

Advice for applying Machine Learning

Contents

Machine Learning diagnostics.....	3
Model selection.....	4
Biasn and Variance.....	6
Bias,variance and regularization.....	7
Establishing a baseline level of performance.....	8
Learning curves.....	9
Bias/Variance applied to neural networks.....	11
Error analysis.....	14
Dealing with skewed datasets.....	14

Summary

So after implementation, the algorithm is making unacceptably large errors: what to do next? Two key concepts are high bias and high variance, which are explained in detail in the following pages.

But, as a summary, the following rectangle contains the basic strategies to tackle this:

To fix high bias/underfitting:

- **Try getting additional features:** For example, if you try to predict the price of a house just by the size, you might be ignoring other crucial factors like zone, age of the house, etc, and your model will have high bias. By adding extra features, in this case the data will fit much better.
- **Make your model more complex:**
 - **Try adding polynomial features:** Is similar to adding additional features - a linear function might not be able to fit the data well, and by adding polynomial features the function will become more complex and fit the training data better.
 - **Try decreasing the regularization parameter λ :** It will increase the size of the weights and lead to a more complex function able to fit the data better.

To fix high variance/overfitting:

- **Get more training examples:** Will help making the function more complex, improving how it fits the training set while also generalizing better.
- **Simplify your model:**
 - **Try a smaller set of features:** Sometimes having too many features makes the algorithm try to fit very complex functions, increasing the variance. If you suspect your algorithm has plenty of features which are not helpful or redundant, then eliminating them will remove the flexibility of the algorithm to fit the data.
 - **Try increasing the regularization parameter λ :** This will reduce the weights, making the function smoother and reducing its complexity and high variance.

"The concept of high bias and high variance takes a short time to learn, but takes a whole life to master"

Let's explore this further!

Machine Learning diagnostics

We split the data into 2 sets: training and test.

With **linear regression**, we fit the model with the training set (and with regularization, if applicable) and compute both, the training error and test error (note: when computing the error, the regularization is not included):

Train/test procedure for linear regression (with squared error cost)

Fit parameters by minimizing cost function $J(\vec{w}, b)$

$$\rightarrow J(\vec{w}, b) = \left[\frac{1}{2m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m_{\text{train}}} \sum_{j=1}^n w_j^2 \right]$$

Compute test error:

$$J_{\text{test}}(\vec{w}, b) = \frac{1}{2m_{\text{test}}} \left[\sum_{i=1}^{m_{\text{test}}} (f_{\vec{w},b}(\vec{x}_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2 \right] \quad \cancel{\sum_{j=1}^n w_j^2}$$

Compute training error:

$$J_{\text{train}}(\vec{w}, b) = \frac{1}{2m_{\text{train}}} \left[\sum_{i=1}^{m_{\text{train}}} (f_{\vec{w},b}(\vec{x}_{\text{train}}^{(i)}) - y_{\text{train}}^{(i)})^2 \right]$$

With **classification**, we can apply very similar ideas:

Train/test procedure for classification problem

0/1

Fit parameters by minimizing $J(\vec{w}, b)$ to find \vec{w}, b

E.g.,

$$J(\vec{w}, b) = -\frac{1}{m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} \left[y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m_{\text{train}}} \sum_{j=1}^n w_j^2$$

Compute test error:

$$J_{\text{test}}(\vec{w}, b) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left[y_{\text{test}}^{(i)} \log(f_{\vec{w},b}(\vec{x}_{\text{test}}^{(i)})) + (1 - y_{\text{test}}^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}_{\text{test}}^{(i)})) \right]$$

Compute train error:

$$J_{\text{train}}(\vec{w}, b) = -\frac{1}{m_{\text{train}}} \sum_{i=1}^{m_{\text{train}}} \left[y_{\text{train}}^{(i)} \log(f_{\vec{w},b}(\vec{x}_{\text{train}}^{(i)})) + (1 - y_{\text{train}}^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}_{\text{train}}^{(i)})) \right]$$

Log loss will work well, although in classification is more commonly used the proportion of the data that has been misclassified:

Train/test procedure for classification problem

We want the test error to be as low as possible, in which case the model generalizes well.

fraction of the test set and the fraction of the train set that the algorithm has misclassified.

count $\hat{y} \neq y$

$$\hat{y} = \begin{cases} 1 & \text{if } f_{\vec{w},b}(\vec{x}^{(i)}) \geq 0.5 \\ 0 & \text{if } f_{\vec{w},b}(\vec{x}^{(i)}) < 0.5 \end{cases}$$

$J_{\text{test}}(\vec{w}, b)$ is the fraction of the test set that has been misclassified.

$J_{\text{train}}(\vec{w}, b)$ is the fraction of the train set that has been misclassified.

