

10. Advice for applying Machine Learning:

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices.

$$\rightarrow J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

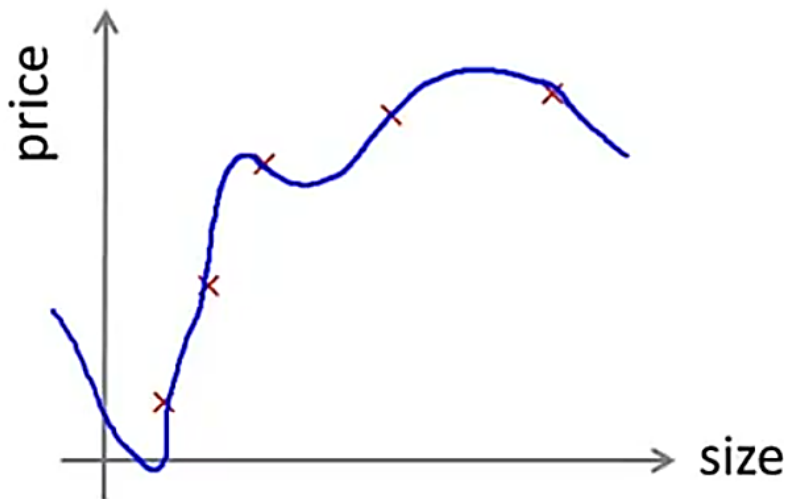
- \rightarrow - Get more training examples
- Try smaller sets of features $x_1, x_2, x_3, \dots, x_{100}$
- \rightarrow - Try getting additional features
- Try adding polynomial features ($\underline{x_1^2}, \underline{x_2^2}, \underline{x_1 x_2}$, etc.)
- Try decreasing λ
- Try increasing λ

Machine learning diagnostic:

Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.

EVALUATING A HYPOTHESIS:



$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Due to overfitting, the learned hypothesis may fail to generalize to new examples not in training set

One way to evaluate the hypothesis:

To evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: **a training set and a test set.**

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	} Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	
<hr/>		
1427	199	} Test set
1380	212	
1494	243	

70%

30%

$$\begin{pmatrix} (x^{(1)}, y^{(1)}) \\ (x^{(2)}, y^{(2)}) \\ \vdots \\ (x^{(m)}, y^{(m)}) \end{pmatrix}$$

$$\begin{pmatrix} (x_{test}^{(1)}, y_{test}^{(1)}) \\ (x_{test}^{(2)}, y_{test}^{(2)}) \\ \vdots \\ (x_{test}^{(m_{test})}, y_{test}^{(m_{test})}) \end{pmatrix}$$

$m_{test} = \text{no. of test example}$
 $(x_{test}^{(i)}, y_{test}^{(i)})$

- Train the model only using the 70% of input dataset, and use the remaining 30% as test set: and check if the predictions match the given values
- When dividing the dataset into training set and test set, make sure the data is not ordered in any way. Randomize it. Random is better.

Procedure:

Training/testing procedure for linear regression

- Learn parameter θ from training data (minimizing training error $J(\theta)$) 70%

- Compute test set error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left(\underbrace{h_{\theta}(x_{\text{test}}^{(i)})}_{\uparrow} - \underbrace{y_{\text{test}}^{(i)}}_{\uparrow} \right)^2$$

Training/testing procedure for logistic regression

- Learn parameter θ from training data
- Compute test set error:

$$\rightarrow J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$$

- Misclassification error (0/1 misclassification error):

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq \underline{0.5}, y = \underline{0} \\ & \text{or if } h_{\theta}(x) < \underline{0.5}, y = \underline{1} \end{cases} \text{ error}$$

$$0 \text{ otherwise}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

Err is 1 → if o/p was actually 1 but hypothesis gave 0 OR
 o/p was 0 but hypothesis gave 1

Otherwise (when hypothesis and o/p match), its 0.

