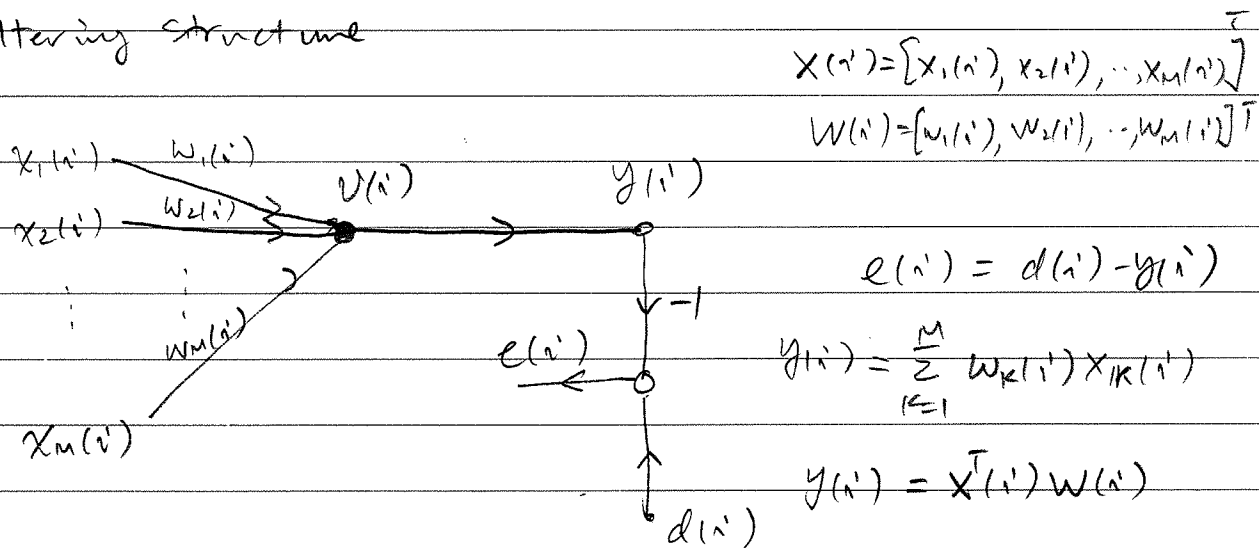


## The Least-Mean-Square Algorithm

- the 1<sup>st</sup> linear adaptive filtering algorithm
- applications include prediction, communication-channel equalization, ... system ID
- inspired by perceptron, linear combiner
- "linear" complexity yet effective in practice
- simple to code
- robust to disturbance
- set the stage for backpropagation

Filtering structure



where  $i = 1, 2, \dots, n$  denotes time index

sample pairs  $\mathcal{T} = \{x(i), d(i); i = 1, \dots, n, \dots\}$

①  $X(i)$  can be a snapshot of data

②  $X(i)$  can be past values of  $X(i-1)$ , ...

e.g.  $y(i) = w_1(i) x_1(i) + w_2(i) x_1(i-1) + \dots + w_m(i) x_1(i-(m-1))$

The Least-Mean Square Algorithm.

Introduce the instantaneous cost function

$$E(\hat{W}) = \frac{1}{2} e^2(n)$$

for a linear neuron,

$$e(n) = d(n) - X^T(n) \hat{W}(n)$$

where  $X(n)$  - input vector

$d(n)$  - desired response

$\hat{W}(n)$  - weight vector.

$$\text{then } \frac{\partial E(\hat{W})}{\partial \hat{W}} = e(n) \frac{\partial e(n)}{\partial \hat{W}}$$

$$\text{and } \frac{\partial e(n)}{\partial \hat{W}} = -X(n)$$

$$\Rightarrow \frac{\partial E(\hat{W})}{\partial \hat{W}(n)} = -X(n) e(n)$$

$$\text{or } \hat{g}(n) = -X(n) e(n)$$

The LMS algorithm:

$$\hat{W}(n+1) = \hat{W}(n) + \eta X(n) e(n)$$

## Wiener Filter

Wiener filter can be spatial

$$y = \sum_{k=1}^M w_k x_k$$

or temporal, e.g. MA filter.

Let  $d$ : desired response

$e$ : error signal.

$$e = d - y$$

$\mathcal{E}$ : mean-square error.

$$\mathcal{E} = \frac{1}{2} E(e^2) = \frac{1}{2} E(d - y)^2$$

The linear optimal filter problem:

Determine the optimal set of weights,  $w_0, w_1, w_2, \dots, w_M$ ,  
for which the mean-squared error  $\mathcal{E}$  is minimized

→ the solution of the linear optimal filter problem  
is the Wiener filter.

