by the dashed ‘green’ curve in Fig. 1 for the univariate case, for $\epsilon = 1 \times 10^{-4}$). For this, we first note that

$$\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{d}(n) = \lambda\mathbf{X}^T(n-1)\mathbf{A}(n-1)\mathbf{d}(n-1) + \mathbf{x}(n)d(n). \quad (3)$$

Next, we write

$$\begin{aligned} \mathbf{P}^{-1}(n) &= \mathbf{X}^T(n)\mathbf{A}(n)\mathbf{X}(n) + \rho\mathbf{D}_0(n) \\ &= \lambda\mathbf{X}^T(n-1)\mathbf{A}(n-1)\mathbf{X}(n-1) + \mathbf{x}(n)\mathbf{x}^T(n) + \rho\mathbf{D}_0(n) \\ &= \lambda\mathbf{P}^{-1}(n-1) + \mathbf{x}(n)\mathbf{x}^T(n) + \rho(\mathbf{D}_0(n) - \lambda\mathbf{D}_0(n-1)). \end{aligned}$$

The above implies

$$\begin{aligned} \lambda\mathbf{P}(n)\mathbf{P}^{-1}(n-1) &= \mathbf{I} - \mathbf{P}(n)\mathbf{x}(n)\mathbf{x}^T(n) \\ &\quad - \rho\mathbf{P}(n)(\mathbf{D}_0(n) - \lambda\mathbf{D}_0(n-1)) \end{aligned} \quad (4)$$

From (3) and (4), the time-recursive version of (2) is then obtained as

$$\begin{aligned} \mathbf{w}(n) &= \mathbf{P}(n)\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{d}(n) \\ &= \mathbf{P}(n)[\lambda\mathbf{X}^T(n-1)\mathbf{A}(n-1)\mathbf{d}(n-1) + \mathbf{x}(n)d(n)] \\ &= \mathbf{P}(n)[\lambda\mathbf{P}^{-1}(n-1)\mathbf{w}(n-1) + \mathbf{x}(n)d(n)] \\ &= [\mathbf{I} - \mathbf{P}(n)\mathbf{x}(n)\mathbf{x}^T(n)]\mathbf{w}(n-1) \\ &\quad - \rho\mathbf{P}(n)(\mathbf{D}_0(n) - \lambda\mathbf{D}_0(n-1))\mathbf{w}(n-1) + \mathbf{P}(n)\mathbf{x}(n)d(n) \\ &= [\mathbf{I} - \rho\mathbf{P}(n)(\mathbf{D}_0(n) - \lambda\mathbf{D}_0(n-1))]\mathbf{w}(n-1) \\ &\quad + \mathbf{P}(n)\mathbf{x}(n)e_p(n), \end{aligned} \quad (5)$$

where $e_p(n) = d(n) - \mathbf{w}^T(n-1)\mathbf{x}(n)$ is the so-called a priori error.

Next we try to evaluate $\mathbf{P}(n)$, which we write as $[\mathbf{Q}^{-1}(n) + \rho\mathbf{D}_0(n)]^{-1}$, where $\mathbf{Q}(n) = (\mathbf{X}^T(n)\mathbf{A}(n)\mathbf{X}(n))^{-1}$. To compute the inverse, we invoke the Woodbury matrix inversion lemma [4] which states that given \mathbf{A} and \mathbf{C} to be two $l \times l$ and $r \times r$ invertible matrices, and \mathbf{U} and \mathbf{V} be two $l \times r$ and $r \times l$ matrices, respectively, the following holds:

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}. \quad (6)$$

Taking $\mathbf{A} = \rho\mathbf{D}_0(n)$, $\mathbf{C} = \mathbf{Q}^{-1}(n)$, and \mathbf{U} and \mathbf{V} both to be identity matrices, we have

$$\begin{aligned} \mathbf{P}(n) &= (\mathbf{Q}^{-1}(n) + \rho\mathbf{D}_0(n))^{-1} \\ &= \frac{1}{\rho}\mathbf{D}_0^{-1}(n) - \frac{1}{\rho}\mathbf{D}_0^{-1}(n)[\mathbf{Q}(n) + \frac{1}{\rho}\mathbf{D}_0^{-1}(n)]^{-1}\frac{1}{\rho}\mathbf{D}_0^{-1}(n). \end{aligned}$$

Simplifying

$$\mathbf{P}(n) = [\mathbf{I} - \mathbf{D}_{0,\text{inv}}(n)\{\mathbf{Q}(n) + \mathbf{D}_{0,\text{inv}}(n)\}^{-1}]\mathbf{D}_{0,\text{inv}}(n), \quad (7)$$

where $\mathbf{D}_{0,\text{inv}}(n) = (\rho\mathbf{D}_0(n))^{-1} = (1/\rho)\mathbf{D}_0^{-1}(n)$.