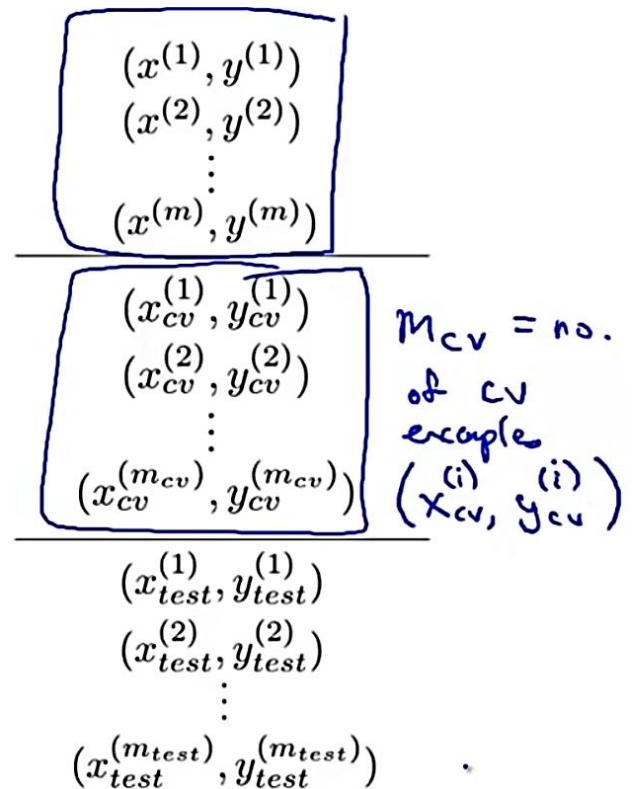
Test error – Gives us the proportion of the test data that was misclassified

HOW TO SELECT MODEL:

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	60% } Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	20% } Cross validation set (cv)
1427	199	
1380	212	20% } test set
1494	243	



Train/validation/test error

Training error:

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad J(\theta)$$

Cross Validation error:

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

We might generally expect $J_{cv}(\theta)$ To be lower than $J_{test}(\theta)$ because:

An extra parameter (d, the degree of the polynomial) has been fit to the cross-validation set.

Model selection

1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(4)})$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(4)})$

➤ Here, superscript is the no of model chosen (degree of polynomial)

➔ Pick the model with lowest $J_{cv}(\theta)$: let's say $d=4$ is the one

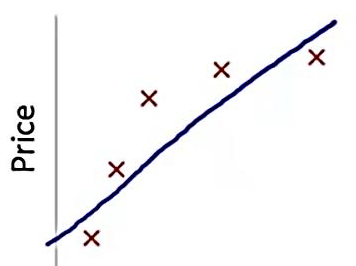
Pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4 \leftarrow$

Estimate generalization error for test set $J_{test}(\theta^{(4)})$

BIAS AND VARIANCE:

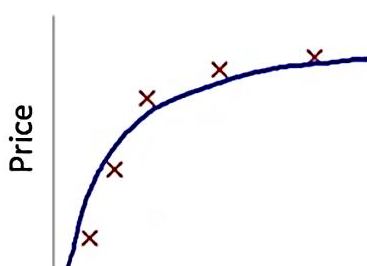
High bias is underfitting and high variance is overfitting. Ideally, we need to find a golden mean between these two.

We need to distinguish whether bias or variance is the problem contributing to bad predictions.



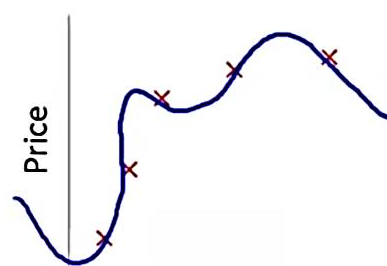
$$\text{Size} \\ \theta_0 + \theta_1 x$$

High bias
(underfit)
 $d=1$



$$\text{Size} \\ \theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"
 $d=2$