Model selection

Dividing into training and test sets when selecting a model (for example, the degree of a polynomial) is a flawed approach, as our test error will be optimistic compared with the generalization error (this is, the error for completely new or unknown data). This is because out test data was used to select out model, which means that this data fits particularly well with this polynomial, and some skewedness is inevitable. Andrew Ng explains it as follows:

Model selection (choosing a model)

$d=1$ 1. $f_{\bar{w},b}(\vec{x}) = w_1x + b \rightarrow w^{<1>, b^{<1>} \rightarrow J_{test}(w^{<1>, b^{<1>})$
 $d=2$ 2. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + b \rightarrow w^{<2>, b^{<2>} \rightarrow J_{test}(w^{<2>, b^{<2>})$
 $d=3$ 3. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + w_3x^3 + b \rightarrow w^{<3>, b^{<3>} \rightarrow J_{test}(w^{<3>, b^{<3>})$
 \vdots
 $d=10$ 10. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + \dots + w_{10}x^{10} + b \rightarrow J_{test}(w^{<10>, b^{<10>})$
 Choose $w_1x + \dots + w_5x^5 + b$ $d=5$ $J_{test}(w^{<5>, b^{<5>})$

How well does the model perform? Report test set error $J_{test}(w^{<5>, b^{<5>})$?

The problem: $J_{test}(w^{<5>, b^{<5>})$ is likely to be an optimistic estimate of generalization error (ie. $J_{test}(w^{<5>, b^{<5>}) < \text{generalization error}$). Because an extra parameter d (degree of polynomial) was chosen using the test set.

w, b are overly optimistic estimate of generalization error on training data.

Solution? Add another set: we split the data into training, cross-validation and test sets (Andrew does 60% training, 20% cv, 20% test).

Training/cross validation/test set

size	price		validation set development set dev set	
2104	400	→ training set 60%	→ $(x^{(1)}, y^{(1)})$ \vdots $(x^{(m_{train})}, y^{(m_{train})})$	$m_{train} = 6$
1600	330			
2400	369			
1416	232			
3000	540			
1985	300			
1534	315	→ cross validation 20%	→ $(x_{cv}^{(1)}, y_{cv}^{(1)})$ \vdots $(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$	$m_{cv} = 2$
1427	199			
1380	212	→ test set 20%	→ $(x_{test}^{(1)}, y_{test}^{(1)})$ \vdots $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$	$m_{test} = 2$
1494	243			

Andrew explains the cross-validation set is also called in other ways; he likes dev set because is shorter.

Training/cross validation/test set

Training error: $J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})^2 \right]$

Cross validation error: $J_{cv}(\vec{w}, b) = \frac{1}{2m_{cv}} \left[\sum_{i=1}^{m_{cv}} (f_{\vec{w},b}(\vec{x}_{cv}^{(i)}) - y_{cv}^{(i)})^2 \right]$ (validation error, dev error)

Test error: $J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w},b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right]$

We then use:

1. Training data: to fit all the different models
2. Validation/Dev data: To choose the model with lowest dev error.
3. Test data: The test error is a fair value for our generalization error (the one we expect for new data)

Model selection

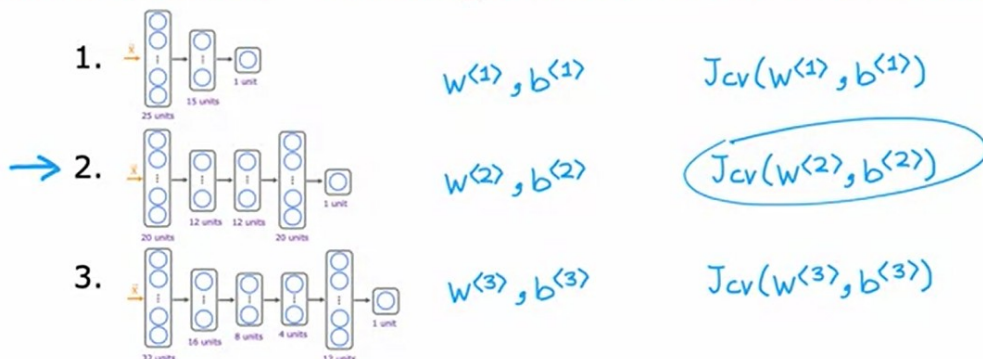
$d=1$	1. $f_{\vec{w},b}(\vec{x}) = w_1x + b$	$w^{<1>}, b^{<1>}$	$\rightarrow J_{cv}(w^{<1>}, b^{<1>})$
$d=2$	2. $f_{\vec{w},b}(\vec{x}) = w_1x + w_2x^2 + b$		$\rightarrow J_{cv}(w^{<2>}, b^{<2>})$
$d=3$	3. $f_{\vec{w},b}(\vec{x}) = w_1x + w_2x^2 + w_3x^3 + b$		\vdots
\vdots	\vdots		\vdots
$d=10$	10. $f_{\vec{w},b}(\vec{x}) = w_1x + w_2x^2 + \dots + w_{10}x^{10} + b$		$J_{cv}(w^{<10>}, b^{<10>})$

\rightarrow Pick $w_1x + \dots + w_4x^4 + b$ ($J_{cv}(w^{<4>}, b^{<4>})$)

Estimate generalization error using test the set: $J_{test}(w^{<4>}, b^{<4>})$

This can be applied for other models:

