

**Stability of the Steepest-Descent Algorithm**

Since the SD algorithm can be viewed as a feedback model, it is subject to the possibility of becoming unstable, i.e.

$$\lim_{n \rightarrow \infty} \mathbf{w}(n) \neq \mathbf{w}_o \text{ and hence } \lim_{n \rightarrow \infty} J[\mathbf{w}(n)] \neq J_{\min}.$$

We will show that the stability performance of the SD algorithm is determined by two factors: the step-size parameter,  $\mu$  and the largest eigenvalue  $\lambda_{\max}$  of the correlation matrix,  $\mathbf{R}$ .

We begin with the SD recursion (2)

$$\begin{aligned} \mathbf{w}(n+1) &= \mathbf{w}(n) + \mu[\mathbf{p} - \mathbf{R}\mathbf{w}(n)] \\ &= \mathbf{w}(n) + \mu[\mathbf{R}\mathbf{w}_o - \mathbf{R}\mathbf{w}(n)] \end{aligned}$$

Subtracting the optimal weight vector (translation of coordinate system) from both sides yields

$$\begin{aligned} \mathbf{w}(n+1) - \mathbf{w}_o &= [\mathbf{w}(n) - \mathbf{w}_o] - \mu[\mathbf{R}\mathbf{w}(n) - \mathbf{R}\mathbf{w}_o] \\ &= (\mathbf{I} - \mu\mathbf{R})[\mathbf{w}(n) - \mathbf{w}_o] \end{aligned}$$

Defining the misalignment vector (weight-error vector) at time  $n$

$$\mathbf{c}(n) = \mathbf{w}(n) - \mathbf{w}_o$$

yields

$$\begin{aligned} \mathbf{c}(n+1) &= \mathbf{w}(n) - \mathbf{w}_o + \mu[\mathbf{R}\mathbf{w}_o - \mathbf{R}\mathbf{w}(n)] \\ &= (\mathbf{I} - \mu\mathbf{R})\mathbf{c}(n) \end{aligned}$$

In order for stability of the algorithm and hence *convergence* of  $\mathbf{w}(n)$  to  $\mathbf{w}_o$  we need to drive  $\mathbf{c}(n)$  to  $\mathbf{0}$ . To continue the analysis, we employ the diagonalization of  $\mathbf{R}$  to rewrite the above as

$$\mathbf{c}(n+1) = (\mathbf{I} - \mu\mathbf{\Theta}\mathbf{\Lambda}\mathbf{\Theta}^H)\boldsymbol{\chi}(n).$$

We apply our unitary transformation (rotation of coordinate system) so that we can decouple the weight error elements from each other. We have

$$\begin{aligned} \mathbf{Q}^H \mathbf{c}(n+1) &= \mathbf{\Theta}^H (\mathbf{I} - \mu\mathbf{\Theta}\mathbf{\Lambda}\mathbf{\Theta}^H)\boldsymbol{\chi}(n) \\ &= \left( \mathbf{\Theta}^H - \mu \mathbf{\Theta}^H \mathbf{\Theta} \mathbf{\Lambda} \mathbf{\Theta}^H \right) \boldsymbol{\chi}(n) \\ &= (\mathbf{I} - \mu\mathbf{\Lambda})\mathbf{\Theta}^H \boldsymbol{\chi}(n) \end{aligned}$$

As before, we denote the translated and rotated (transformed) vectors as

$$\begin{aligned} \mathbf{v}(n) &= \mathbf{Q}^H \mathbf{c}(n) \\ &= \mathbf{Q}^H [\mathbf{w}(n) - \mathbf{w}_o] \end{aligned}$$

which leads to

$$\mathbf{v}(n+1) = (\mathbf{I} - \mu\mathbf{\Lambda})\mathbf{v}(n).$$

We examine a single weight error element,  $v_i(n)$  (which due to the coordinate transformation and decoupling is only a function of itself)

$$\begin{aligned} v_k(n+1) &= (1 - \mu\lambda_k)v_k(n) \\ &= (1 - \mu\lambda_k)^n v_k(0) \end{aligned} \quad (\text{kth natural mode})$$

We remember that for stability, we wish to drive the misalignment  $\mathbf{c}(n)$  to  $\mathbf{0}$  and consequently each  $v_k(n)$  to 0, for  $k = 1, \dots, M$ , as quickly as possible.

Clearly for each mode to die out we require

$$-1 < 1 - \mu\lambda_k < 1$$

or

$$0 < \mu < \frac{2}{\lambda_k}.$$

We note that

$$\frac{2}{\lambda_{\mu\alpha\xi}} \leq \frac{2}{\lambda_k} \leq \frac{2}{\lambda_{\mu\nu}}$$

thus we must select the  $\mu$  which guarantees all modes die out. That  $\mu$  is given by

$$\boxed{0 < \mu < \frac{2}{\lambda_{\mu\alpha\xi}}}.$$

We finally note that the mode which takes longest to die out is that associated with  $\lambda_{\min}$

$$v_k(n+1) = (1 - \mu\lambda_{\min})^n v_k(0).$$

We therefore conclude that when the eigenvalues are spread over a large range (or equivalently the condition number,  $\chi = \frac{\lambda_{\max}}{\lambda_{\min}}$  is large) the settling time of the SD algorithm is limited by the smallest eigenvalues.

