Machine Learning

Advice for applying machine learning

Deciding what to try next

**Debugging a learning algorithm:**

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

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

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?

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

**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.

Advice for applying
machine learning

Evaluating a
hypothesis

Machine Learning

# Evaluating your hypothesis

price

size

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

Fails to generalize to new examples not in training set.

$x_1 =$ size of house
$x_2 =$ no. of bedrooms
$x_3 =$ no. of floors
$x_4 =$ age of house
$x_5 =$ average income in neighborhood
$x_6 =$ kitchen size
$\vdots$
$x_{100}$

# Evaluating your hypothesis

Dataset:

| Size | Price |
|------|-------|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| 1985 | 300 |
| 1534 | 315 |
| 1427 | 199 |
| 1380 | 212 |
| 1494 | 243 |

70%

30%

Training set

Test Set

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

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

$m_{test}$ = no. of test example

$(x_{test}^{(i)}, y_{test}^{(i)})$

# Training/testing procedure for linear regression

→ - Learn parameter $\theta$ from training data (minimizing training error $J(\theta)$)

      70%

- Compute test set error:

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

# Training/testing procedure for logistic regression

- Learn parameter $\theta$ from training data
- Compute test set error:

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

- Misclassification error (0/1 misclassification error):

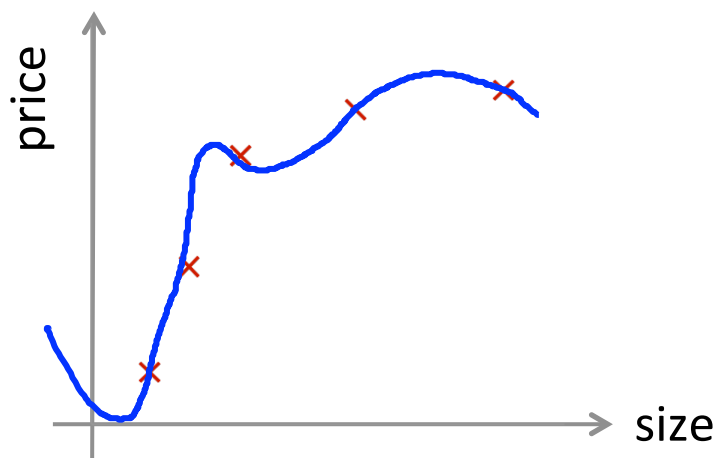Machine Learning

# Advice for applying machine learning

## Model selection and training/validation/test sets

# Overfitting example



price

size

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

Once parameters $\theta_0, \theta_1, \ldots, \theta_4$ were fit to some set of data (training set), the error of the parameters as measured on that data (the training error $J(\theta)$) is likely to be lower than the actual generalization error.

**Model selection**

$d=1$  1. $\to h_\theta(x) = \theta_0 + \theta_1 x$ $\longrightarrow$ $\Theta^{(1)}$ $\longrightarrow$ $J_{test}(\Theta^{(1)})$

$d=2$  2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$ $\longrightarrow$ $\Theta^{(2)}$ $\longrightarrow$ $J_{test}(\Theta^{(2)})$

$d=3$  3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$ $\longrightarrow$ $\Theta^{(3)}$ $\to$ $J_{test}(\Theta^{(3)})$

$\vdots$

$d=10$  10. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$ $\to$ $\Theta^{(10)}$ $\to$ $J_{test}(\Theta^{(10)})$

$d$ = degree of polynomial

Choose $\theta_0 + \ldots \theta_5 x^5$ $\Leftarrow$

You should not use the test set to choose the regularization parameter, as you will then have an artificially low value for test error and it will not give a good estimate of generalization error. The cross validation lets us find the "just right" setting of the regularization parameter given the fixed model parameters learned from the training set. We can then use this to find the test error without risking an optimistic estimate of generalization error.

