

Lecture 2

1

During this lecture you will learn about

- The *Least Mean Squares* algorithm (LMS)
- Convergence analysis of the LMS
- Equalizer (Kanalutjämnare)

Background

2

The method of the *Steepest descent* that was studied at the last lecture is a recursive algorithm for calculation of the Wiener filter when the statistics of the signals are known (knowledge about \mathbf{R} och \mathbf{p}).

The problem is that this information is often unknown!

LMS is a method that is based on the same principles as the method of the Steepest descent, but where the statistics is estimated continuously.

Since the statistics is estimated continuously, the LMS algorithm can adapt to changes in the signal statistics; The LMS algorithm is thus an adaptive filter.

Because of estimated statistics the gradient becomes noisy. The LMS algorithm belongs to a group of methods referred to as *stochastic gradient* methods, while the method of the Steepest descent belongs to the group *deterministic gradient* methods.

The Least Mean Square (LMS) algorithm

3

We want to create an algorithm that minimizes $E\{|e(n)|^2\}$, just like the SD, but based on unknown statistics.

A strategy that then can be used is to use estimates of the autocorrelation matrix \mathbf{R} and the cross correlation vector \mathbf{p} . If instantaneous estimates are chosen,

$$\hat{\mathbf{R}}(n) = \mathbf{u}(n)\mathbf{u}^H(n)$$

$$\hat{\mathbf{p}}(n) = \mathbf{u}(n)d^*(n)$$

the resulting method is the *Least Mean Squares* algorithm.

The Least Mean Square (LMS) algorithm

4

For the SD, the update of the filter weights is given by

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \frac{1}{2}\mu[-\nabla J(n)]$$

where $\nabla J(n) = -2\mathbf{p} + 2\mathbf{R}\mathbf{w}(n)$.

In the LMS we use the estimates $\hat{\mathbf{R}}$ och $\hat{\mathbf{p}}$ to calculate $\hat{\nabla} J(n)$. Thus, also the updated filter vector becomes an estimate. It is therefore denoted $\hat{\mathbf{w}}(n)$;

$$\hat{\nabla} J(n) = -2\hat{\mathbf{p}}(n) + 2\hat{\mathbf{R}}(n)\hat{\mathbf{w}}(n)$$

The Least Mean Square (LMS) algorithm 5

For the LMS algorithm, the update equation of the filter weights becomes

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \mu \mathbf{u}(n) e^*(n)$$

where $e^*(n) = d^*(n) - \mathbf{u}^H(n) \hat{\mathbf{w}}(n)$.

Compare with the corresponding equation for the SD:

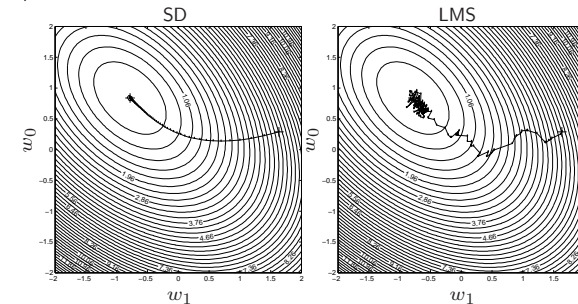
$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mu E\{\mathbf{u}(n) e^*(n)\}$$

where $e^*(n) = d^*(n) - \mathbf{u}^H(n) \mathbf{w}(n)$.

The difference is thus that $E\{\mathbf{u}(n) e^*(n)\}$ is estimated as $\mathbf{u}(n) e^*(n)$. This leads to *gradient noise*.

The Least Mean Square (LMS) algorithm 6

Illustration of *gradient noise*. \mathbf{R} och \mathbf{p} in the example from Lecture 1. $\mu = 0.1$.



Because of the noise in the gradient the LMS never reaches J_{min} , which the SD does.

Convergence analysis of the LMS, overview 7

Consider the update equation

$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \mu \mathbf{u}(n) e^*(n)$$

We want to know how fast and towards what the algorithm converges, in terms of $J(n)$ and $\hat{\mathbf{w}}(n)$.

Strategy:

1. Introduce the filter weight error vector $\boldsymbol{\epsilon}(n) = \hat{\mathbf{w}}(n) - \mathbf{w}_o$.
2. Express the update equation in $\boldsymbol{\epsilon}(n)$.
3. Express $J(n)$ in $\boldsymbol{\epsilon}(n)$. This leads to $\mathbf{K}(n) = E\{\boldsymbol{\epsilon}(n) \boldsymbol{\epsilon}^H(n)\}$.
4. Calculate the difference equation in $\mathbf{K}(n)$, which controls the convergence.
5. Make a variable change from $\mathbf{K}(n)$ to $\mathbf{X}(n)$ which converges equally fast.