Model selection – choosing a neural network architecture



Pick $w^{<2>}, b^{<2>}$

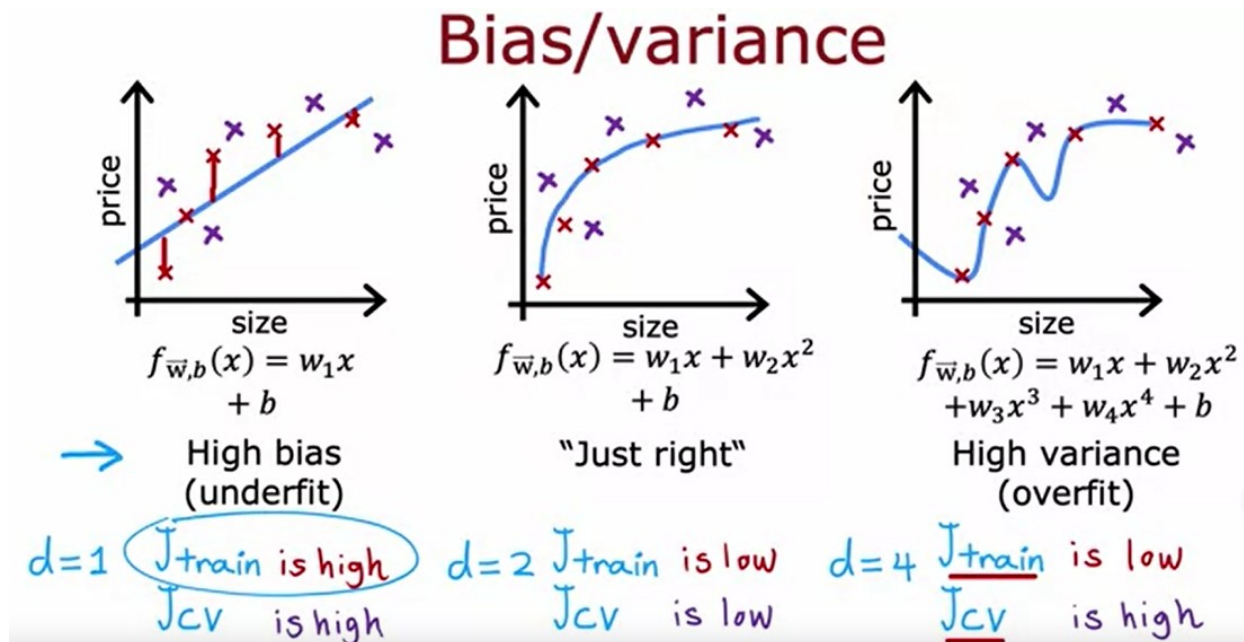
Estimate generalization error using the test set: $J_{test}(w^{<2>}, b^{<2>})$

Bias and Variance

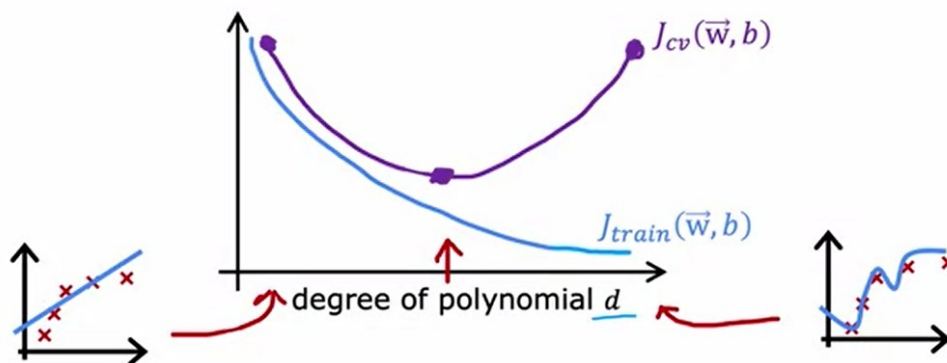
This section is presented with the assumption there is no regularization.

High bias: the model does not perform well even in the training set.

High variance: the model performs very well with the training set, but fails at generalizing.



Understanding bias and variance



The lowest point of the J_{cv} curve (the "valley" spot) is where we find the best performing polynomial. Before that, the model is prone to high bias (it underfits), and after it to high variance (it overfits).

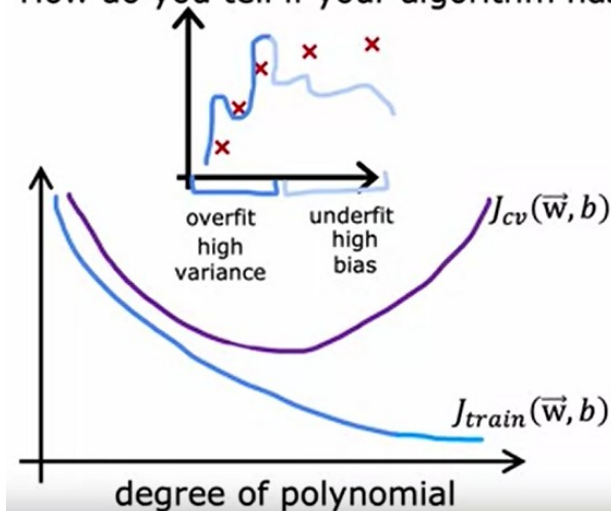
Sometimes you can have simultaneously high bias and high variance; you won't see this happen as much in linear regression, but if you are training neural networks, there are some applications where you might have high bias and high variance. To recognize this, you see the training error is high, but the cv error is even much higher. An explanation is that for

some data the model clearly underfits but, for the other, it overfits (which is very unlikely to happen in linear regression).