$\Theta_0, \Theta_1 \ldots$

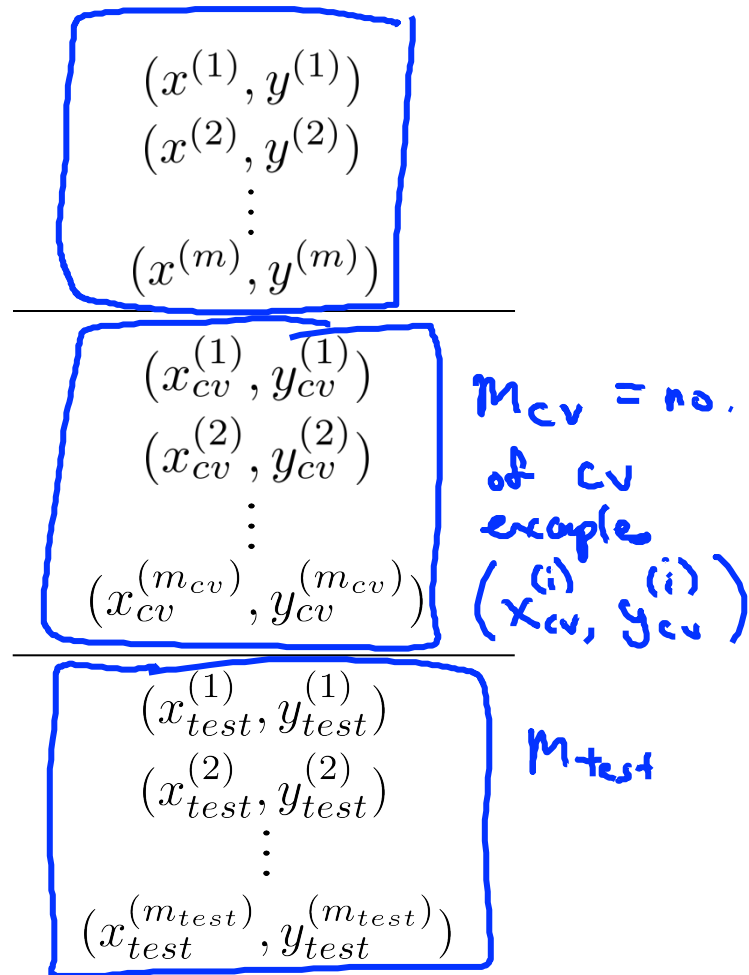How well does the model generalize? Report test set error $J_{test}(\theta^{(5)})$.

$\Theta^{(5)}$

Problem: $J_{test}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error. I.e. our extra parameter ($d$ = degree of polynomial) is fit to test set.

Andrew Ng

# Evaluating your hypothesis

## Dataset:

| Size | Price |
|------|-------|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| 1985 | 300 |
| 1534 | 315 |
| 1427 | 199 |
| 1380 | 212 |
| 1494 | 243 |

$60\%$ } Traing set

$20\%$ } Cross validation set (CV)

$20\%$ } test set

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

$$(x_{cv}^{(1)}, y_{cv}^{(1)})$$
$$(x_{cv}^{(2)}, y_{cv}^{(2)})$$
$$\vdots$$
$$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$$

$M_{cv} = $ no. of CV example $(x_{cv}^{(i)}, y_{cv}^{(i)})$

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

$M_{test}$

Andrew Ng

# Train/validation/test error

Training error:

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

Cross Validation error:

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

Test error:

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

# Model selection

$d=1$   1.   $h_\theta(x) = \theta_0 + \theta_1 x$   $\longrightarrow$   $\min_\theta J(\theta) \rightarrow \Theta^{(1)} \longrightarrow J_{cv}(\Theta^{(1)})$

$d=2$   2.   $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$   $\longrightarrow$   $\Theta^{(2)} \longrightarrow J_{cv}(\Theta^{(2)})$

$d=3$   3.   $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$   $\longrightarrow$   $\Theta^{(3)}$