1. The filter weight error vector 8

The LMS contains feedback, and just like for the SD there is a risk that the algorithm diverges, if the stepsize μ is not chosen appropriately.

In order to investigate the stability, the filter weight error vector is introduced $\boldsymbol{\epsilon}(n)$:

$$\boldsymbol{\epsilon}(n) = \hat{\mathbf{w}}(n) - \mathbf{w}_o \quad (\mathbf{w}_o = \mathbf{R}^{-1} \mathbf{p} \text{ Wiener-Hopf})$$

Note that $\boldsymbol{\epsilon}(n)$ corresponds to $\mathbf{c}(n) = \mathbf{w}(n) - \mathbf{w}_o$ for the *Steepest descent*, but that $\boldsymbol{\epsilon}(n)$ is stochastic because of $\hat{\mathbf{w}}(n)$.

2. Update equation in $\epsilon(n)$

9

The update equation for $\epsilon(n)$ at time $n+1$ can be expressed recursively as

$$\begin{aligned}\epsilon(n+1) &= [\mathbf{I} - \mu \hat{\mathbf{R}}(n)]\epsilon(n) + \mu[\hat{\mathbf{p}}(n) - \hat{\mathbf{R}}(n)\mathbf{w}_o] \\ &= [\mathbf{I} - \mu \mathbf{u}(n)\mathbf{u}^H(n)]\epsilon(n) + \mu \mathbf{u}(n)e_o^*(n)\end{aligned}$$

where $e_o^*(n) = d(n) - \mathbf{w}_o^H \mathbf{u}(n)$.

Compare to that of the SD:

$$\mathbf{c}(n+1) = (\mathbf{I} - \mu \mathbf{R})\mathbf{c}(n)$$

3. Express $J(n)$ in $\epsilon(n)$

10

The convergence analysis of the LMS is considerably more complicated than that for the SD. Therefore, two tools are needed (approximations).

These are

- Independence theory
- Direct-averaging method

3. Express $J(n)$ in $\epsilon(n)$, cont.

11

Independence Theory

1. The input vectors of different time instants n , $\mathbf{u}(1), \mathbf{u}(2), \dots, \mathbf{u}(n)$, is independent of each other (and thereby uncorrelated).
2. The input signal vector at time n is independent of the desired signal at all earlier time instants, $\mathbf{u}(n)$ independent of $d(1), d(2), \dots, d(n-1)$.
3. The desired signal at time n , $d(n)$, is independent of all earlier desired signals, $d(1), d(2), \dots, d(n-1)$, but dependent of the input signal vector at time n , $\mathbf{u}(n)$.
4. The signals $d(n)$ and $\mathbf{u}(n)$ at the same time instant n is mutually normally distributed.

3. Express $J(n)$ in $\epsilon(n)$, cont.

12

Direct-averaging

Direct-averaging means that in the update equation for ϵ ,

$$\epsilon(n+1) = [\mathbf{I} - \mu \mathbf{u}(n)\mathbf{u}^H(n)]\epsilon(n) + \mu \mathbf{u}(n)e_o^*(n),$$

the instantaneous estimate $\hat{\mathbf{R}}(n) = \mathbf{u}(n)\mathbf{u}^H(n)$ is substituted with the expectation over the ensemble, $\mathbf{R} = E\{\mathbf{u}(n)\mathbf{u}^H(n)\}$, such that

$$\epsilon(n+1) = (\mathbf{I} - \mu \mathbf{R})\epsilon(n) + \mu \mathbf{u}(n)e_o^*(n)$$

results.

3. Express $J(n)$ i $\epsilon(n)$, cont.

13

For the SD we could show that

$$J(n) = J_{min} + [\mathbf{w}(n) - \mathbf{w}_o]^H \mathbf{R} [\mathbf{w}(n) - \mathbf{w}_o]$$

and that the eigenvalues of \mathbf{R} controls the convergence speed.

