

Regularization

- Many strategies used in machine learning are explicitly designed to reduce the test error, possibly at the expense of increased training error. These strategies are known collectively as regularization.
- In section 5.2.2, we defined regularization as “any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.”
- We almost always find that the best fitting model is a large model that has been regularized appropriately.
- Many regularization approaches are based on limiting the capacity of models by adding a parameter norm penalty $\Omega(\theta)$ to the objective function J .
- We denote the regularized objective function by \tilde{J} :

$$\tilde{J}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\theta), \quad (7.1)$$

where $\alpha \in [0, \infty)$ is a hyperparameter that weights the relative contribution of the norm penalty term, Ω , relative to the standard objective function J .

- Setting α to 0 results in no regularization. Larger values of α correspond to more regularization.

Parameter norm penalties

- We typically choose to use a parameter norm penalty Ω that penalizes only the weights of the affine transformation at each layer and leaves the biases unregularized.
- Each weight specifies how two variables interact. Fitting the weight well requires observing both variables in a variety of conditions.
- Each bias controls only a single variable. This means that we do not induce too much variance by leaving the biases unregularized. The biases typically require less data than the weights to fit accurately. Also, regularizing the bias parameters can introduce a significant amount of underfitting.
- In the context of neural networks, it is sometimes desirable to use a separate penalty with a different α coefficient for each layer of the network.
- Because it can be expensive to search for the correct value of multiple hyperparameters, it is still reasonable to use the same weight decay at all layers just to reduce the size of search space.

L2 parameter regularization

- The $L2$ parameter norm penalty is commonly known as weight decay. This regularization strategy drives the weights closer to the origin by adding a regularization term $\Omega(\theta) = \frac{1}{2} \|\mathbf{w}\|^2$ to the objective function.
- In other academic communities, $L2$ regularization is also known as ridge regression or Tikhonov regularization.
- We can gain some insight into the behavior of weight decay regularization by studying the gradient of the regularized objective function.
- We assume no bias parameter, so θ is just \mathbf{w} . Such a model has the following total objective function:

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y}), \quad (7.2)$$

with the corresponding parameter gradient

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}). \quad (7.3)$$

- To take a single gradient step to update the weights, we perform this update:

$$\mathbf{w} \leftarrow \mathbf{w} - \epsilon (\alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})). \quad (7.4)$$

Written another way, the update is

$$\mathbf{w} \leftarrow (1 - \epsilon \alpha) \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}). \quad (7.5)$$

- The addition of the weight decay term has modified the learning rule to multiplicatively shrink the weight vector by a constant factor on each step before gradient update.
- This describes what happens in a single step. But what happens over the entire course of training?
- We further simplify the analysis by making a quadratic approximation to the objective function in the neighborhood of the value of the weights that obtains minimal unregularized training cost, $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$.
- If the objective function is truly quadratic, as in the case of fitting a linear regression model with mean squared error, then the approximation is perfect.
- The approximation J^\wedge 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 J^\wedge occurs where its gradient

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

is equal to 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 J^\wedge . We use the variable $\tilde{\mathbf{w}}$ to represent the location of the minimum.