          $\vdots$                                           $J_{cv}(\Theta^{(4)})$

$d=10$   10.   $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$   $\longrightarrow$   $\Theta^{(10)} \longrightarrow J_{cv}(\Theta^{(10)})$

$$\underline{d = 4} \longrightarrow \nearrow$$

Pick $\theta_0 + \theta_1 x_1 + \cdots + \theta_4 x^4$ $\longleftarrow$

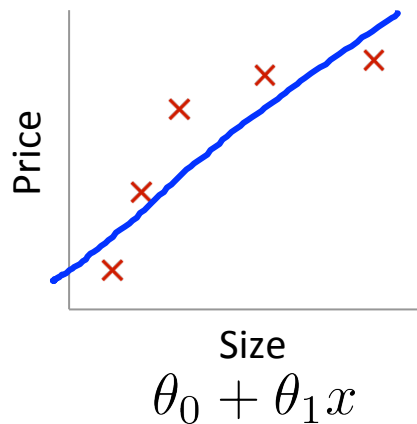Estimate generalization error   for test set $\underline{J_{test}(\theta^{(4)})}$ $\longleftarrow$

# Bias/variance



$$\theta_0 + \theta_1 x$$

High bias
(underfit)

$d=1$

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"

$d=2$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance
(overfit)
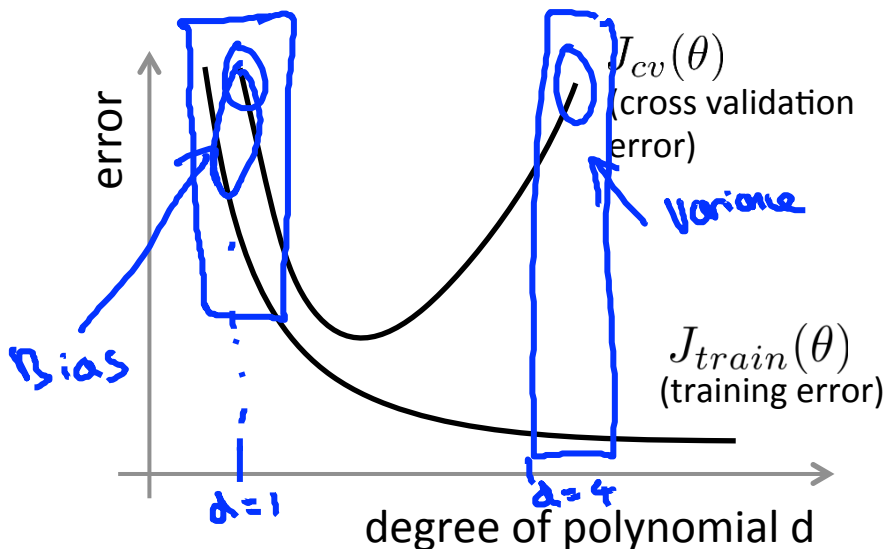
$d=4$

Andrew Ng

# Bias/variance

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

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



$\leftarrow J_{cv}(\theta)$ $\left( \text{or } J_{test}(\theta) \right)$

$J_{train}(\theta)$

error

degree of polynomial d

$d=1$ $d=2$

Price / Size

Price / Size

# 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?



error

$J_{cv}(\theta)$
(cross validation error)

Variance

Bias

$d=1$

$d=4$

$J_{train}(\theta)$
(training error)

degree of polynomial d

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

$J_{cv}(\theta) >> J_{train}(\theta)$

$>>$

Machine Learning

Advice for applying machine learning

Regularization and bias/variance

# Linear regression with regularization

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

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



Large $\lambda$
High bias (underfit)
$\lambda = 10000. \ \theta_1 \approx 0, \theta_2 \approx 0, \dots$
$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$
"Just right"

