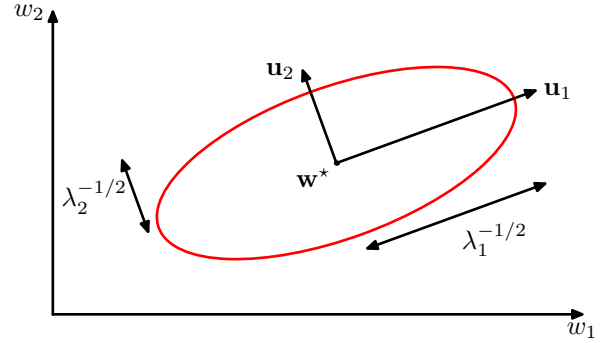


Figure 7.2 In the neighbourhood of a minimum \mathbf{w}^* , the error function can be approximated by a quadratic. Contours of constant error are then ellipses whose axes are aligned with the eigenvectors \mathbf{u}_i of the Hessian matrix, with lengths that are inversely proportional to the square roots of the corresponding eigenvalues λ_i .



Because the eigenvectors $\{\mathbf{u}_i\}$ form a complete set, an arbitrary vector \mathbf{v} can be written in the form

$$\mathbf{v} = \sum_i c_i \mathbf{u}_i. \quad (7.13)$$

From (7.8) and (7.9), we then have

$$\mathbf{v}^T \mathbf{H} \mathbf{v} = \sum_i c_i^2 \lambda_i \quad (7.14)$$

Exercise 7.2

and so \mathbf{H} will be positive definite if, and only if, all its eigenvalues are positive. Thus, a necessary and sufficient condition for \mathbf{w}^* to be a local minimum is that the gradient of the error function should vanish at \mathbf{w}^* and the Hessian matrix evaluated at \mathbf{w}^* should be positive definite. In the new coordinate system, whose basis vectors are given by the eigenvectors $\{\mathbf{u}_i\}$, the contours of constant $E(\mathbf{w})$ are axis-aligned ellipses centred on the origin, as illustrated in [Figure 7.2](#).

Exercise 7.3

Exercise 7.6

7.2. Gradient Descent Optimization

There is little hope of finding an analytical solution to the equation $\nabla E(\mathbf{w}) = 0$ for an error function as complex as one defined by a neural network, and so we resort to iterative numerical procedures. The optimization of continuous nonlinear functions is a widely studied problem, and there exists an extensive literature on how to solve it efficiently. Most techniques involve choosing some initial value $\mathbf{w}^{(0)}$ for the weight vector and then moving through weight space in a succession of steps of the form

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} + \Delta \mathbf{w}^{(\tau-1)} \quad (7.15)$$

where τ labels the iteration step. Different algorithms involve different choices for the weight vector update $\Delta \mathbf{w}^{(\tau)}$.

Because of the complex shape of the error surface for all but the simplest neural networks, the solution found will depend, among other things, on the particular choice of initial parameter values $\mathbf{w}^{(0)}$. To find a sufficiently good solution, it may