For the LMS, $\epsilon(n) = \hat{\mathbf{w}}(n) - \mathbf{w}_o$,

$$\begin{aligned} J(n) &= E\{|e(n)|^2\} = E\{|d(n) - \hat{\mathbf{w}}^H(n)\mathbf{u}(n)|^2\} \\ &= E\{|e_o(n) - \epsilon^H(n)\mathbf{u}(n)|^2\} = \\ &\stackrel{i}{=} E\{\underbrace{|e_o(n)|^2}_{J_{min}}\} + E\{\epsilon^H(n)\underbrace{\mathbf{u}(n)\mathbf{u}^H(n)}_{\hat{\mathbf{R}}(n)}\epsilon(n)\} \end{aligned}$$

Here, (i) indicates that the independence assumption is used. Since ϵ and $\hat{\mathbf{R}}$ are estimates, and in addition not independent, the expectation operator cannot just be "moved into $\hat{\mathbf{R}}$ ".

3. Express $J(n)$ i $\epsilon(n)$, cont.

14

The behavior of $J_{ex}(n) = J(n) - J_{min} = E\{\epsilon^H(n)\hat{\mathbf{R}}\epsilon(n)\}$ determines the convergence.

Since the following holds for the trace of a matrix,

1. $\text{tr}(\text{scalar}) = \text{scalar}$
2. $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$
3. $E\{\text{tr}(\mathbf{A})\} = \text{tr}(E\{\mathbf{A}\})$

$J_{ex}(n)$ can be written as

$$\begin{aligned} J_{ex}(n) &\stackrel{1,2}{=} E\{\text{tr}[\hat{\mathbf{R}}(n)\epsilon(n)\epsilon^H(n)]\} \\ &\stackrel{3}{=} \text{tr}[E\{\hat{\mathbf{R}}(n)\epsilon(n)\epsilon^H(n)\}] \\ &\stackrel{i}{=} \text{tr}[E\{\hat{\mathbf{R}}(n)\}E\{\epsilon(n)\epsilon^H(n)\}] \\ &= \text{tr}[\mathbf{RK}(n)] \end{aligned}$$

The convergence thus depends on $\mathbf{K}(n) = E\{\epsilon(n)\epsilon^H(n)\}$.

4. Difference equation of $\mathbf{K}(n)$

15

The update equation for the filter weight error is, as shown previously,

$$\epsilon(n+1) = [\mathbf{I} - \mu\mathbf{u}(n)\mathbf{u}^H(n)]\epsilon(n) + \mu\mathbf{u}(n)e_o^*(n)$$

The solution to this stochastic difference equation is for small μ close to the solution for

$$\epsilon(n+1) = (\mathbf{I} - \mu\mathbf{R})\epsilon(n) + \mu\mathbf{u}(n)e_o^*(n),$$

which is resulting from *Direct-averaging*.

This equation is still difficult to solve, but now the behavior of $\mathbf{K}(n)$ can be described by the following difference equation

$$\mathbf{K}(n+1) = (\mathbf{I} - \mu\mathbf{R})\mathbf{K}(n)(\mathbf{I} - \mu\mathbf{R}) + \mu^2 J_{min}\mathbf{R}$$

5. Changing of variables from $\mathbf{K}(n)$ to $\mathbf{X}(n)$

16

The variable changes

$$\begin{aligned} \mathbf{Q}^H \mathbf{R} \mathbf{Q} &= \mathbf{\Lambda} \\ \mathbf{Q}^H \mathbf{K}(n) \mathbf{Q} &= \mathbf{X}(n) \end{aligned}$$

where $\mathbf{\Lambda}$ is the diagonal eigenvalue matrix of \mathbf{R} , and where $\mathbf{X}(n)$ normally is not diagonal, yields

$$J_{ex}(n) = \text{tr}(\mathbf{RK}(n)) = \text{tr}(\mathbf{\Lambda X}(n))$$

The convergence can be described by

$$\mathbf{X}(n+1) = (\mathbf{I} - \mu\mathbf{\Lambda})\mathbf{X}(n)(\mathbf{I} - \mu\mathbf{\Lambda}) + \mu^2 J_{min}\mathbf{\Lambda}$$

5. Changing of variables from $\mathbf{K}(n)$ to $\mathbf{X}(n)$, cont. 17

$J_{ex}(n)$ depends only on the diagonal elements of $\mathbf{X}(n)$ (because of the trace $\text{tr}[\mathbf{A}\mathbf{X}(n)]$). We can therefore write $J_{ex}(n) = \sum_{i=1}^M \lambda_i x_i(n)$.