Small $\lambda$
High variance (overfit)
$\lambda = 0$

Andrew Ng

# Choosing the regularization parameter $\lambda$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

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

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

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

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

$J(\theta)$

$J_{train}$
$J_{cv}$
$J_{test}$

**Choosing the regularization parameter $\lambda$**

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$
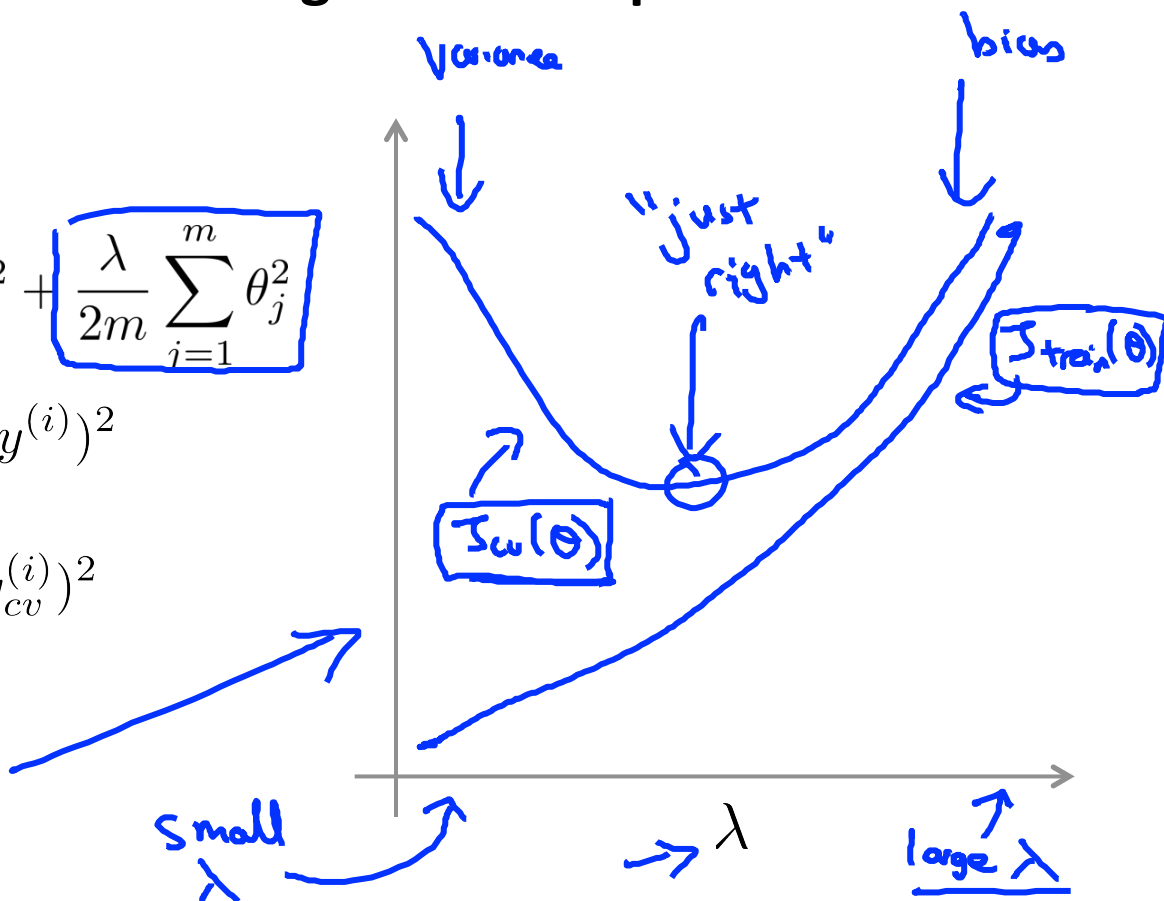
1. Try $\lambda = 0$    $\min_\theta J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$

2. Try $\lambda = 0.01$    $\min_\theta J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$

3. Try $\lambda = 0.02$       $\theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$

4. Try $\lambda = 0.04$

5. Try $\lambda = 0.08$       $\vdots$   $\theta^{(5)}$    $J_{cv}(\theta^{(5)})$

    $\vdots$

12. Try $\lambda = 10$       $\theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$

$\uparrow$ $\overline{10.24}$

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

# Bias/variance as a function of the regularization parameter $\lambda$

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

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

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



Variance

bias

"just right"

