



Figure 7.4: An illustration of the effect of early stopping. (Left) The solid contour lines indicate the contours of the negative log-likelihood. The dashed line indicates the trajectory taken by SGD beginning from the origin. Rather than stopping at the point  $w^*$  that minimizes the cost, early stopping results in the trajectory stopping at an earlier point  $\tilde{w}$ . (Right) An illustration of the effect of  $L^2$  regularization for comparison. The dashed circles indicate the contours of the  $L^2$  penalty, which causes the minimum of the total cost to lie nearer the origin than the minimum of the unregularized cost.

We are going to study the trajectory followed by the parameter vector during training. For simplicity, let us set the initial parameter vector to the origin,<sup>3</sup> that is  $w^{(0)} = \mathbf{0}$ . Let us study the approximate behavior of gradient descent on  $J$  by analyzing gradient descent on  $\hat{J}$ :

$$w^{(\tau)} = w^{(\tau-1)} - \epsilon \nabla_w \hat{J}(w^{(\tau-1)}) \quad (7.35)$$

$$= w^{(\tau-1)} - \epsilon H(w^{(\tau-1)} - w^*) \quad (7.36)$$

$$w^{(\tau)} - w^* = (I - \epsilon H)(w^{(\tau-1)} - w^*). \quad (7.37)$$

Let us now rewrite this expression in the space of the eigenvectors of  $H$ , exploiting the eigendecomposition of  $H$ :  $H = Q\Lambda Q^\top$ , where  $\Lambda$  is a diagonal matrix and  $Q$  is an orthonormal basis of eigenvectors.

$$w^{(\tau)} - w^* = (I - \epsilon Q\Lambda Q^\top)(w^{(\tau-1)} - w^*) \quad (7.38)$$

$$Q^\top(w^{(\tau)} - w^*) = (I - \epsilon \Lambda)Q^\top(w^{(\tau-1)} - w^*) \quad (7.39)$$

<sup>3</sup>For neural networks, to obtain symmetry breaking between hidden units, we cannot initialize all the parameters to  $\mathbf{0}$ , as discussed in section 6.2. However, the argument holds for any other initial value  $w^{(0)}$ .