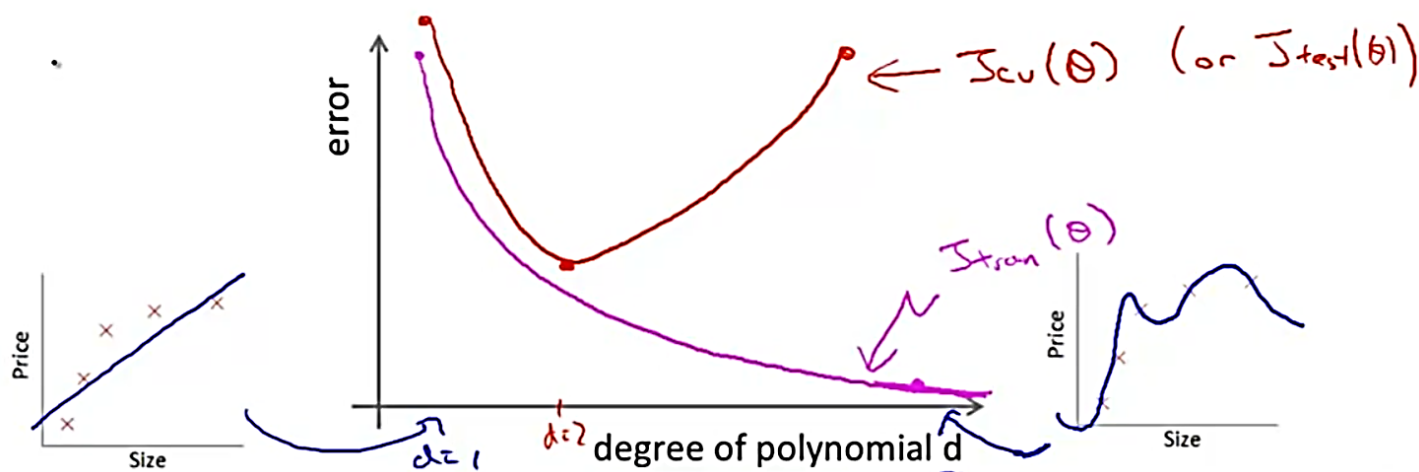


$$\text{Size} \\ \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance
(overfit)
 $d=4$

Training error: $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (or $J_{test}(\theta)$)



The training error will tend to decrease as we increase the degree d of the polynomial.

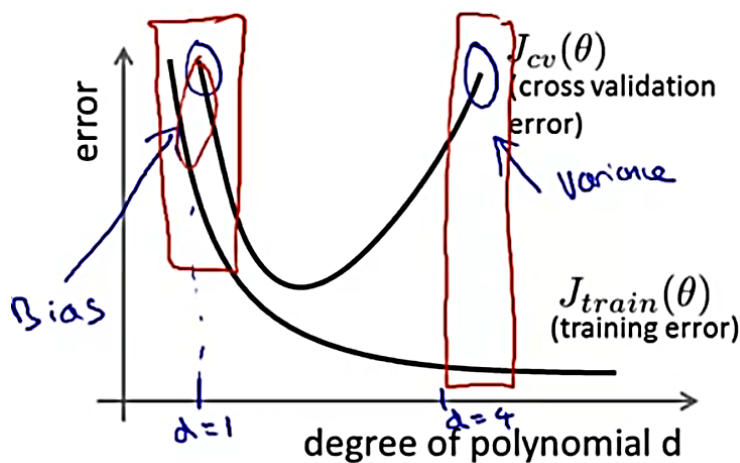
At the same time, the cross-validation error will tend to decrease as we increase d up to a point, and then it will increase as d is increased, forming a convex curve.

DETERMINING WHETHER ITS BIAS OR VARIANCE:

Diagnosing: means determining/localizing the problem

Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



Bias (underfit):
 $\rightarrow J_{train}(\theta)$ will be high
 $J_{cv}(\theta) \approx J_{train}(\theta)$

Variance (overfit):
 $\rightarrow J_{train}(\theta)$ will be low
 $J_{cv}(\theta) \gg J_{train}(\theta)$

High bias (underfitting): both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ will be high. Also, $J_{CV}(\Theta) \approx J_{train}(\Theta)$.

High variance (overfitting): $J_{train}(\Theta)$ will be low and $J_{CV}(\Theta)$ will be much greater than $J_{train}(\Theta)$.

