# Advice for applying machine learning

## Deciding what to try next

But sometimes getting more training data doesn't actually help and in the next few videos we will see why, and we will see how you can avoid spending a lot of time collecting more training data in settings where it is just not going to help.

Try getting additional features也可能很花時間.

Machine Learning

**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?

後面會講 甚麼情況下 用以下的甚麼方法

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

Unfortunately, the most common method that people use to pick one of these is to go by gut feeling. In which what many people will do is sort of randomly pick one of these options

**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

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

Andrew Ng

# Training/testing procedure for linear regression

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

70%

- Compute test set error:

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

# Training/testing procedure for logistic regression

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

$$J_{test}(\theta) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h_\theta(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log h_\theta(x_{test}^{(i)})$$
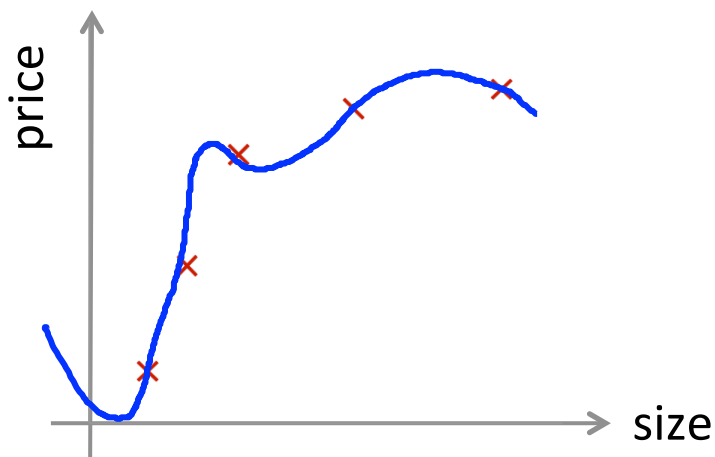
- Misclassification error (0/1 misclassification error):

Machine Learning

# Advice for applying machine learning

## Model selection and training/validation/test sets

**Overfitting example**

We've already seen a lot of times the problem of overfitting, in which just because a learning algorithm fits a training set well, that doesn't mean it's a good hypothesis. More generally, this is why the training set's error is not a good predictor for how well the hypothesis will do on new example.



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

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.

Andrew Ng

in order to select one of these models, I could
then see which model has the lowest test set error.

$d$ = degree of polynomial

## Model selection

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

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

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

$\vdots$

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

Choose $\boxed{\theta_0 + \ldots \theta_5 x^5} \leftarrow$

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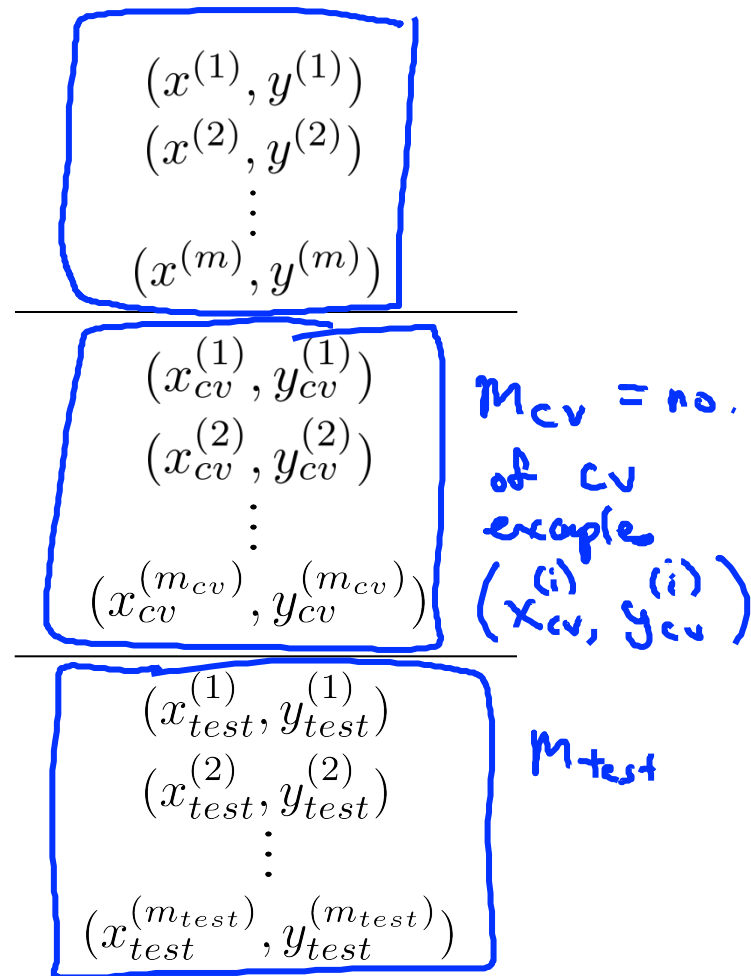
$\Theta^{(5)}$

