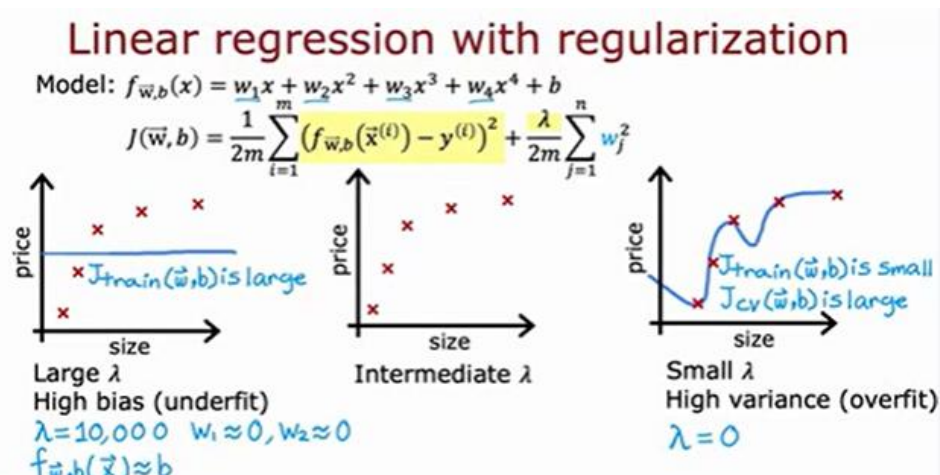


Regularization and bias/variance

variance of your learning algorithm, and therefore its overall performance. In this video, let's examine how regularization, specifically the choice of the regularization parameter λ , affects the bias and variance, and thus the overall performance of the algorithm. This will be useful when you want to choose a good value for λ for your algorithm's regularization parameter.

Let's take a look. In this example, I'll use a polynomial of order four, but we'll fit this model using regularization, where the value of λ is the regularization parameter that controls the extent to which you trade off between keeping the parameters w small and fitting the training data correctly. Let's start with the example of setting λ to a very large value (VALUE). Let's say λ equals 10,000.

If you did that, you'd end up with a model that looks something like this. Because if λ were very, very large, then the algorithm would be strongly incentivized to keep those parameters w very small. And so you end up with w_1, w_2 , really all those parameters being very close to zero. And so the model ends up being that f of x is just approximately b , a constant value, which is why you end up with a model like this. And this model clearly has a high bias and it doesn't fit the training data because it doesn't even work well on the training set. And J_{TRAIN} is large. Let's look at the other extreme.

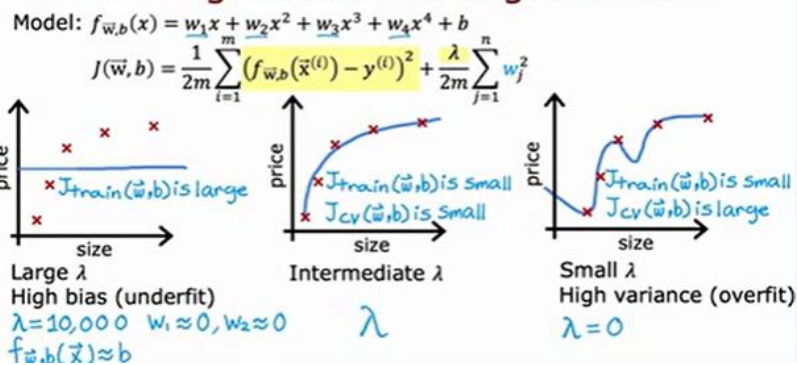


Let's say you set λ to a very small value (VALUE). So with a small value of λ , in fact, let's go to the extreme by setting λ equal to zero. With this choice of λ , there is no regularization, and so we are simply fitting a four-part polynomial without regularization, and you end up with this. And so on.

Thank you. We saw earlier that the data is overfitted, and so what we saw previously is that when you have a model like this, the J_{train} is small, but J_{CV} is much larger than the J_{train} , or J_{CV} is large, and this indicates that we have a large variance relative to the physical data.