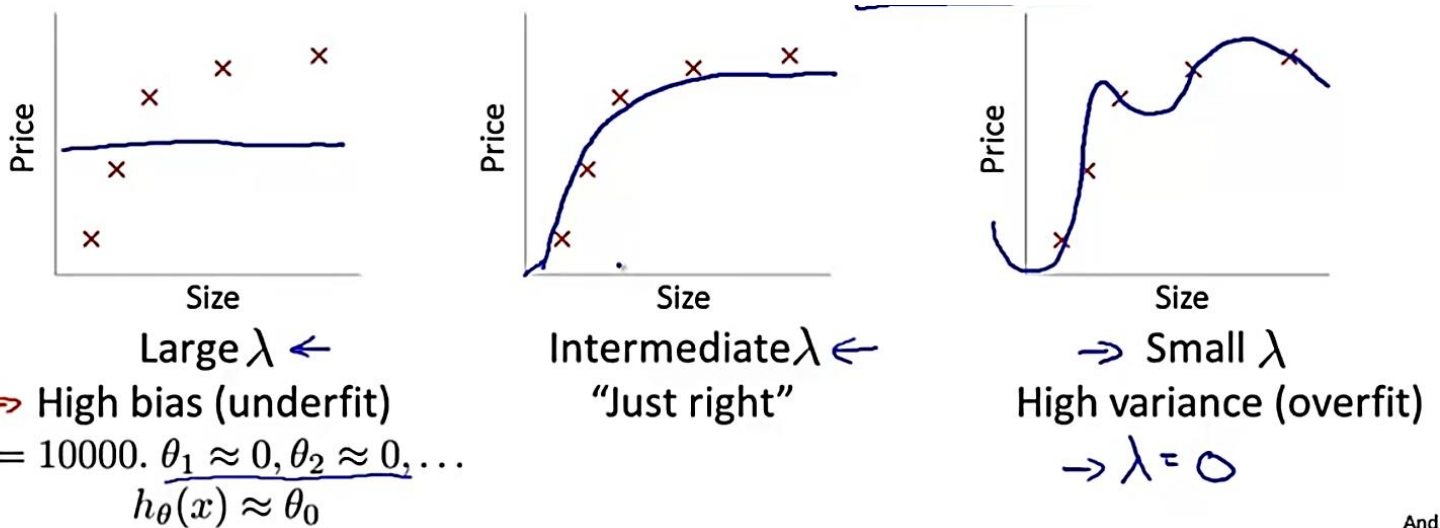
Bias and Variance vs Regularization:

Linear regression with regularization

Model: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

As λ approaches 0, hypothesis tends to overfit



And

Choosing the regularization parameter λ :

1. Try $\lambda = 0 \leftarrow \uparrow \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. Try $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. Try $\lambda = 0.02 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
4. Try $\lambda = 0.04$
5. Try $\lambda = 0.08 \rightarrow \theta^{(5)} \rightarrow J_{cv}(\theta^{(5)})$
- \vdots
12. Try $\lambda = 10 \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
 $\uparrow \quad 10.24$