$\boxed{J_{train}(\theta)}$

$\boxed{J_{cv}(\theta)}$

small $\lambda$

$\lambda$

large $\lambda$

# Advice for applying machine learning

## Learning curves

Machine Learning

# Learning curves

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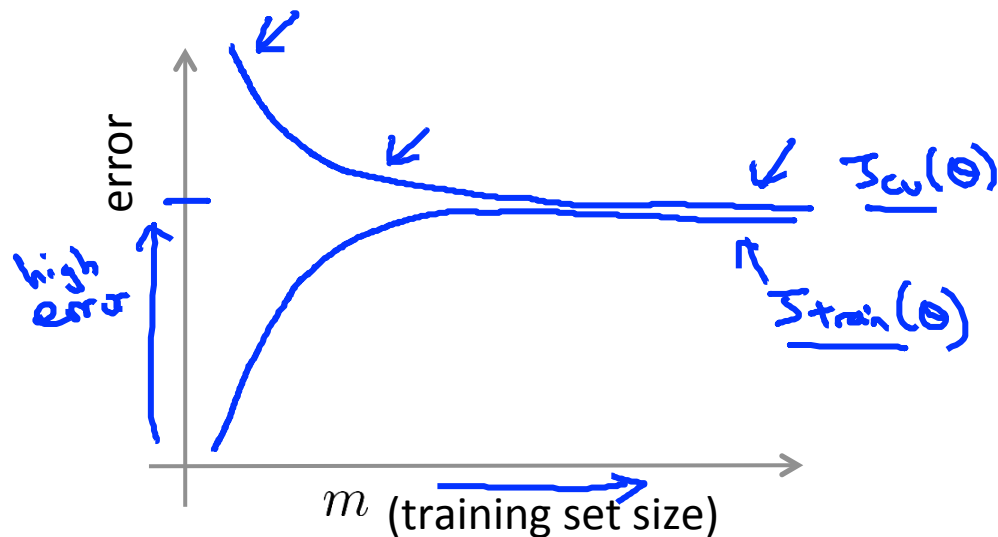
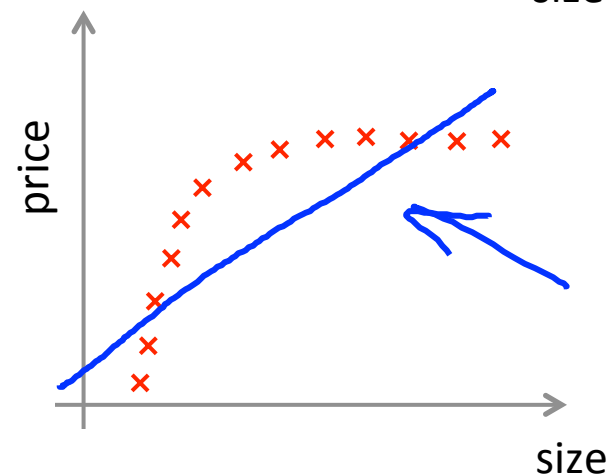$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

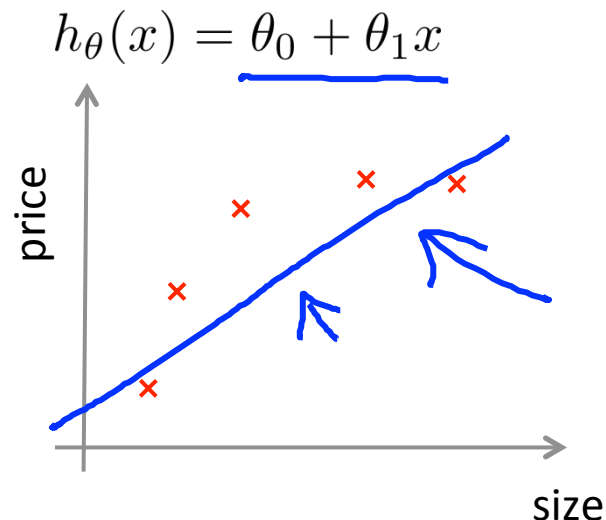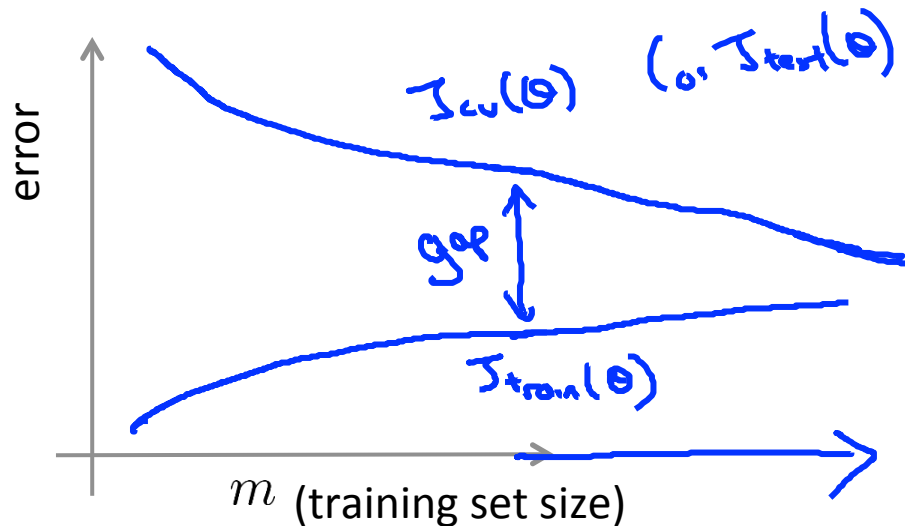$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$m=1$

