

the sum of squared errors:

$$(\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}). \quad (7.14)$$

When we add  $L^2$  regularization, the objective function changes to

$$(\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \frac{1}{2}\alpha \mathbf{w}^\top \mathbf{w}. \quad (7.15)$$

This changes the normal equations for the solution from

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (7.16)$$

to

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (7.17)$$

The matrix  $\mathbf{X}^\top \mathbf{X}$  in equation 7.16 is proportional to the covariance matrix  $\frac{1}{m} \mathbf{X}^\top \mathbf{X}$ . Using  $L^2$  regularization replaces this matrix with  $(\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1}$  in equation 7.17. The new matrix is the same as the original one, but with the addition of  $\alpha$  to the diagonal. The diagonal entries of this matrix correspond to the variance of each input feature. We can see that  $L^2$  regularization causes the learning algorithm to “perceive” the input  $\mathbf{X}$  as having higher variance, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

### 7.1.2 $L^1$ Regularization

While  $L^2$  weight decay is the most common form of weight decay, there are other ways to penalize the size of the model parameters. Another option is to use  $L^1$  regularization.

Formally,  $L^1$  regularization on the model parameter  $\mathbf{w}$  is defined as:

$$\Omega(\boldsymbol{\theta}) = \|\mathbf{w}\|_1 = \sum_i |w_i|, \quad (7.18)$$

that is, as the sum of absolute values of the individual parameters.<sup>2</sup> We will now discuss the effect of  $L^1$  regularization on the simple linear regression model, with no bias parameter, that we studied in our analysis of  $L^2$  regularization. In particular, we are interested in delineating the differences between  $L^1$  and  $L^2$  forms

---

<sup>2</sup>As with  $L^2$  regularization, we could regularize the parameters towards a value that is not zero, but instead towards some parameter value  $\mathbf{w}^{(o)}$ . In that case the  $L^1$  regularization would introduce the term  $\Omega(\boldsymbol{\theta}) = \|\mathbf{w} - \mathbf{w}^{(o)}\|_1 = \sum_i |w_i - w_i^{(o)}|$ .