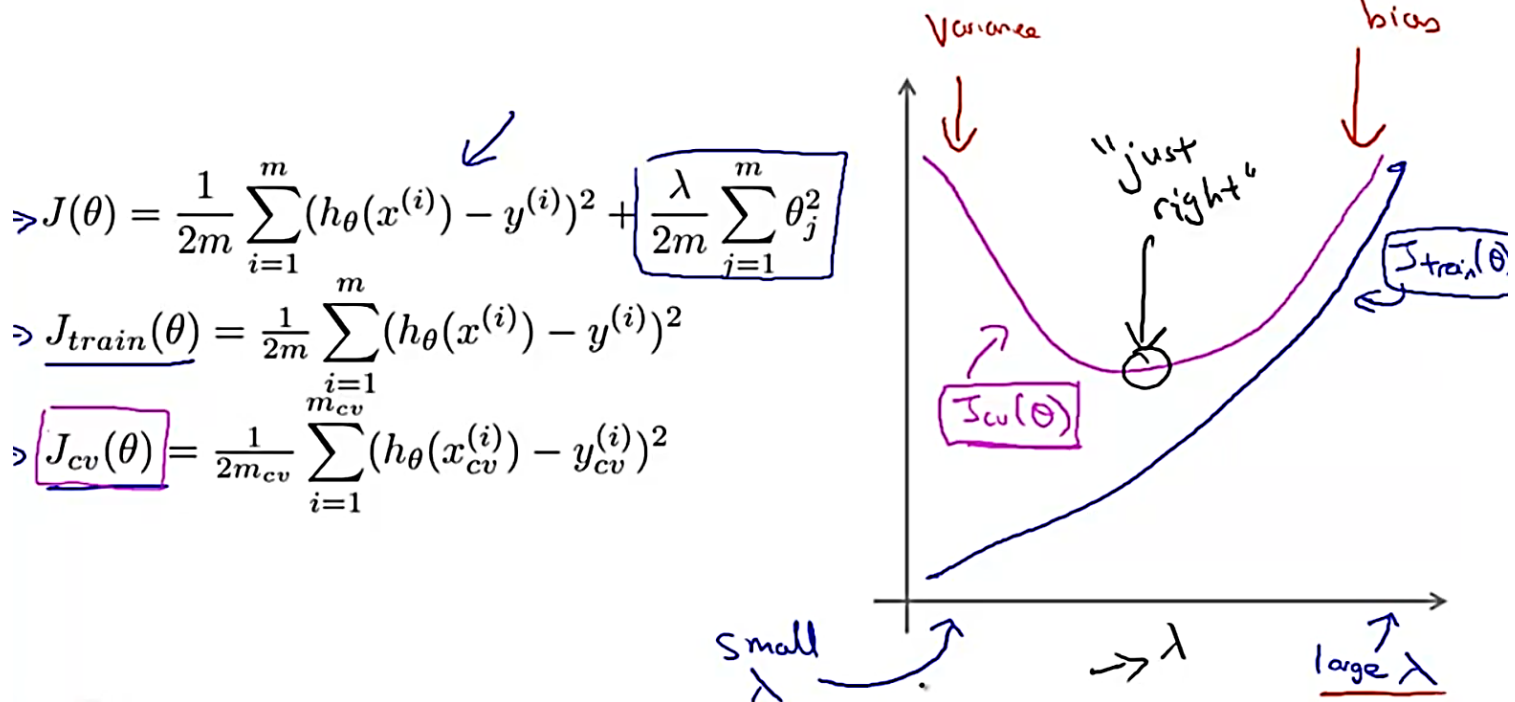
Pick the one which gives lowest J_{cv} (say 5th one)

Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$

1. Create a list of lambdas (i.e.
 $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$);
2. Create a set of models with different degrees or any other variants.
3. Iterate through the λ s and for each λ go through all the models to learn some θ .

4. Compute the cross validation error using the learned Θ (computed with λ) on the $J_{CV}(\Theta)$ **without** regularization or $\lambda = 0$.
5. Select the best combo that produces the lowest error on the cross validation set.
6. Using the best combo Θ and λ , apply it on $J_{test}(\Theta)$ to see if it has a good generalization of the problem.

Bias/variance as a function of the regularization parameter λ

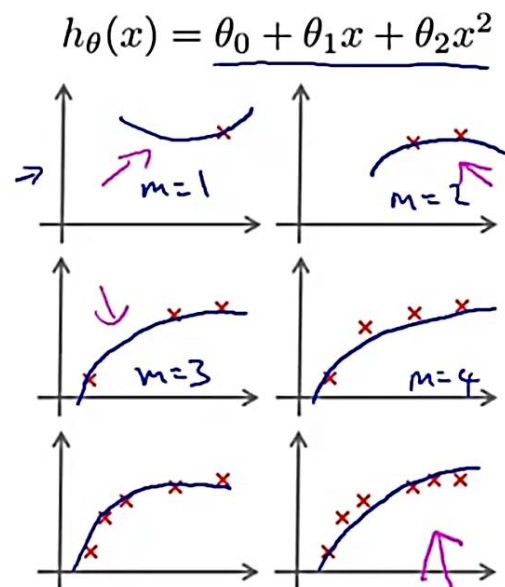
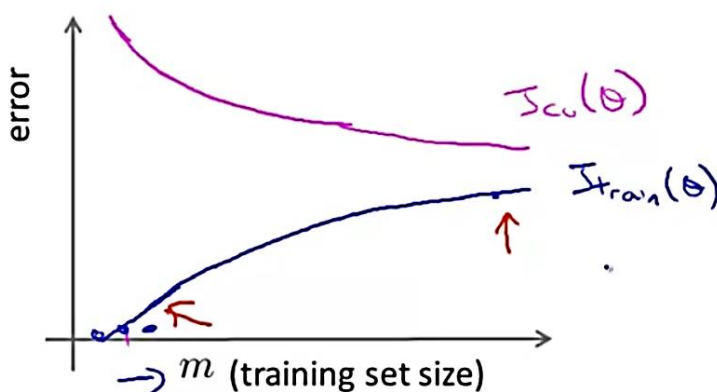