This is explained in the next image, with a small picture replicating what in linear regression would mean simultaneous high bias and high variance (a bit of an artificial example):

Diagnosing bias and variance

How do you tell if your algorithm has a bias or variance problem?



High bias (underfit)

→ J_{train} will be high
($J_{train} \approx J_{cv}$)

High variance (overfit)

→ $J_{cv} \gg J_{train}$
(J_{train} may be low)

High bias and high variance

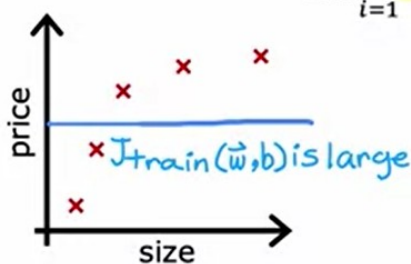
→ J_{train} will be high
and $J_{cv} \gg J_{train}$

Bias, variance and regularization

Linear regression with regularization

Model: $f_{\vec{w},b}(x) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

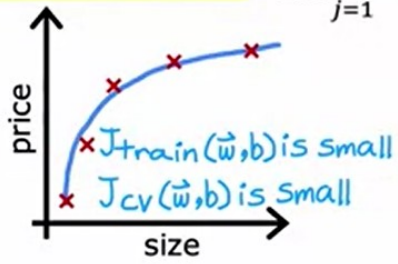


Large λ

High bias (underfit)

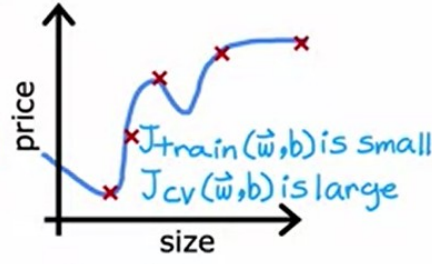
$\lambda = 10,000$ $w_1 \approx 0, w_2 \approx 0$

$f_{\vec{w},b}(\vec{x}) \approx b$



Intermediate λ

λ



Small λ

High variance (overfit)

$\lambda = 0$

Choosing the regularization parameter λ

Model: $f_{\vec{w},b}(x) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$

- \rightarrow 1. Try $\lambda = 0 \rightarrow \min_{\vec{w},b} J(\vec{w},b) \rightarrow w^{<1>,b^{<1>} \rightarrow J_{cv}(w^{<1>,b^{<1>})$
 \rightarrow 2. Try $\lambda = 0.01 \rightarrow w^{<2>,b^{<2>} \rightarrow J_{cv}(w^{<2>,b^{<2>})$
 \rightarrow 3. Try $\lambda = 0.02 \rightarrow J_{cv}(w^{<3>,b^{<3>})$
 \rightarrow 4. Try $\lambda = 0.04 \rightarrow J_{cv}(w^{<4>,b^{<4>})$
 \rightarrow 5. Try $\lambda = 0.08 \rightarrow J_{cv}(w^{<5>,b^{<5>})$
 \vdots
 \rightarrow 12. Try $\lambda \approx 10 \rightarrow w^{<12>,b^{<12>} \rightarrow J_{cv}(w^{<12>,b^{<12>})$