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

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

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

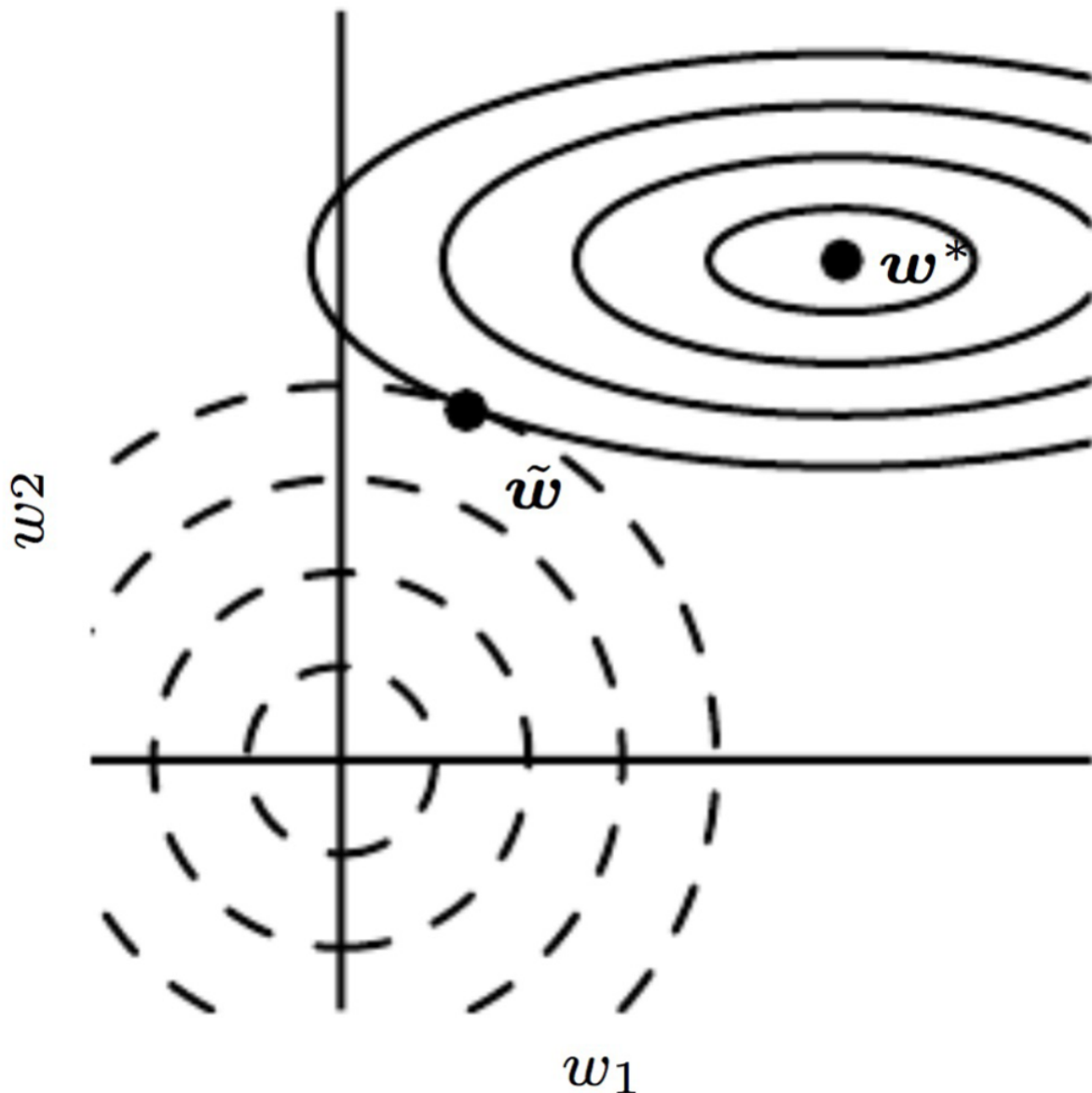
- As α approaches 0, the regularized solution $\tilde{\mathbf{w}}$ approaches \mathbf{w}^* . But what happens as α grows?
- Because \mathbf{H} is real and symmetric, we can decompose it into a diagonal matrix $\mathbf{\Lambda}$ and an orthonormal basis of eigenvectors \mathbf{Q} , such that $\mathbf{H} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$.
- 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 $\lambda_i/(\lambda_i + \alpha)$.
 - 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. Yet components with $\lambda_i \ll \alpha$ will be shrunk to have nearly zero magnitude.
 - Only directions along which the parameters contribute significantly to reducing the objective function are preserved relatively intact. In directions that do not contribute to reducing the objective function, a small eigenvalue of the Hessian tells us that movement in this direction will not significantly increase the gradient.
 - Components of the weight vector corresponding to such unimportant directions are decayed away through the use of the regularization throughout training.
-



- The solid ellipses represent contours of equal value of the unregularized objective. The dotted circles represent contours of equal value of the $L2$ regularizer.
- At the point \tilde{w} , these competing objectives reach an equilibrium. In the first dimension, the eigenvalue of the Hessian of J is small. The objective function does not increase much when moving horizontally away from w^* . Because the objective function does not express a strong preference along this direction, the regularizer has a strong effect on this axis. The regularizer pulls w_1 close to zero.
- In the second dimension, the objective function is very sensitive to movements away from w^* . The corresponding eigenvalue is large, indicating high curvature. As a result, weight decay affects the position of w_2 relatively little.
- For linear regression, the cost function is 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 α to the diagonal. The diagonal entries of this matrix correspond to the variance of each input feature.
- 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.

L^1 regularization

- While L^2 weight decay is the most common form of weight decay, another option is to use L^1 regularization.
- Formally, L^1 regularization on the model parameter \mathbf{w} is

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

that is, as the sum of absolute values of the individual parameters.

- As with L^2 weight decay, L^1 weight decay controls the strength of the regularization by scaling the penalty Ω using a positive hyperparameter α .
- Thus, the regularized objective function $J(\mathbf{w}; \mathbf{X}, \mathbf{y})$ is

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \|\mathbf{w}\|_1 + J(\mathbf{w}; \mathbf{X}, \mathbf{y}), \quad (7.19)$$

with the corresponding gradient (actually, sub gradient)

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \text{sign}(\mathbf{w}) + \nabla_{\mathbf{w}} J(\mathbf{X}, \mathbf{y}; \mathbf{w}), \quad (7.20)$$

where $\text{sign}(\mathbf{w})$ is the sign of \mathbf{w} applied element-wise.

- By inspecting equation 7.20, we can see that the regularization contribution to the gradient no longer scales linearly with each w_i ; instead it is a constant factor with a sign equal to $\text{sign}(w_i)$.