$m=2$

$m=3$

$m=4$



error

$J_{cv}(\theta)$

$J_{train}(\theta)$

$m$ (training set size)

**High bias**



$$h_\theta(x) = \theta_0 + \theta_1 x$$

error

high error

$J_{cv}(\theta)$

$J_{train}(\theta)$

$m$ (training set size)

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

price

size

price

size

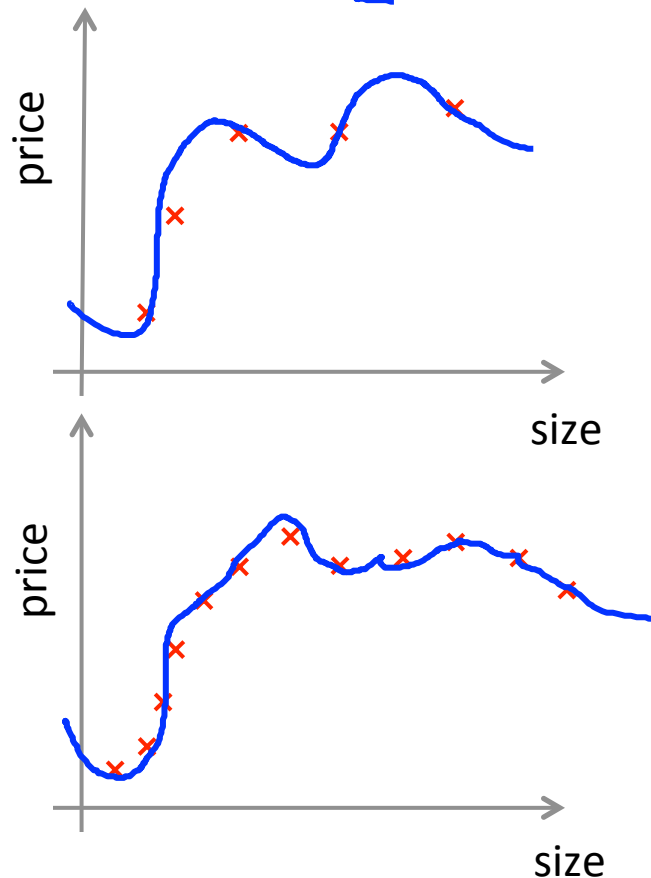**High variance**



$$h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{100} x^{100}$$
(and small $\lambda$)

error

$J_{cv}(\theta)$ $(or J_{test}(\theta))$

gap

$J_{train}(\theta)$

$m$ (training set size)

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

price

size

price

size

Machine Learning

Advice for applying machine learning
_____

Deciding what to try next (revisited)
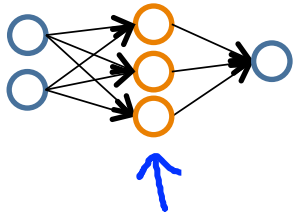
**Debugging a learning algorithm:**

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 $(x_1^2, x_2^2, x_1 x_2, \text{etc})$ → *fixes high bias.*
- Try decreasing $\lambda$ → *fixes high bias*
- Try increasing $\lambda$ → *fixes high variance*

# Neural networks and overfitting

Adding more layers will increase model complexity, making the variance problem worse.

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

"Large" neural network (more parameters; more prone to overfitting)



Computationally cheaper

Computationally more expensive.

Use regularization ($\lambda$) to address overfitting.

$J_{c_0}(\vec{s})$

Andrew Ng