Detailed derivations:

$$\mathcal{E} = \frac{1}{2} E(e^2) = \frac{1}{2} E(d - y)^2$$

$$= \frac{1}{2} E(d^2) - E\left(\sum_{k=1}^M w_k x_k d\right) + \frac{1}{2} E\left(\sum_{j=1}^M \sum_{k=1}^M w_j w_k x_j x_k\right)$$

$$= \frac{1}{2} E(d^2) - \sum_{k=1}^M w_k E(x_k d) + \frac{1}{2} \sum_{j=1}^M \sum_{k=1}^M w_j w_k E(x_j x_k)$$

let  $r_d = E(d^2)$  : MS value of desired response

$r_{dx}(k) = E(d x_k)$ ,  $k=1, \dots, M$ , cross-correlation

$r_x(j, k) = E(x_j x_k)$ ,  $j, k=1, \dots, M$ , auto-correlation

Therefore :

$$\mathcal{E} = \frac{1}{2} r_d - \sum_{k=1}^M w_k r_{dx}(k) + \frac{1}{2} \sum_{j=1}^M \sum_{k=1}^M w_j w_k r_x(j, k)$$

Necessary condition of optimal  $\mathcal{E}$ ,

$$\nabla_{w_k} \mathcal{E} = 0, \quad k=1, 2, \dots, M$$

where

$$\nabla_{w_k} \mathcal{E} = \frac{\partial \mathcal{E}}{\partial w_k} \quad k=1, \dots, M$$

$$= -r_{dx}(k) + \sum_{j=1}^M w_j r_x(j, k)$$

The Wiener filter :

$$\sum_{j=1}^M w_j r_x(j, k) = r_{dx}(k), \quad k=1, \dots, M$$

In matrix form

$$R_{xx} w_0 = R_{xd}$$

$$w_0 = R_{xx}^{-1} R_{xd}$$

Remark:  $R_{xx}$  is a Toeplitz matrix

using Levinson-Durbin recursion to solve,  $O(Nw^2)$

Alternatively, one can use steepest descent,

$$\text{Since } \nabla_{w_k} \mathcal{E}(n) = -r_{dx}(k) + \sum_{j=1}^M w_j(n) r_x(j, k)$$

$$\Delta w_k(n) = -\eta \frac{\partial \mathcal{E}}{\partial w_k(n)} = -\eta \nabla_{w_k} \mathcal{E}(n)$$

$k=1, \dots, M$

$$\text{or } w_k(n+1) = w_k(n) + \Delta w_k(n)$$

$$= w_k(n) - \eta \nabla_{w_k} \mathcal{E}(n)$$

$$= w_k(n) + \eta \left[ r_{dx}(k) - \sum_{j=1}^M w_j(n) r_x(j, k) \right],$$

$k=1, \dots, M$

Reminder in the above,

$$\mathcal{E}(n) = \frac{1}{2} \sigma[e^2(n)]$$

If we let

$$\hat{r}_x(j, k; n) = x_j(n) x_k(n)$$

$$\hat{r}_{dx}(k; n) = x_k(n) d(n).$$

$$\Rightarrow \hat{w}_k(n+1) = \hat{w}_k(n) + \eta \left[ x_k(n) d(n) - \sum_{j=1}^M \hat{w}_j(n) x_j(n) x_k(n) \right]$$

$$= \hat{w}_k(n) + \eta \left[ d(n) - \sum_{j=1}^M \hat{w}_j(n) x_j(n) \right] x_k(n)$$

$$\text{The LMS } \hat{w}_k(n) = \hat{w}_k(n) + \eta \left[ d(n) - y(n) \right] x_k(n)$$

$k = 1, \dots, M$

where

$$y(n) = \sum_{j=1}^M \hat{w}_j(n) x_j(n)$$

## Remarks on LMS algorithm

- "approximate", however very simple
  - minimizing instantaneous error  $\varepsilon(n)$  where  

$$\varepsilon(n) = \frac{1}{2} e^2(n)$$
  - no memory requirement
  - optimal in  $H^\infty$  - cater to worst-case scenarios
  - Slow convergence especially if poor condition # of  $R_{xx}$ , or sensitive to condition # of  $R_{xx}$
  - Convergence of the LMS depends on statistical property of  $x(n)$  and learning rate  $\eta$ .
- convergence in the mean

$$E[\hat{w}(n)] \rightarrow w_0 \quad \text{as } n \rightarrow \infty$$

convergence in MS sense

$$E[e^2(n)] \rightarrow \text{constant} \quad \text{as } n \rightarrow \infty$$

conv. in MS  $\begin{matrix} \rightarrow \\ \leftarrow \end{matrix}$  conv. in mean

Sufficient condition for convergence,

$$0 < \eta < \frac{2}{\lambda_{\max}}$$

relaxed

$$0 < \eta < \frac{2}{\text{tr}[R_x]}$$



where  $R_{xx} = E[\tilde{x}(n)\tilde{x}^T(n)]$

$\lambda_{\max}$  - largest eigenvalue of  $R_{xx}$

•  $\eta$  in practice (annealing schedules)

LMS:  $\eta(n) = \eta_0 = \text{small constant for all } n$

Robbins-Monro:  $\eta(n) = \frac{c}{n}$ , where  $c$  is constant  
reference Lin & S., in Neural Comp 1998

Darkeu-Moody:  $\eta(n) = \frac{\eta_0}{1 + n/\bar{\tau}}$  where  $\eta_0, \bar{\tau}$  are  
chosen constants

can be more aggressive initially,  
approaches Robbins-Monro as  $n \rightarrow \infty$ .