

can be written as

$$\nabla E = \sum_i \alpha_i \lambda_i \mathbf{u}_i. \quad (7.24)$$

Again using (7.10) we can express the change in the weight vector in terms of corresponding changes in the coefficients $\{\alpha_i\}$:

$$\Delta \mathbf{w} = \sum_i \Delta \alpha_i \mathbf{u}_i. \quad (7.25)$$

Combining (7.24) with (7.25) and the gradient descent formula (7.16) and using the orthonormality relation (7.9) for the eigenvectors of the Hessian, we obtain the following expression for the change in α_i at each step of the gradient descent algorithm:

$$\Delta \alpha_i = -\eta \lambda_i \alpha_i \quad (7.26)$$

Exercise 7.10

from which it follows that

$$\alpha_i^{\text{new}} = (1 - \eta \lambda_i) \alpha_i^{\text{old}} \quad (7.27)$$

where ‘old’ and ‘new’ denote values before and after a weight update. Using the orthonormality relation (7.9) for the eigenvectors together with (7.10), we have

$$\mathbf{u}_i^T (\mathbf{w} - \mathbf{w}^*) = \alpha_i \quad (7.28)$$

and so α_i can be interpreted as the distance to the minimum along the direction \mathbf{u}_i . From (7.27) we see that these distances evolve independently such that, at each step, the distance along the direction of \mathbf{u}_i is multiplied by a factor $(1 - \eta \lambda_i)$. After a total of T steps we have

$$\alpha_i^{(T)} = (1 - \eta \lambda_i)^T \alpha_i^{(0)}. \quad (7.29)$$

It follows that, provided $|1 - \eta \lambda_i| < 1$, the limit $T \rightarrow \infty$ leads to $\alpha_i = 0$, which from (7.28) shows that $\mathbf{w} = \mathbf{w}^*$ and so the weight vector has reached the minimum of the error.

Note that (7.29) demonstrates that gradient descent leads to linear convergence in the neighbourhood of a minimum. Also, convergence to the stationary point requires that all the λ_i be positive, which in turn implies that the stationary point is indeed a minimum. By making η larger we can make the factor $(1 - \eta \lambda_i)$ smaller and hence improve the speed of convergence. There is a limit to how large η can be made, however. We can permit $(1 - \eta \lambda_i)$ to go negative (which gives oscillating values of α_i), but we must ensure that $|1 - \eta \lambda_i| < 1$ otherwise the α_i values will diverge. This limits the value of η to $\eta < 2/\lambda_{\max}$ where λ_{\max} is the largest of the eigenvalues. The rate of convergence, however, is dominated by the smallest eigenvalue, so with η set to its largest permitted value, the convergence along the direction corresponding to the smallest eigenvalue (the long axis of the ellipse in Figure 7.3) will be governed by

$$\left(1 - \frac{2\lambda_{\min}}{\lambda_{\max}}\right) \quad (7.30)$$