It is seen from (7) that we still need a matrix inversion, i.e. $\{\mathbf{Q}(n) + \mathbf{D}_{0,\text{inv}}(n)\}^{-1}$ to evaluate $\mathbf{P}(n)$. Similar to [3], in the sequel, we show that this matrix inversion can be performed alternatively using an iterative update rule exploiting the structures of the individual matrices. First recall that $\mathbf{D}_{0,\text{inv}}(n)$ is a diagonal matrix with the i th diagonal element given by $(1/\rho)[w_i^2(n-1) + \epsilon]$. For large n , it is reasonable to expect that $\mathbf{D}_{0,\text{inv}}(n)$ will have approximately $r = \|\mathbf{w}_0\|_0$ number of significant diagonal entries. Using a threshold δ (which is a very small positive constant), we retain only those diagonal elements $[\mathbf{D}_{0,\text{inv}}(n)]_{i,i}$ for which $|w_i(n)| > \delta$. Define the index set $S = \{0, 1, \dots, N-1\}$ and the subset $NZ \subset S$ such that for every $i \in NZ$, $|w_i(n)| > \delta$ and for every $j \in S \setminus NZ$, $|w_j(n)| \leq \delta$. Let $|NZ| = l$, which, for large n , will be close to $r = \|\mathbf{w}_0\|_0$. We arrange the indices belonging to NZ in a $l \times 1$ vector $\boldsymbol{\theta}_n$, i.e. for every $i \in \{1, 2, \dots, l\}$, $\theta_n(i) \in NZ$ and if $i \neq j$, $\theta_n(i) \neq \theta_n(j)$. Then $\mathbf{D}_{0,\text{inv}}(n)$ can be approximated as $\mathbf{D}_{0,\text{inv}}(n) \simeq \sum_{i=1}^l \mathbf{D}_{0,\text{inv}}^i(n)$, where $\mathbf{D}_{0,\text{inv}}^i(n)$ is a rank one diagonal matrix with a non-zero entry at the $\theta_n(i)$ th diagonal position. Defining $\mathbf{Q}_j(n) = \mathbf{Q}(n) + \sum_{i=1}^j \mathbf{D}_{0,\text{inv}}^i(n)$ for $j = 1, \dots, l$, and also, $\mathbf{Q}_0(n) = \mathbf{Q}(n)$, the task then boils down to computing the inverse of $\mathbf{Q}_j(n)$. We do this job recursively. Suppose, the inverse of $\mathbf{Q}_j(n)$ has been computed for some $j \in \{0, 1, \dots, l-1\}$. Then, to compute the inverse of $\mathbf{Q}_{j+1}(n) = \mathbf{Q}_j(n) + \mathbf{D}_{0,\text{inv}}^{j+1}(n)$, we again use the aforesaid matrix inversion lemma [4] for which we choose $\mathbf{A} = \mathbf{Q}_j(n)$, $\mathbf{C} = 1$, \mathbf{u} : an $N \times 1$ vector with all elements zero except having $+\sqrt{(1/\rho)(w_{\theta_n(j+1)}^2(n-1) + \epsilon)} \simeq \sqrt{(1/\rho)}w_{\theta_n(j+1)}(n-1)$ at the $\theta_n(j+1)$ th position and $\mathbf{v} = \mathbf{u}^T$. (Here, the scalar C and the vectors \mathbf{u} and \mathbf{v} are equivalent to the matrices \mathbf{C} , \mathbf{U} , and \mathbf{V} , respectively, of (6).)

Then, we obtain from the aforementioned matrix inversion lemma

$$\begin{aligned} \mathbf{Q}_{j+1}^{-1}(n) &= [\mathbf{I} - \mathbf{Q}_j^{-1}(n)\mathbf{u}[1 + \mathbf{u}^T\mathbf{Q}_j^{-1}(n)\mathbf{u}]^{-1}\mathbf{u}^T]\mathbf{Q}_j^{-1}(n) \\ &= [\mathbf{I} - \frac{1}{1 + \text{Tr}(\mathbf{D}_{0,\text{inv}}^{j+1}(n)\mathbf{Q}_j^{-1}(n))}\mathbf{Q}_j^{-1}(n)\mathbf{D}_{0,\text{inv}}^{j+1}(n)]\mathbf{Q}_j^{-1}(n) \\ &= \mathbf{Q}_j^{-1}(n) - \frac{w_{\theta_n(j+1)}^2(n-1)}{\rho + w_{\theta_n(j+1)}^2(n-1)q_{j+1}(n)}\mathbf{q}_{j+1}(n)\mathbf{q}_{j+1}^T(n), \end{aligned}$$

for $j = 0, 1, 2, \dots, l-1$, where $\mathbf{q}_{j+1}(n)$ and $q_{j+1}(n)$ are the $\theta_n(j+1)$ th column and $\theta_n(j+1)$ th diagonal element of the matrix $\mathbf{Q}_j^{-1}(n)$, respectively.

For highly sparse \mathbf{w}_0 , $l \ll N$. Thus, the number of iterations required to compute $\mathbf{Q}_l^{-1}(n)$ can be kept considerably small.