$\boxed{\Theta_0, \Theta_1 \ldots}$

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

Andrew Ng

To address this problem, in a model selection setting, if we want to evaluate a hypothesis, this is what we usually do instead.

# 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%  Training 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 examples
$(x_{cv}^{(i)}, y_{cv}^{(i)})$

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

$m_{test}$

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

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

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

$J_{cv}(\theta^{(4)})$

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

注意, 這裡的cross validation沒像Gatech ML note中那樣輪換弄.

$d = 4$

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

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

Instead of using the test set to select the model, we're instead going to use the validation set, or the cross validation set, to select the model. Concretely, we're going to first take our first hypothesis, take this first model, and say, minimize the cross function, and this would give me some parameter vector theta for the new model. Instead of testing these hypotheses on the test set, I'm instead going to test them on the cross validation set. And then I'm going to pick the hypothesis with the lowest cross validation error. And so this degree of polynomial

Andrew Ng

Advice for applying machine learning

Diagnosing bias vs. variance

Machine Learning

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



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

$J_{train}(\theta)$

degree of polynomial d

d=1   d=2

# 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

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

$\gg$

Andrew Ng

Machine Learning

Advice for applying machine learning
_____

Regularization and bias/variance

# Linear regression with regularization

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

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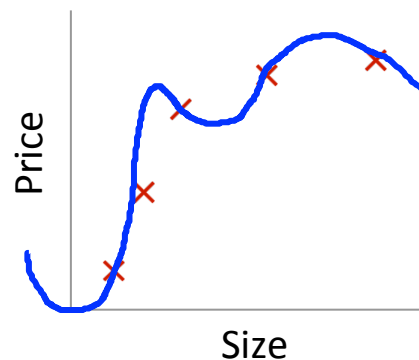


Large $\lambda$

High bias (underfit)

$\lambda = 10000.\ \theta_1 \approx 0, \theta_2 \approx 0, \ldots$

$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$

"Just right"

Small $\lambda$

High variance (overfit)