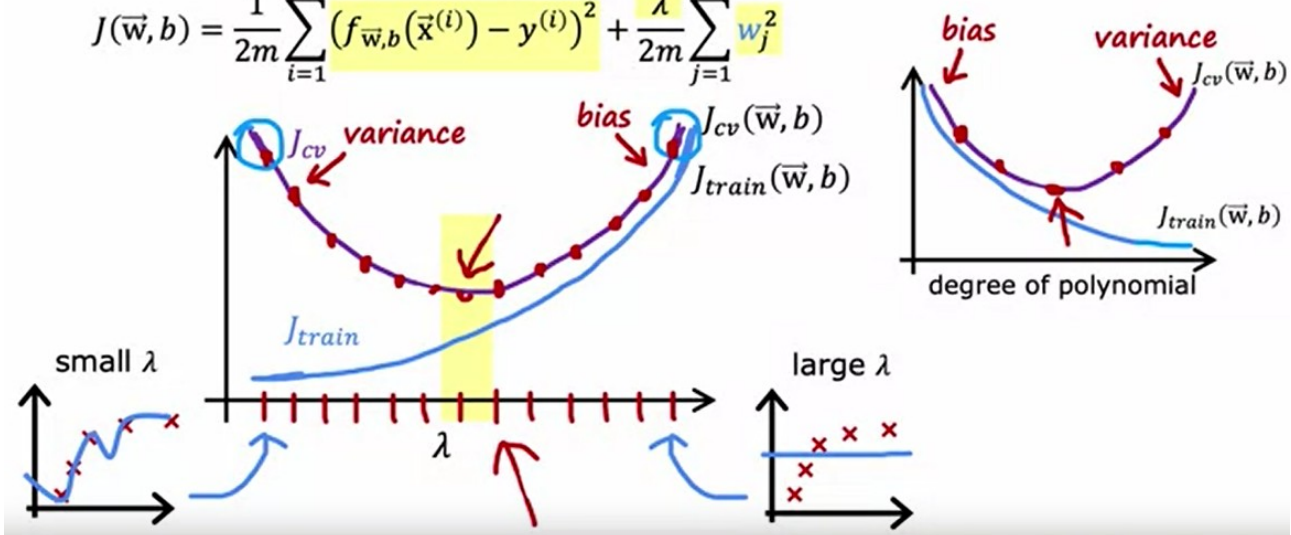
Pick $w^{<5>,b^{<5>}$

Report test error: $J_{test}(w^{<5>,b^{<5>})$

So, how does the training and cv error vary as a function of the parameter λ ?

Bias and variance as a function of regularization parameter λ

$$J(\vec{w},b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$



When λ is very small the model has high variance and when is big it has high bias; it looks similar to a mirrored version of the plot for the optimal degree of polynomial.

Establishing a baseline level of performance

Andrew explains how a first approach to evaluate the performance of a model is to compare its error to the one a human expert could make. For example, in speech recognition:

Speech recognition example



Human level performance	: 10.6%	
Training error J_{train}	: 10.8%	$\uparrow 0.2\%$
Cross validation error J_{cv}	: 14.8%	$\uparrow 4.0\%$



We could think that, in this example, a model with a training error of 10.8% is very high, but actually humans have an error rate of 10.6% when transcribing these messages (for example, messages with lots of noise, or cut in the middle, etc). But then, if the cv error is of 14.8%, that is actually quite high; in this case, we might have a high variance problem.

We need to establish a baseline level of performance, which can be achieved by:

Establishing a baseline level of performance

What is the level of error you can reasonably hope to get to?

- • Human level performance
- • Competing algorithms performance
- • Guess based on experience

So coming back to the previous example:

Bias/variance examples

Baseline performance	: 10.6%	↑ 0.2%	10.6%	↑ 4.4%	10.6%	↑ 4.4%
Training error (J_{train})	: 10.8%	↓ 0.2%	15.0%	↓ 4.4%	15.0%	↓ 4.4%
Cross validation error (J_{cv})	: 14.8%	↓ 4.0%	15.5%	↓ 0.5%	19.7%	↓ 4.7%
			high variance	high bias	high bias	high variance

So instead of asking ourselves: is the error high? We want to be able to ask: **is the error high compared to human performance/competition results? (for example)**

Learning curves

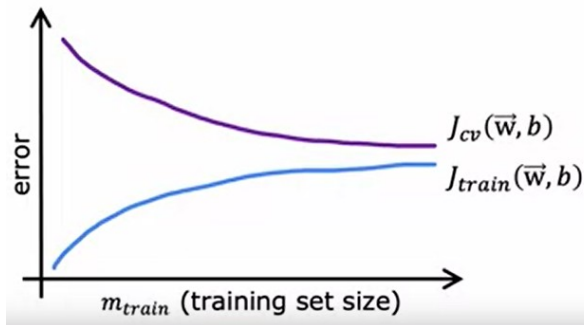
When you have a very small number of training examples, is relative easy to get zero or very small training error, but as the examples increase it becomes more difficult to fit all those examples, so the error will increase.

On the other side, the cv error will typically be higher than the training error, as the parameters are expected to be at least a bit better for the training examples. When plotting this, we will have a curve similar to this one:

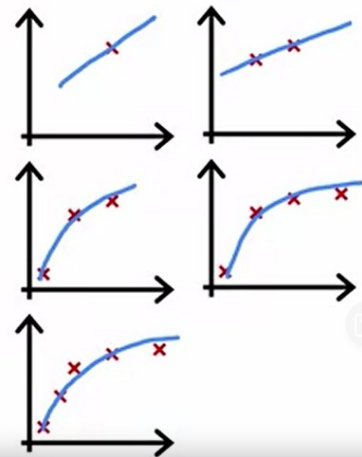
Learning curves

J_{train} = training error

J_{cv} = cross validation error

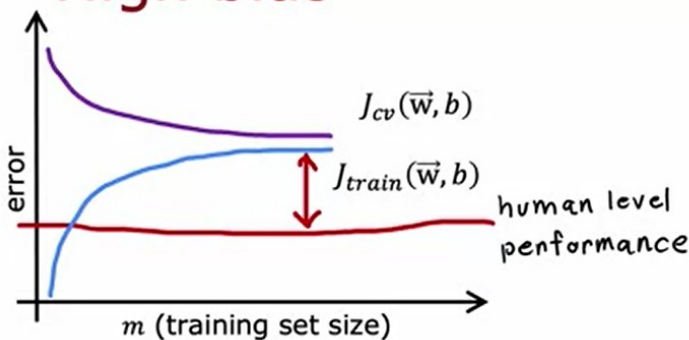


$$f_{\bar{w},b}(x) = w_1x + w_2x^2 + b$$



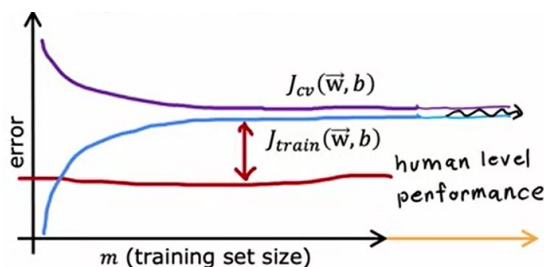
This curve looks like this for high bias:

High bias



The training error and cv error are both much higher than the expected performance, with both being quite close to each other. Generally we would need a more complex function or strategy to improve performance.

One serious mistake we can do in this case is to, in order to improve the performance of the model, increase the training size... which won't affect performance as this is what will happen:

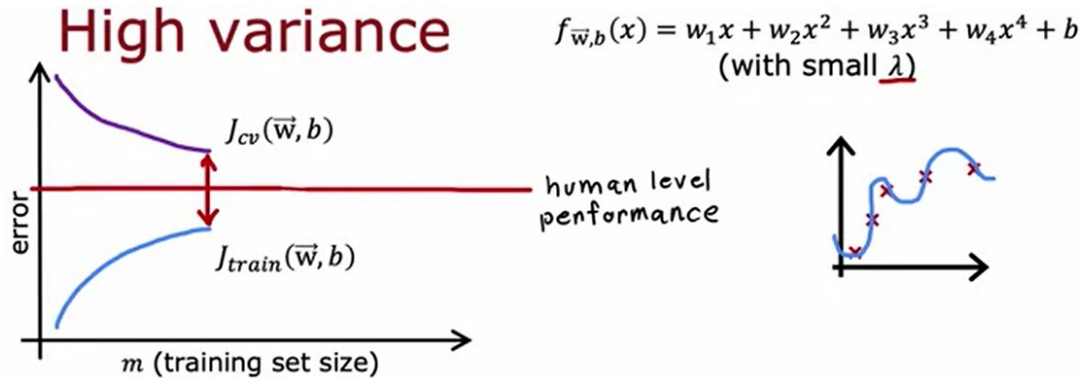


Both errors have reached a plateau and they will not just go down by increasing the examples.

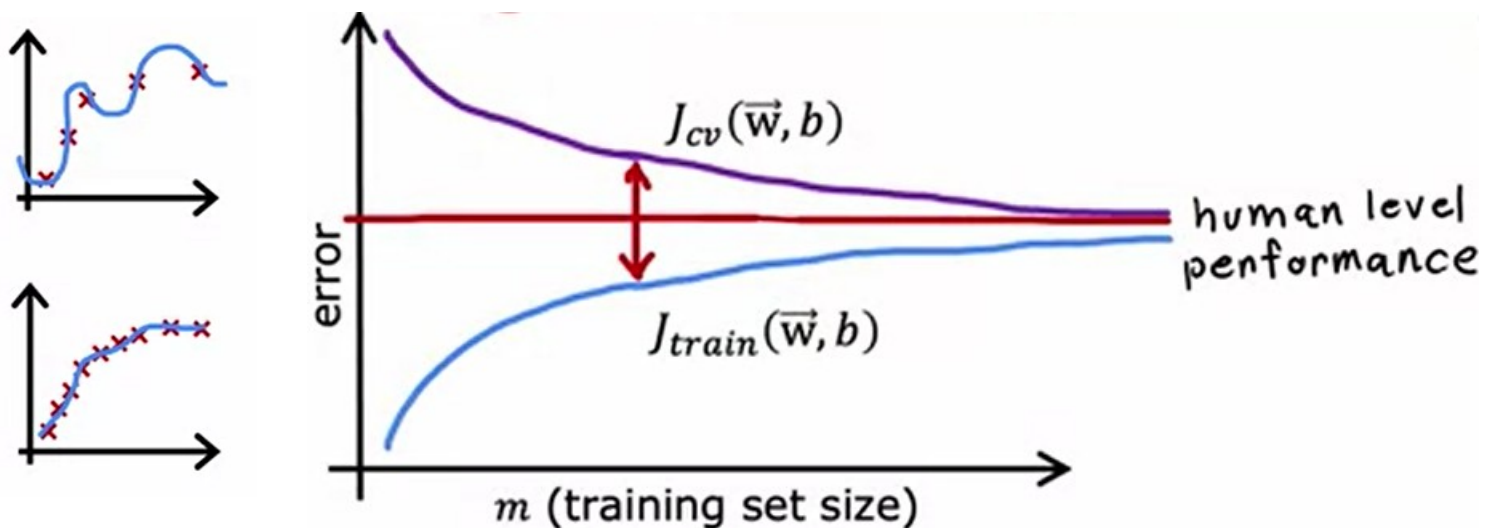
If your algorithm have high bias, bringing more examples will not improve performance by itself.