- One consequence of this form of the gradient is that we will not necessarily see clean algebraic solutions to quadratic approximations of $J(\mathbf{X}, \mathbf{y}; \mathbf{w})$ as we did for L_2 regularization.
- Our simple linear model has a quadratic cost function that we can represent via its Taylor series. Alternately, we could imagine that this is a truncated Taylor series approximating the cost function of a more sophisticated model.
- The gradient in this setting is given by

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*), \quad (7.21)$$

where, again, \mathbf{H} is the Hessian matrix of J with respect to \mathbf{w} evaluated at \mathbf{w}^* .

- Because the L_1 penalty does not admit clean algebraic expressions in the case of a fully general Hessian, we will further assume that the Hessian is diagonal, $\mathbf{H} = \text{diag}(H_{1,1}, \dots, H_{n,n})$, where each $H_{i,i} > 0$.
- This assumption holds if the data for the linear regression problem has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.
- Our quadratic approximation of the L_1 regularized objective function decomposes into a sum over the parameters:

$$\hat{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = J(\mathbf{w}^*; \mathbf{X}, \mathbf{y}) + \sum_i \left[\frac{1}{2} H_{i,i} (\mathbf{w}_i - \mathbf{w}_i^*)^2 + \alpha |\mathbf{w}_i| \right]. \quad (7.22)$$

- The problem of minimizing this approximate cost function has an analytical solution (for each dimension i), with the following form:

$$\mathbf{w}_i = \text{sign}(\mathbf{w}_i^*) \max \left\{ |\mathbf{w}_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}. \quad (7.23)$$

- Consider the situation where $\mathbf{w}_i^* > 0$ for all i . There are two possible outcomes:
 1. The case where $\mathbf{w}_i^* \leq \alpha/H_{i,i}$. Here the optimal value of \mathbf{w}_i under the regularized objective is simply $\mathbf{w}_i = 0$. This occurs because the contribution of $J(\mathbf{w}; \mathbf{X}, \mathbf{y})$ to the regularized objective $\hat{J}(\mathbf{w}; \mathbf{X}, \mathbf{y})$ is overwhelmed—in direction i —by the L_1 regularization, which pushes the value of \mathbf{w}_i to zero.
 2. The case where $\mathbf{w}_i^* > \alpha/H_{i,i}$. In this case, the regularization does not move the optimal value of \mathbf{w}_i to zero but instead just shifts it in that direction by a distance equal to $\alpha/H_{i,i}$.
- A similar process happens when $\mathbf{w}_i < 0$, but with the L_1 penalty making \mathbf{w}_i less negative by $\alpha/H_{i,i}$, or 0.
- In comparison to L_2 regularization, L_1 regularization results in a solution that is more sparse. Sparsity in this context refers to the fact that some parameters have an optimal value of zero.
- The sparsity property induced by L_1 regularization has been used extensively as a feature selection mechanism (choosing which subset of the available features should be used).
- The well known LASSO (least absolute shrinkage and selection operator) model integrates an L_1 penalty with a linear model and a least-squares cost function. The L_1 penalty causes a subset of the weights to become zero, suggesting that the corresponding features may safely be discarded.

Dataset augmentation

- The best way to make a machine learning model generalize better is to train it on more data. In practice, the amount of data we have is limited. One way to solve this problem is to create fake data and add it to the training set.
- This approach is easiest for classification. The main task facing a classifier is to be invariant to a wide variety of transformations. We can generate new (x, y) pairs easily just by transforming the x inputs in our training set.
- This approach is not as readily applicable to many other tasks. For example, it is difficult to generate new fake data for a density estimation task unless we have already solved the density estimation problem.
- Dataset augmentation is particularly effective for a specific classification problem: object recognition.
- Operations like translating the training images a few pixels in each direction can often greatly improve generalization, even if the model has already been designed to be partially translation invariant. Many other operations, such as rotating the image or scaling the image, have also proved quite effective.
- One must be careful not to apply transformations that would change the correct class. For example, optical character recognition tasks require recognizing the difference between “b” and “d” and the difference between “6” and “9,” so horizontal flips and 180° rotations are not appropriate.
- Injecting noise in the input to a neural network can also be seen as a form of data augmentation.
- For many classification and even some regression tasks, the task should still be possible to solve even if small random noise is added to the input. One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs.
- Noise injection also works when the noise is applied to the hidden units, which can be seen as doing dataset augmentation at multiple levels of abstraction.
- Dropout can be seen as a process of constructing new inputs by multiplying by noise.

Noise robustness

- For some models, the addition of noise with infinitesimal variance at the input of the model is equivalent to imposing a penalty on the norm of the weights. In the general case, noise injection can be much more powerful than simply shrinking the parameters, especially when the noise is added to the hidden units.
- Another way that noise has been used is by adding it to the weights. This technique has been used primarily in the context of RNN.
- Assume that with each input presentation we include a random perturbation $\epsilon W \sim N(\epsilon; \mathbf{0}, \eta I)$ of the network weights. We use a standard l -layer MLP, and denote the perturbed model as $y_Y(x)$. The objective function is $4W$
- For small η , the minimization of J with added weight noise (with covariance ηI) is equivalent to minimization of J with an additional regularization term:
$$\eta \mathbb{E} \sum (x, y) \nabla W y_Y(x)!$$