BIAS/VARIANCE vs m (NO OF TRAINING EXAMPLES):

Learning curves

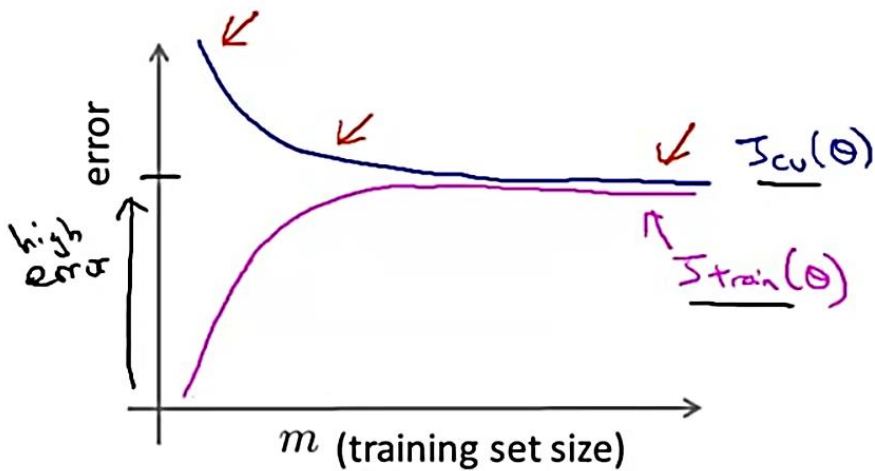
$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

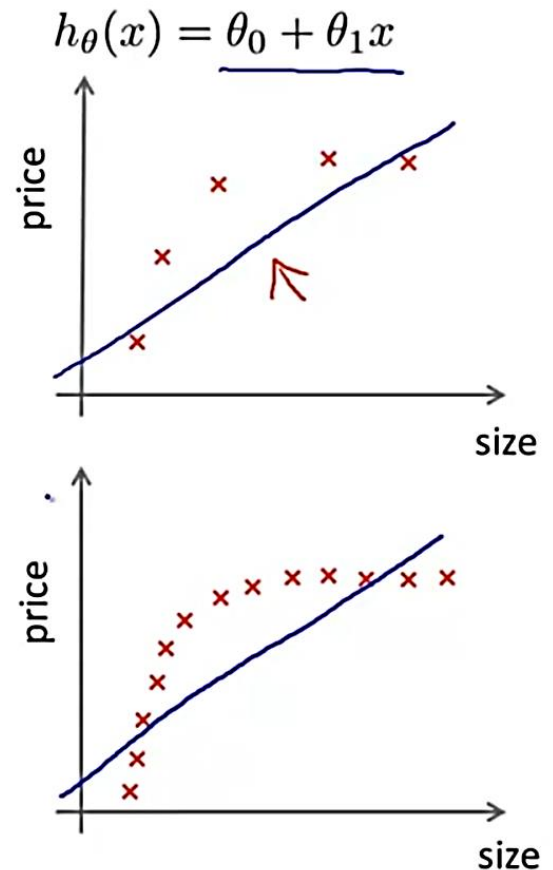


J_{cv} tend to decrease as \rightarrow more the no of training examples, better is the generalization to new data

