

mean squared error, then the approximation is perfect. The approximation \hat{J} is given by

$$\hat{J}(\boldsymbol{\theta}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*), \quad (7.6)$$

where \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* . There is no first-order term in this quadratic approximation, because \mathbf{w}^* is defined to be a minimum, where the gradient vanishes. Likewise, because \mathbf{w}^* is the location of a minimum of J , we can conclude that \mathbf{H} is positive semidefinite.

The minimum of \hat{J} occurs where its gradient

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (7.7)$$

is equal to $\mathbf{0}$.

To study the effect of weight decay, we modify equation 7.7 by adding the weight decay gradient. We can now solve for the minimum of the regularized version of \hat{J} . We use the variable $\tilde{\mathbf{w}}$ to represent the location of the minimum.

$$\alpha \tilde{\mathbf{w}} + \mathbf{H}(\tilde{\mathbf{w}} - \mathbf{w}^*) = \mathbf{0} \quad (7.8)$$

$$(\mathbf{H} + \alpha \mathbf{I}) \tilde{\mathbf{w}} = \mathbf{H} \mathbf{w}^* \quad (7.9)$$

$$\tilde{\mathbf{w}} = (\mathbf{H} + \alpha \mathbf{I})^{-1} \mathbf{H} \mathbf{w}^*. \quad (7.10)$$

As α approaches 0, the regularized solution $\tilde{\mathbf{w}}$ approaches \mathbf{w}^* . But what happens as α grows? Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors, \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$. Applying the decomposition to equation 7.10, we obtain:

$$\tilde{\mathbf{w}} = (\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top + \alpha \mathbf{I})^{-1} \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{w}^* \quad (7.11)$$

$$= \left[\mathbf{Q} (\mathbf{\Lambda} + \alpha \mathbf{I}) \mathbf{Q}^\top \right]^{-1} \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{w}^* \quad (7.12)$$

$$= \mathbf{Q} (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{w}^*. \quad (7.13)$$

We see that the effect of weight decay is to rescale \mathbf{w}^* along the axes defined by the eigenvectors of \mathbf{H} . Specifically, the component of \mathbf{w}^* that is aligned with the i -th eigenvector of \mathbf{H} is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$. (You may wish to review how this kind of scaling works, first explained in figure 2.3).

Along the directions where the eigenvalues of \mathbf{H} are relatively large, for example, where $\lambda_i \gg \alpha$, the effect of regularization is relatively small. However, components with $\lambda_i \ll \alpha$ will be shrunk to have nearly zero magnitude. This effect is illustrated in figure 7.1.