The update of the diagonal elements is controlled by

$$x_i(n+1) = (1 - \mu\lambda_i)x_i(n) + \mu^2 J_{min}\lambda_i$$

This is an inhomogenous difference equation. This means that the solution will contain a transient part, and a stationary part that depends on J_{min} and μ . The cost function can therefore be written

$$J(n) = J_{min} + J_{ex}(\infty) + J_{trans}(n)$$

where $J_{min} + J_{ex}(\infty)$ och $J_{trans}(n)$ is the stationary and transient solutions, respectively.

At the best, the LMS can thus reach $J_{min} + J_{ex}(\infty)$.

Properties of the LMS 18

- Convergence for the same interval as the SD:

$$0 < \mu < \frac{2}{\lambda_{max}}$$
- $J_{ex}(\infty) = J_{min} \sum_{i=1}^M \frac{\mu\lambda_i}{2 - \mu\lambda_i}$
- Misadjustment $\mathcal{M} = \frac{J_{ex}(\infty)}{J_{min}}$ is a measure of how close the the optimal solution the LMS (in mean-square sense).

Rules of thumb LMS 19

The eigenvalues of \mathbf{R} are rarely known. Therefore, a set of rules of thumbs is commonly used, that is based on the *tap-input power*, $\sum_{k=0}^{M-1} E\{|u(n-k)|^2\} = Mr(0) = \text{tr}(\mathbf{R})$.

- The stepsize $0 < \mu < \frac{2}{\sum_{k=0}^{M-1} E\{|u(n-k)|^2\}}$
- Misadjustment $\mathcal{M} \approx \frac{\mu}{2} \sum_{k=0}^{M-1} E\{|u(n-k)|^2\}$
- Average time constant $\tau_{mse,av} \approx \frac{1}{2\mu\lambda_{av}}$ där
 $\lambda_{av} = \frac{1}{M} \sum_{i=1}^M \lambda_i$.
 When λ_i is unknown, the following relationship is useful.
- Approximate relationship \mathcal{M} and $\tau_{mse,av}$ $\mathcal{M} \approx \frac{M}{4\tau_{mse,av}}$

Here, $\tau_{mse,av}$ denotes the time it takes for $J(n)$ to decay a factor of e^{-1} .

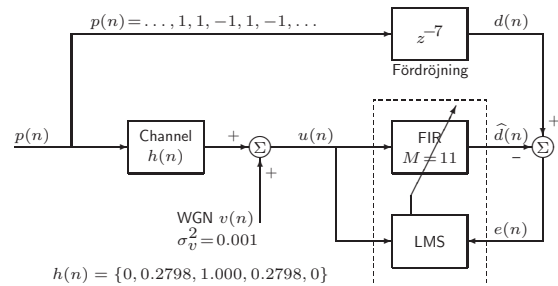
Summary LMS 20

Summary of the iterations:

- Set start values for the filter coefficients, $\hat{\mathbf{w}}(0)$
- Calculate the error $e(n) = d(n) - \hat{\mathbf{w}}^H(n)\mathbf{u}(n)$
- Update the filter coefficients

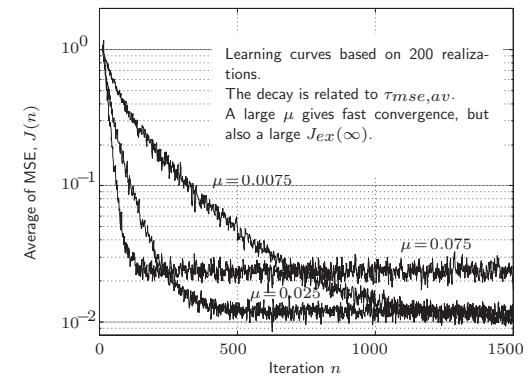
$$\hat{\mathbf{w}}(n+1) = \hat{\mathbf{w}}(n) + \mu\mathbf{u}(n)[d^*(n) - \mathbf{u}^H(n)\hat{\mathbf{w}}(n)]$$
- Repeat steps 2 and 3

Example: Channel equalizer in training phase 21



Sent information in the form of a pn -sequence is affected by a channel ($h(n)$). An adaptive inverse filter is used to compensate for the channel such that the total effect can be approximated by a delay. White Gaussian noise is added during the transmission (WGN).

Example, cont: Learning Curve 22



The curves illustrates learning curves for two different μ .

What to read? 23

- Least Mean Squares, Haykin chap.5.1–5.3, 5.4 (see lecture), 5.5–5.7

Exercises: 4.1, 4.2, 1.2, 4.3, 4.5, (4.4)

Computer exercise, theme: Implement the LMS algorithm in Matlab, and generate the results in Haykin chap.5.7