### Transient Behavior of the Coefficient Vector

We have defined the misalignment vector or weight error vector as

$$\mathbf{c}(n) = \boldsymbol{\omega}(n) - \boldsymbol{\omega}_o$$

and the misalignment vector under coordinate rotation (transformed weight error vector) as

$$\begin{aligned} v(n) &= \mathbf{Q}^H \mathbf{c}(n) \\ &= \mathbf{Q}^H [\mathbf{w}(n) - \mathbf{w}_o] \end{aligned}$$

Rearranging, we have

$$\mathbf{Q}^H \mathbf{w}(n) = \mathbf{Q}^H \mathbf{w}_o + v(n)$$

which when premultiplied by  $\mathbf{Q}$  gives us (since  $\mathbf{Q}\mathbf{Q}^H = \mathbf{I}$ )

$$\begin{aligned}
\mathbf{w}(n) &= \boldsymbol{\omega}_o + \boldsymbol{\Theta} \boldsymbol{v}(n) \\
&= \boldsymbol{\omega}_o + [\boldsymbol{\theta}_1 \ \Lambda \ \boldsymbol{\theta}_M] \begin{bmatrix} \boldsymbol{v}_1(n) \\ \mathbf{M} \\ \boldsymbol{v}_M(n) \end{bmatrix} \\
&= \boldsymbol{\omega}_o + \sum_{k=1}^M \boldsymbol{\theta}_k \boldsymbol{v}_k(n)
\end{aligned}$$

Recall the equation for the  $k$ th mode (backsolved)

$$\boldsymbol{v}_k(n) = (1 - \mu \lambda_k)^n \boldsymbol{v}_k(0)$$

and consider the  $i$ th coefficient of  $\mathbf{w}(n)$  (substituting the modal equation)

$$w_i(n) = \omega_{oi} + \sum_{k=1}^M \theta_{ki} \boldsymbol{v}_k(0) (1 - \mu \lambda_k)^n.$$

This expression provides a description of the transient behavior of the  $i$ th coefficient of  $\mathbf{w}(n)$ . Each coefficient of  $\mathbf{w}(n)$  consists of a weighted sum of exponentials. If we choose

$$0 < \mu < \frac{2}{\lambda_{\mu \alpha \xi}}$$

then we guarantee that

$$-1 < (1 - \mu \lambda_k) < 1, \quad 1 \leq k \leq M$$

and hence

$$\lim_{n \rightarrow \infty} (1 - \mu \lambda_k)^n = 0, \quad 1 \leq k \leq M.$$

This of course implies that

$$\lim_{n \rightarrow \infty} \omega_i(n) = \omega_{oi}$$

and

$$\boxed{\lim_{n \rightarrow \infty} \boldsymbol{\omega}(n) = \boldsymbol{\omega}_o}.$$

We note that our coefficient vector converges to the Wiener filter regardless of  $\mathbf{w}(0)$ .

### Transient Behavior of the Mean-Squared Error

Since our error surface is quadratic with a single, global minimum and  $\lim_{n \rightarrow \infty} \boldsymbol{\omega}(n) = \boldsymbol{\omega}_o$  we are guaranteed to have

$$\boxed{\lim_{n \rightarrow \infty} \mathcal{J}(n) = \mathcal{J}_{\text{MSE}}}.$$

This can also be derived independently using our canonical form (Haykin 5.57) of the MSE derived for the Wiener filter

$$\begin{aligned}
 J(n) &= \mathcal{J}_{\text{MSE}} + \sum_{\kappa=1}^M \lambda_{\kappa} |v_{\kappa}(n)|^2 \\
 &= \mathcal{J}_{\text{MSE}} + \sum_{\kappa=1}^M \lambda_{\kappa} \left| (1 - \mu \lambda_{\kappa})^n v_{\kappa}(0) \right|^2 \\
 &= \mathcal{J}_{\text{MSE}} + \sum_{\kappa=1}^M \lambda_{\kappa} (1 - \mu \lambda_{\kappa})^{2n} |v_{\kappa}(0)|^2
 \end{aligned}$$

If  $0 < \mu < \frac{2}{\lambda_{\mu\alpha\xi}}$  then

$$\lim_{n \rightarrow \infty} \mathcal{J}(n) = \mathcal{J}_{\text{MSE}}$$

as we've already proved.

---

To get a feel for the progress our search algorithm makes in finding  $\mathbf{w}_o$ , we typically plot

- $\mathbf{w}(n)$  vs.  $n$  [usually  $\mathbf{w}(n)$  is length 2] or  $\mathbf{v}(n)$  vs.  $n$ . Such a plot is called a weight trajectory or weight error trajectory. Usually we plot  $w_0(n)$  and  $w_1(n)$  as a function of  $n$  along with the loci of  $w_0(n)$  and  $w_1(n)$  for fixed values of the MSE (contours).
- $J(n)$  vs.  $n$ . Such a plot is called the learning curve or simply the MSE plot.