$\lambda = 0$

Andrew Ng

# Choosing the regularization parameter $\lambda$

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

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

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

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

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

$J(\theta)$

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

Andrew Ng

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

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

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

<span style="color:orange">通過使J(theta)最小, 來求得theta^(1), 同理求得theta^(2)等, 然後用它們去算J_CV, 選出使J_CV最小的theta</span>
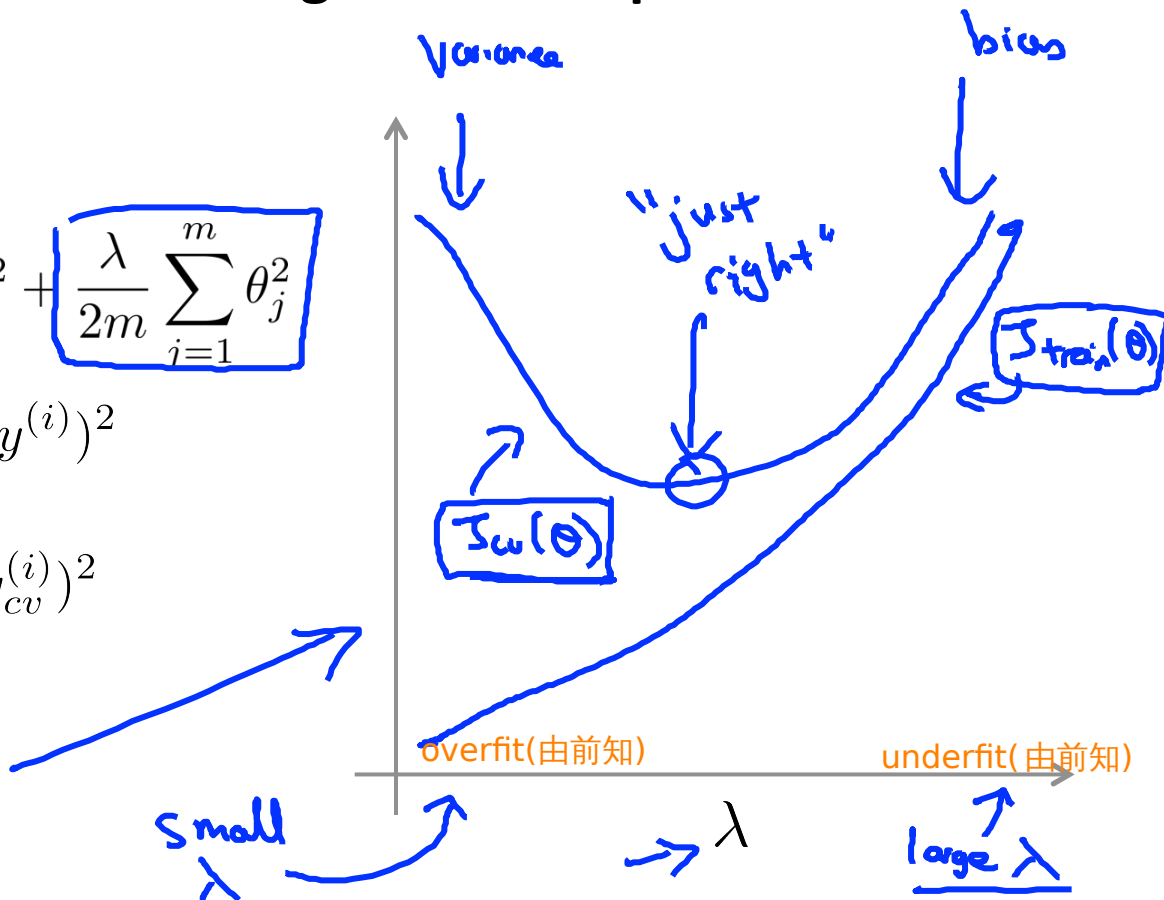
1. Try $\lambda = 0$ $\leftarrow$ $\longrightarrow$ $\min\limits_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$

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

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

4. Try $\lambda = 0.04$

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

   $\vdots$

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

   $\uparrow$ $\overline{10.24}$

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

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

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

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

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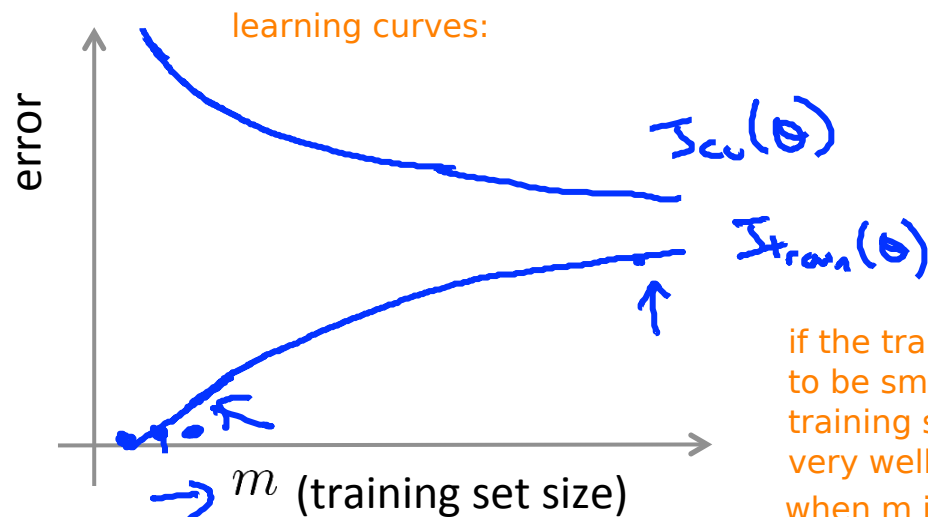
Variance

bias

"just right"

$J_{train}(\theta)$

$J_{cv}(\theta)$

overfit(由前知)          underfit( 由前知)

Small $\lambda$

$\lambda$

large $\lambda$

Advice for applying
machine learning

Learning curves

Machine Learning

# Learning curves

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

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

learning curves:



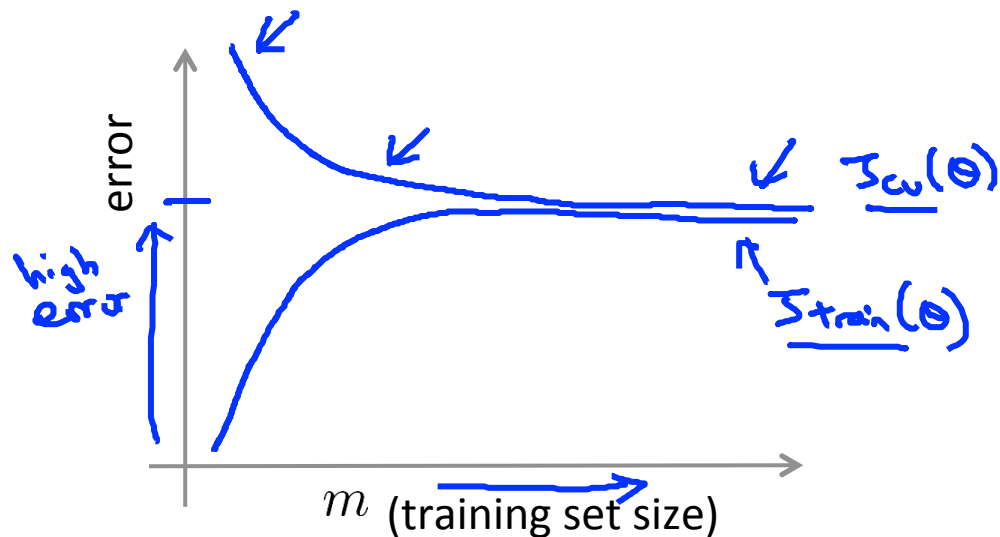$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

no error  fit quite well

m=1   m=2

m=3   m=4

if the training set size is small then the training error is going to be small as well. Because you know, we have a small training set is going to be very easy to fit your training set very well may be even perfectly

when m is larger then gets harder all the training examples perfectly and so your training set error becomes more larger.
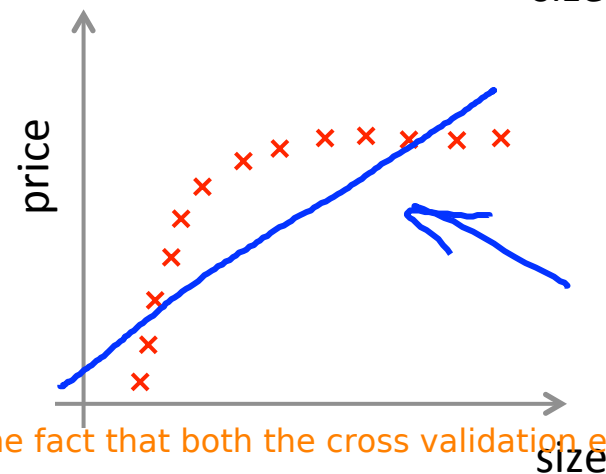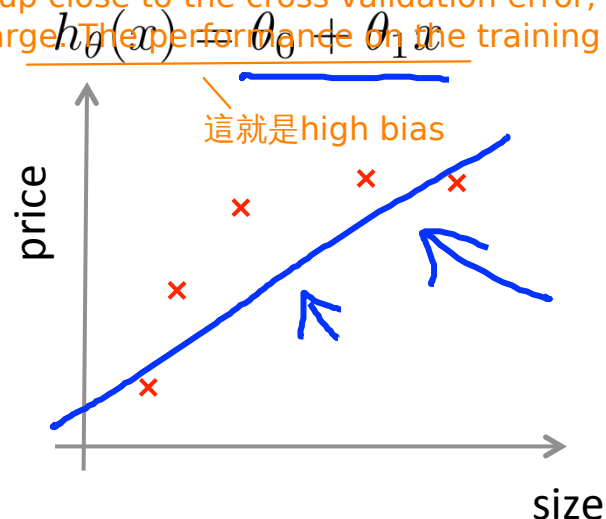
Andrew Ng