What stands out in the high variance case is how the cv error will be much higher than the training error:

High variance



In the case of high variance, increasing the training set can help a lot: the training error will go up, while the cv error will hopefully go down:



One important downside of learning curves is that they are computationally expensive, so they are not used that often; even so, Andrew encourages to keep this picture in mind to imagine what your algorithm is doing and what high bias and high variance looks like.

Bias/Variance applied to neural networks

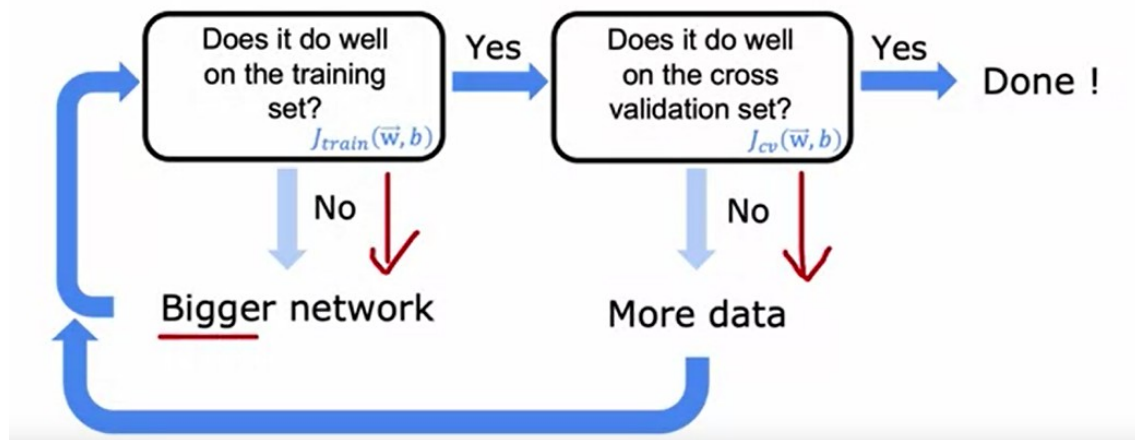
One of the things that explains the success of neural networks is how, when they are combined with big data, high bias and variance also becomes addressed.

Large neural networks, when trained on small to moderate size dataset, are low bias machines. If you make your neural network large enough, you can almost always fit your training set well.

A general routine to address high variance and high bias in neural networks is as follows:

Neural networks and bias variance

Large neural networks are low bias machines



One problem with this approach is that as neural networks become larger, they also become computationally very expensive; the use of GPUs for speeding up the calculations has been crucial in the development of these models, but eventually neural networks can become so large that are not feasible.

Another problem is if at some point you can not get more data.

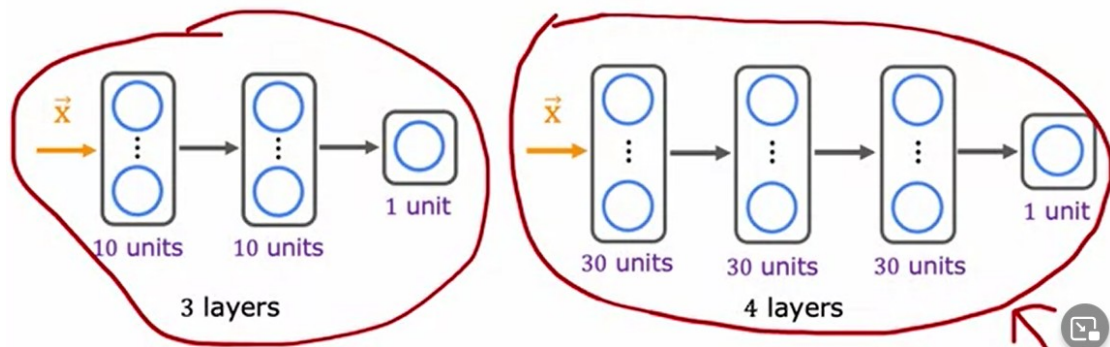
But the reason behind the rise of deep learning all these previous years is that with big enough data and computational power you can have very good performance in plenty of applications.

One possible question is: what if my neural network is too big? Will that create a high variance problem?

It turns out that a large network with well chosen regularization will usually do as well or better than a smaller one.

It almost never hurts to go to a larger neural network so long as you regularize appropriately. Although the computational cost will increase.

Neural networks and regularization



A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{\text{all weights } \mathbf{W}} (w^2) \quad b$$

Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

Error analysis

By manually examining a set of examples that the algorithm is misclassifying, it will often create inspiration about what to try next. The strategy is:

1. Get the data with errors:
 - Classification problems: Missclassified examples.
 - Regression problems: Examples with an error over a certain threshold.
2. If the error examples are too many, or you don't have enough time to go through them -> choose a random set of a size more manageable.
3. Manually go through all the examples and categorize them based on common traits.
4. Analyze the groups:
 - Focus in the ones with lots of errors (don't invest time if, let's say, there is a group with 2 examples out of 1000 total error examples)
 - Take into consideration that some of these categories can be overlapping, or what is the same: not mutually exclusive.