High bias



If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

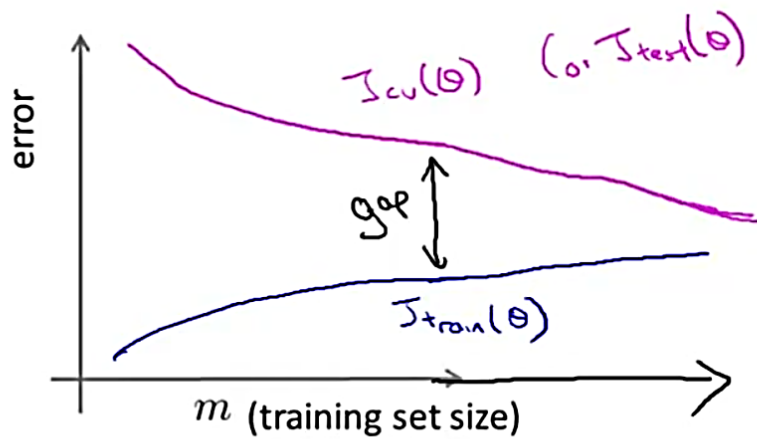


For high bias, as the training examples increases $\rightarrow J_{cv}$ and J_{train} come close to each other and both have a pretty high value

Typical **learning curve** for high bias (at fixed model complexity):



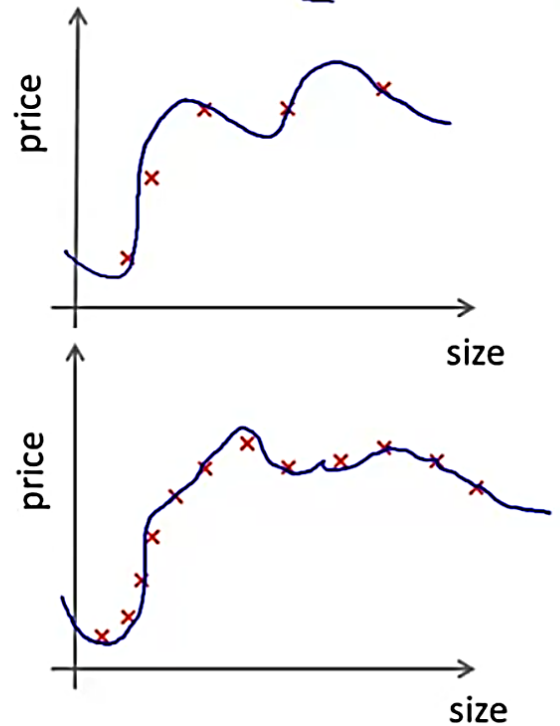
High variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help. ←

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad (11:50 / 11:53)$$

(and small λ)



The J_{train} remains small as the data is perfectly fitted to trained model (overfitted). But, the J_{cv} is quite high due to lack of generalization to new data

For high variance, the gap is quite large, but as we further increase the no of training examples, it is actually helpful.

Typical **learning curve** for high variance (at fixed model complexity):



Debugging a learning algo:

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

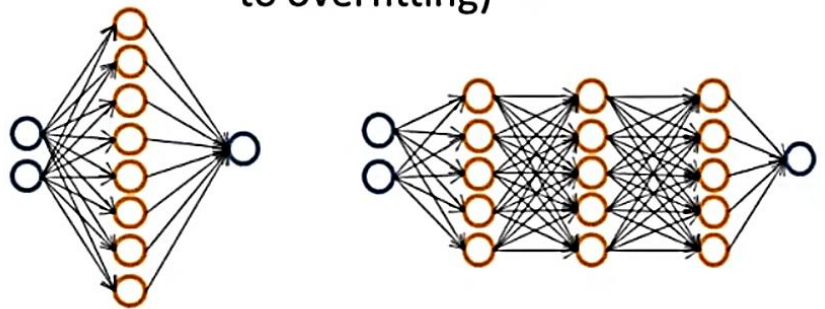
Neural networks and overfitting

“Small” neural network
(fewer parameters; more
prone to underfitting)



Computationally cheaper

“Large” neural network
(more parameters; more prone
to overfitting)



Computationally more expensive.

Use regularization (λ) to address overfitting.

We can choose any type of NN, but a good choice is to choose a large NN with many hidden units or with many hidden layers, and use regularization in it to prevent overfitting.

One good approach is to:

- Try various types of NN
- Select the one with lowest J_{cv}

Model Complexity Effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

→ Why is the recommended approach to perform error analysis using the cross-validation data used to compute $J_{cv}(\theta)$ rather than the test data used to compute $J_{test}(\theta)$?

If we develop new features by examining the test set, then we may end up choosing features that work well specifically for the test set, so $J_{test}(\theta)$ is no longer a good estimate of how well we generalize to new examples.