## High bias

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

這就是high bias



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

the problem with high bias is reflected in the fact that both the cross validation error and the training error are high
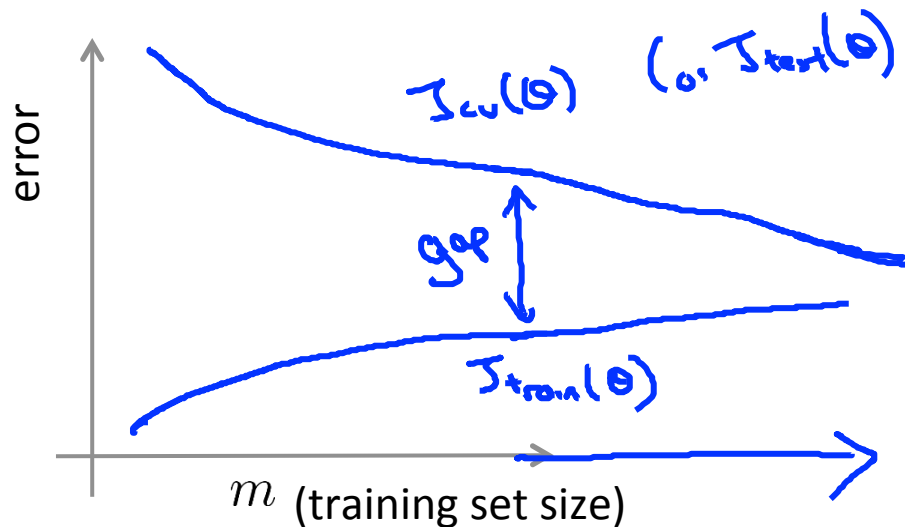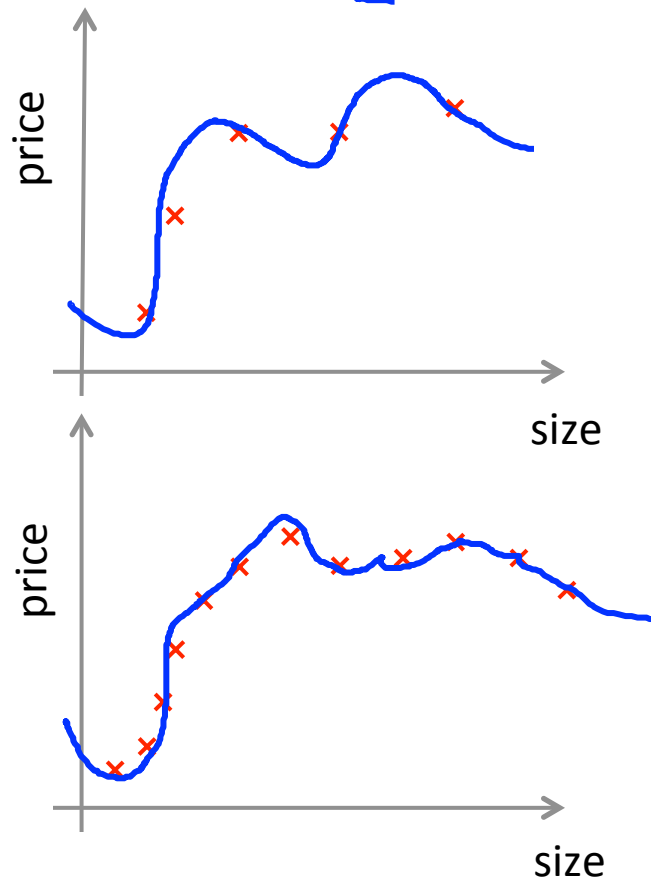
Andrew Ng

# High variance

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

(and small $\lambda$)



$J_{cv}(\theta)$   (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

Andrew Ng

Machine Learning

Advice for applying machine learning
_____

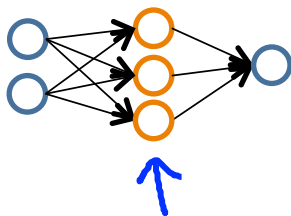Deciding what to try next (revisited)

**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 $\rightarrow$ *fixes high variance*
- Try smaller sets of features $\rightarrow$ *fixes high variance*
- Try getting additional features $\rightarrow$ *fixes high bias*
- Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, \text{etc})$ $\rightarrow$ *fixes high bias.*
- Try decreasing $\lambda$ $\rightarrow$ *fixes high bias*
- Try increasing $\lambda$ $\rightarrow$ *fixes high variance*

Andrew Ng

# Neural networks and overfitting

"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.

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