This might encourage us to:

- Look for new features.
- Look for new examples more focused in a particular trait.
- See that certain type of errors are too rare to try to fix them.

Dealing with skewed datasets

If you have in your data, for example, just 0.5% positive cases (let's say, a rare disease), then an algorithm that just predicts $y = 0$ always will have a 99.5% accuracy and 0.5% error. So if you have an algorithm with 1% error, or 1.2%, or 1.6%, it is very difficult to actually know which one performs better.

How to deal with this? We use precision/recall as indicators.

Precision/recall

$y = 1$ in presence of rare class we want to detect.

Actual Class		1	0
Predicted Class	1	True positive 15	False positive 5
	0	False negative 10	True negative 70
		↓ 25	↓ 75

Precision:

(of all patients where we predicted $y = 1$, what fraction actually have the rare disease?)

$$\frac{\text{True positives}}{\text{\#predicted positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False pos}} = \frac{15}{15 + 5} = 0.75$$

Recall:

(of all patients that actually have the rare disease, what fraction did we correctly detect as having it?)

$$\frac{\text{True positives}}{\text{\#actual positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}} = \frac{15}{15 + 10} = 0.6$$

With this metrics, the algorithm that predicts $y = 0$ for all examples would have recall and precision 0 - confirming it as an extremely poor model.

High precision means that if it diagnoses a patient with the rare disease, probably the patient does have it and is an accurate diagnosis.

High recall means that if there is a patient with that rare disease probably the algorithm will correctly diagnose him/her.

We would like an algorithm to have very high values for both, but actually there is a tradeoff between them.

If we want to predict, as the image says, $y = 1$ only when we are very confident (this means, we aim for high precision), increasing the threshold of the logistic regression will increase the precision. This will reduce the recall value, as we might miss more positive cases.

Trading off precision and recall

Logistic regression: $0 < f_{\vec{w},b}(\vec{x}) < 1$
 → Predict 1 if $f_{\vec{w},b}(\vec{x}) \geq \underline{0.5}$ ~~0.7~~ 0.9
 → Predict 0 if $f_{\vec{w},b}(\vec{x}) < \underline{0.5}$ ~~0.7~~ 0.9

$$\text{precision} = \frac{\text{true positives}}{\text{total predicted positive}}$$

$$\text{recall} = \frac{\text{true positives}}{\text{total actual positive}}$$

Suppose we want to predict $y = 1$ (rare disease) only if very confident.

higher precision, lower recall

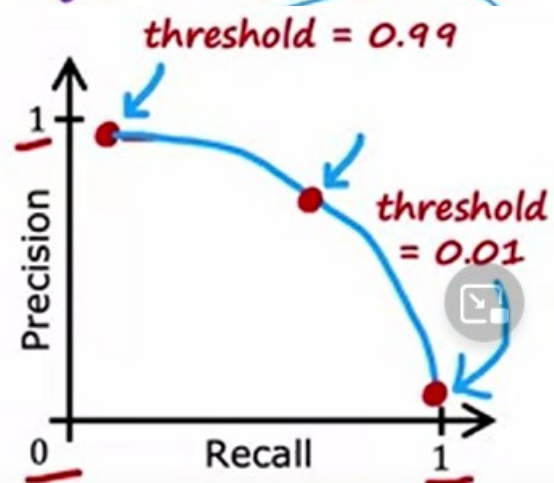
On the other hand, if we want to void missing any patient with the disease, we can do the opposite: reduce the threshold. This will increase the recall while reducing the precision:

Logistic regression: $0 < f_{\vec{w},b}(\vec{x}) < 1$
 → Predict 1 if $f_{\vec{w},b}(\vec{x}) \geq \underline{0.5}$ ~~0.7~~ ~~0.9~~ 0.3
 → Predict 0 if $f_{\vec{w},b}(\vec{x}) < \underline{0.5}$ ~~0.7~~ ~~0.9~~ 0.3

Suppose we want to avoid missing too many case of rare disease (when in doubt predict $y = 1$)

lower precision, higher recall

This is how the tradeoff looks depending on the threshold of the classification algorithm:



To choose an algorithm depending on these 2 values can be tricky; to simplify, we could do the average, but that doesn't work well: we can use the F1 score:

F1 score

How to compare precision/recall numbers?

	Precision (P)	Recall (R)	Average	F ₁ score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.501	0.0392

`print("y=1")`

~~Average = $\frac{P+R}{2}$~~

$$F_1 \text{ score} = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)} = 2 \frac{PR}{P+R}$$

Harmonic mean

The F1 score is also called Harmonic mean, which is a way of calculating an average that emphasizes the smaller values more. We choose the highest F1 score, which is the one with the smaller tradeoff.