Finally, $\mathbf{Q}^{-1}(n)$ is updated time recursively as

$$\mathbf{Q}^{-1}(n) = \lambda\mathbf{Q}^{-1}(n-1) + \mathbf{x}(n)\mathbf{x}^T(n). \quad (8)$$

Numerical simulations: We conduct an adaptive system identification experiment using zero mean, unit variance white Gaussian random process as input and a sparse \mathbf{w}_0 as the system to be identified. The system \mathbf{w}_0 has 128 coefficients, out of which four coefficients with random locations have value 1 each, while rest of the coefficients have zero value. The additive measurement noise is taken to be a zero mean white Gaussian noise with variance 0.01, giving rise to an SNR of 20 dB. Since ZA-RLS-II [3] is known to outperform the other contemporary sparse RLS algorithms (e.g. OCCD-TNWL [5], SPARLS [6] etc.), and behaves identical to ZA-RLS-I [3], we compare the performance of the proposed algorithm with that of ZA-RLS-II [3]. Apart from ZA-RLS-II, the l_0 norm regularised RLS [1], the standard sparsity unaware RLS, and the l_0 -normalised least mean square (NLMS) [2] are also deployed for identifying \mathbf{w}_0 , and their convergence behaviours are compared with that of the proposed one. The forgetting factor for all the

RLS-based filters is chosen as $\lambda = 0.95$. The other parameters for the ZA-RLS-II and the proposed algorithm are chosen as follows: $\rho = 0.0005$, $\epsilon = 0.0001$. The parameters for l_0 -RLS [1] are taken to be $\beta = 50$, $\kappa = 2$. In the case of l_0 -NLMS, parameters are tuned to obtain minimum steady-state mean-square deviation (MSD). We carry out the experiment for 375 iterations and the results are averaged over 100 ensembles. Fig. 2 shows the learning curves (i.e. MSD vs. iteration index) for (i) the proposed algorithm, (ii) ZA-RLS-II [3], (iii) the sparsity unaware RLS, (iv) the l_0 -RLS [1], and (v) the l_0 -NLMS [2] (plotted by dashed 'black', dashed 'green', 'blue', 'magenta', and 'red', respectively). It is seen that the proposed algorithm outperforms all the existing algorithms including the recently proposed ZA-RLS-II [3] (dashed 'green') in terms of steady-state MSD. The convergence speed of the proposed algorithm is also substantially higher than that of the l_0 -RLS [1], the l_0 -NLMS [2], and the RLS, and is at par with that of the ZA-RLS-II [3].

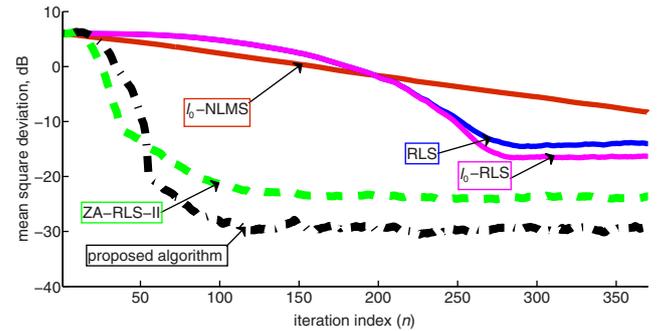


Fig. 2 MSD versus iteration index n for different sparse adaptive filters

Conclusion: In this Letter, we have proposed a new non-convex and a.e. differentiable approximation of the sparsity promoting l_0 norm which can be expressed as a weighted l_2 norm, and have employed it to derive a new RLS-type algorithm which is significantly superior to the existing sparse RLS algorithms including the newly proposed ZA-RLS-II [3] which was derived using a similar weighted l_2 norm-based approximation of the convex l_1 norm. Numerical simulations demonstrated that our proposed algorithm outperforms the existing algorithms including that of [3].

© The Institution of Engineering and Technology 2017

Submitted: 11 September 2017 E-first: 7 November 2017

doi: 10.1049/el.2017.3441

One or more of the Figures in this Letter are available in colour online.

B.K. Das (Department of E&ECE, IIT Kharagpur, West Bengal 721302, India)

References

- 1 Eksiöglu, E.M., and Tanc, A.K.: 'RLS algorithm with convex regularization', *Signal Process. Lett.*, 2011, **18**, (8), pp. 470–473
- 2 Su, G., Jin, J., Gu, Y., et al.: 'Performance analysis of l_0 norm constraint least mean square algorithm', *Trans. Signal Proc.*, 2012, **60**, (5), pp. 2223–2235
- 3 Hong, X., Gao, J., and Chen, S.: 'Zero attracting recursive least squares algorithms', *Trans. Veh. Technol.*, 2017, **66**, (1), pp. 213–221
- 4 Horn, R., and Johnson, C.R.: 'Matrix analysis' (Cambridge University, Cambridge, 1985)
- 5 Angelosante, D., Bazerque, J. A., and Giannakis, G. B.: 'Online adaptive estimation of sparse signals: where RLS meets the l_1 -norm', *Trans. Signal Proc.*, 2010, **58**, (7), pp. 3436–3447
- 6 Babadi, B., Kalouptsidis, N., and Tarokh, V.: 'SPARLS: the sparse RLS algorithm', *Trans. Signal Proc.*, 2010, **58**, (8), pp. 4013–4025