

Assuming that  $\mathbf{w}^{(0)} = 0$  and that  $\epsilon$  is chosen to be small enough to guarantee  $|1 - \epsilon\lambda_i| < 1$ , the parameter trajectory during training after  $\tau$  parameter updates is as follows:

$$\mathbf{Q}^\top \mathbf{w}^{(\tau)} = [\mathbf{I} - (\mathbf{I} - \epsilon\mathbf{\Lambda})^\tau] \mathbf{Q}^\top \mathbf{w}^*. \quad (7.40)$$

Now, the expression for  $\mathbf{Q}^\top \tilde{\mathbf{w}}$  in equation 7.13 for  $L^2$  regularization can be rearranged as:

$$\mathbf{Q}^\top \tilde{\mathbf{w}} = (\mathbf{\Lambda} + \alpha\mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{w}^* \quad (7.41)$$

$$\mathbf{Q}^\top \tilde{\mathbf{w}} = [\mathbf{I} - (\mathbf{\Lambda} + \alpha\mathbf{I})^{-1} \alpha] \mathbf{Q}^\top \mathbf{w}^* \quad (7.42)$$

Comparing equation 7.40 and equation 7.42, we see that if the hyperparameters  $\epsilon$ ,  $\alpha$ , and  $\tau$  are chosen such that

$$(\mathbf{I} - \epsilon\mathbf{\Lambda})^\tau = (\mathbf{\Lambda} + \alpha\mathbf{I})^{-1} \alpha, \quad (7.43)$$

then  $L^2$  regularization and early stopping can be seen to be equivalent (at least under the quadratic approximation of the objective function). Going even further, by taking logarithms and using the series expansion for  $\log(1+x)$ , we can conclude that if all  $\lambda_i$  are small (that is,  $\epsilon\lambda_i \ll 1$  and  $\lambda_i/\alpha \ll 1$ ) then

$$\tau \approx \frac{1}{\epsilon\alpha}, \quad (7.44)$$

$$\alpha \approx \frac{1}{\tau\epsilon}. \quad (7.45)$$

That is, under these assumptions, the number of training iterations  $\tau$  plays a role inversely proportional to the  $L^2$  regularization parameter, and the inverse of  $\tau\epsilon$  plays the role of the weight decay coefficient.

Parameter values corresponding to directions of significant curvature (of the objective function) are regularized less than directions of less curvature. Of course, in the context of early stopping, this really means that parameters that correspond to directions of significant curvature tend to learn early relative to parameters corresponding to directions of less curvature.

The derivations in this section have shown that a trajectory of length  $\tau$  ends at a point that corresponds to a minimum of the  $L^2$ -regularized objective. Early stopping is of course more than the mere restriction of the trajectory length; instead, early stopping typically involves monitoring the validation set error in order to stop the trajectory at a particularly good point in space. Early stopping therefore has the advantage over weight decay that early stopping automatically determines the correct amount of regularization while weight decay requires many training experiments with different values of its hyperparameter.