And this would be the case if you had an intermediate value of λ , not really far above 10,000, but not as small as zero, that you would hopefully get a model that looks like this, which is correct and fits the data well with a small J_{train} and a small J_{CV} . So, if you are trying to decide what a good λ value to use for the regularization parameter is, cross-validation also provides a way to do this.

Linear regression with regularization



Let's see how we could achieve this. And just as a reminder, the problem we're addressing is that if you're fitting a four-part polynomial— that is, the model—and using regularization, how can you choose a good value for lambda? This would be a similar procedure to the one you saw for choosing the degree of polynomial D using cross-validation.

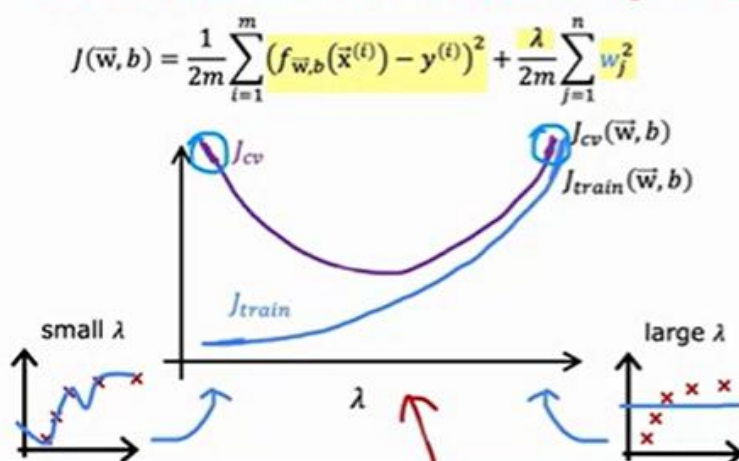
More precisely, let's say we're trying to fit a model using lambda equal to zero, and thus we would minimize the cost function using lambda equal to zero, and we would end up with some parameters w_1 , b_1 , and you can then calculate the cross-validation error j_{cv} of w_1 , b_1 . Now, let's try a different value for lambda. Let's say you're trying lambda equal to 0.01.

Then, minimizing the cost function gives you a second set of parameters, w_2 , b_2 . And you can also see how well this works on the cross-validation set, and so on. Let's continue trying other Lambda values. In this example, I'll try doubling it so that Lambda equals 0.02, which would give you J_{CV} of w_3 , b_3 , and so on.

Let's double the gain , and double the gain again. After doubling it a few times, you get a Lambda approximately equal to 10. This will give you the parameters w_{12} , b_{12} , and J_{CV} , w_{12} , b_{12} . By trying a wide range of possible Lambda values, adjusting the parameters using these different regularization parameters, and then evaluating the performance on the cross-validation set, you can then try to choose the best value for the regularization parameter.

Specifically, if in this example you find that j_{cv} of w_5 b_5 has the lowest value of all these different cross-validation errors, you can then decide to choose this value for lambda and therefore use w_5 b_5 as chosen parameters and finally if you wish to report and an estimate of the generalization error.

Bias and variance as a function of regularization parameter λ



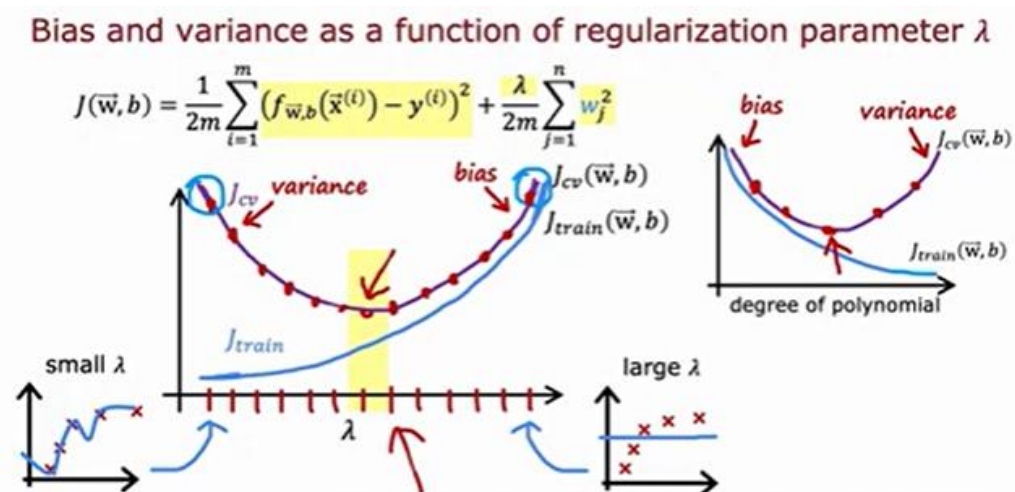
You would then report the tested error, J_{test} of W_5 , B_5 . To further refine our intuition about what this algorithm does, let's examine how the training error and cross-validation error vary depending on the lambda parameter. So, in this figure, I've changed the x-axis again.

Notice that the x-axis here is labeled with the value of the regularization parameter lambda, and if we look at the extreme of lambda equal to zero here on the left, this corresponds to not using regularization. This is where we ended up with this very wavy curve if lambda was small or even zero.

In this case, we have a high-variance model, and therefore J_{TRAIN} will be small and J_{CV} will be large because it performs very well on the training data but performs much less well on the cross-validation data. This extreme right-hand side, with very large lambda values—let's say lambda equals 10,000—ends up fitting a model that looks like this. This therefore has a high bias; it underfits the data, and it turns out that J_{TRAIN} will be high, and J_{CV} will also be high. And in fact, if you look at how J_{train} varies with lambda, you'll see that J_{train} will increase like this.

Because in the optimization cost function, the larger Lambda is, the more the algorithm tries to keep W squared small—that is, the more weight it gives to this regularization term, and therefore less attention is paid to performing well on the training set. This term on the left is J_{TRAIN} , so the more it tries to keep the parameters small, the less well it does to minimize the training error.

That's why, as Lambda increases, the AI training error tends to increase like this. And now, what about the cross-validation error? It turns out that the cross-validation error will look like this. Because we've seen that if Lambda is too small or too large, it doesn't perform well on the cross-validation set. It's either too large here on the left, or it's too small here on the right.



And there will be an intermediate value of Lambda that will make the algorithm work optimally. What cross-validation does is try a large number of different values of Lambda. This is what we saw on the last slide. Try Lambda equal to 0, Lambda equal to 0.01, Lambda equal to 0.02.

Try many different values of Lambda and evaluate the cross-validation error at several of these different points. Then, hopefully, choose a value that has a low cross-validation error, and that will hopefully provide a good model for your application.

If you compare this diagram to the one we had in the previous video where the horizontal axis was the degree of the polynomial, these two diagrams look somewhat similar—not mathematically and not formally, but they are somewhat like mirror images of each other. And that's because when you fit a degree of polynomial, the left part of this curve corresponds to an underfit and a high bias.

The right-hand side corresponded to overfitting and high variance. Whereas in this one, the high variance was on the left and the high bias was on the right. But that's why these two images are somewhat like mirror images of each other. But in both cases, cross-validation—evaluating different values—can help you choose a good t -value or a good λ value.

So that's how the choice of the λ regularization parameter affects the bias, variance, and overall performance of your algorithm. And you've also seen how you can use cross-validation to make a good choice for the λ regularization parameter.

So far, we've talked about how a high training set error, a high J_{TRAIN} , indicates high bias, and how a high cross-validation error, J_{CV} , especially if it's much higher than J_{TRAIN} , indicates a variance problem. But what do the words "high" or "much higher" actually mean?

Let's look at this in the next video, where we'll see how you can examine the J_{Train} and J_{CV} numbers and judge whether they are high or low. It turns out that a further refinement of these ideas—establishing a baseline performance level—will allow your learning algorithm to much more easily assess these J_{Train} and J_{CV} numbers and determine if they are high or low.