1.遇到错误怎么解决：

**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        $X_1, X_2, X_3, \ldots, X_{100}$
$\rightarrow$ - 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.

2. 评估假设：

按照7：3分为训练集和测试集。

# 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%  Training set

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

30%  Test Set

第一组测试样本

$$(x^{(1)}_{test}, y^{(1)}_{test})$$
$$(x^{(2)}_{test}, y^{(2)}_{test})$$
$$\vdots$$
$$(x^{(m_{test})}_{test}, y^{(m_{test})}_{test})$$

$m_{test} = $ no. of test example $(x^{(i)}_{test}, y^{(i)}_{test})$

最好随机选择70%作为训练集。

典型的方法：

# Training/testing procedure for linear regression

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

- Compute test set error:

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

# Training/testing procedure for logistic regression

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

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

# Training/testing procedure for logistic regression

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

$M_{test}$

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

- Misclassification error (0/1 misclassification error):

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

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

那部分测试集中的样本

模型选择：通过分别对不同模型进行训练，得到各个模型测试集误差，现在从中选择误差最小的模型，如果要测试这个模型，就不能在原来的测试集上进行测试，因为这些参数就是在测试集上训练而来的。

## Model selection

$\rightarrow d$ = degree of polynomial

$d=1$ 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. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3 \longrightarrow \Theta^{(3)} \longrightarrow J_{test}(\Theta^{(3)})$

$\vdots$

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

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.

因此采用下面方法：把数据集分为三个部分：训练集，交叉验证集，测试集：

## 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

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

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

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

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

$M_{test}$

用m 下标test来表示测试样本的总数

各种误差：

## Train/validation/test error

Training error:

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

Cross Validation error:

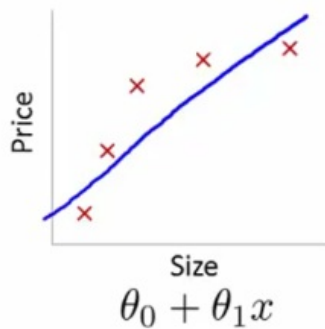$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

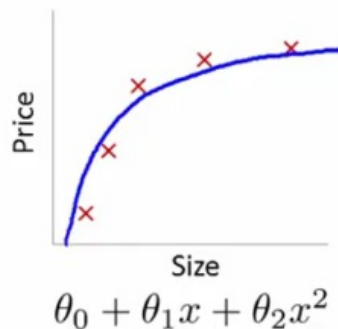$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$
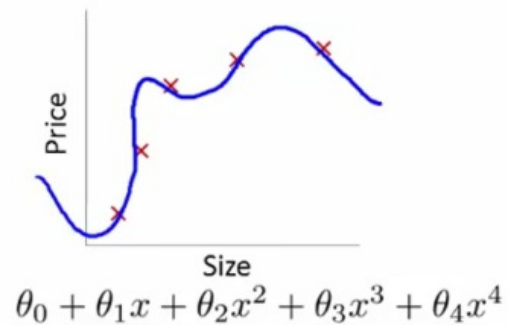
和测试误差

Bias / Variance

# Bias/variance



$\theta_0 + \theta_1 x$

**High bias (underfit)**

$d=1$

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

**"Just right"**

$d=2$

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
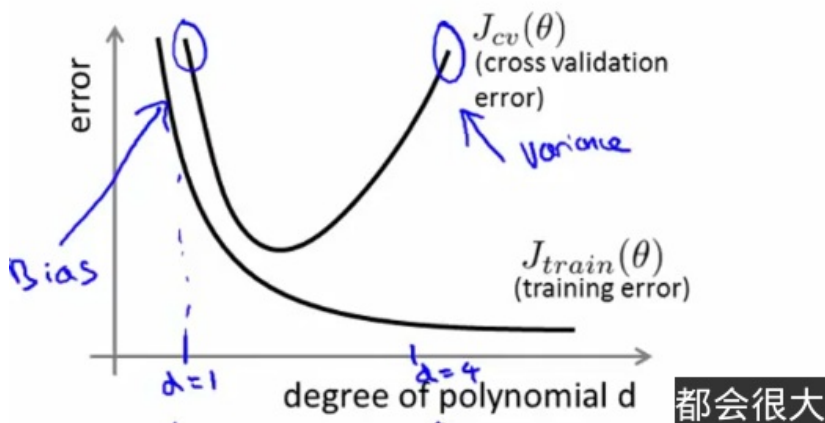
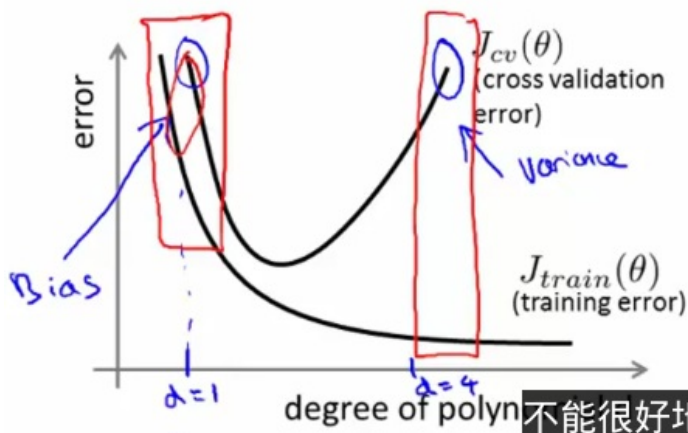**High variance (overfit)**

$d=4$

训练误差和交叉验证误差：

# 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?



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

Variance

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

Bias

$d=1$    degree of polynomial d    $d=4$    都会很大

## 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):
$J_{train}(\theta)$ will be high
$J_{cv}(\theta) \approx J_{train}(\theta)$

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

不能很好地拟合训练集数据

诊断一个模型是处于什么状态：

## Linear regression with regularization

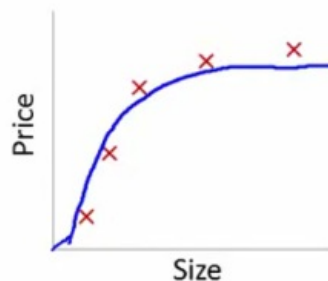Model: $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

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



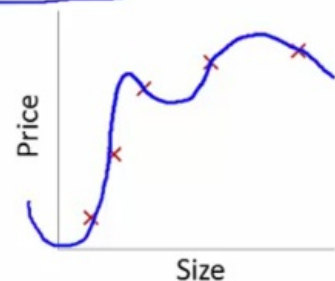Large $\lambda$
High bias (underfit)
$\lambda = 10000.\ \theta_1 \approx 0, \theta_2 \approx 0, \ldots$
$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$
"Just right"
但在这里我就不讨论这些情况了

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

Andre

选择正规划参数lambda：

## Choosing the regularization parameter $\lambda$

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

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

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

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

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

$J(\theta)$

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

这就是模型选择在选取正则化参数 λ 时的应用
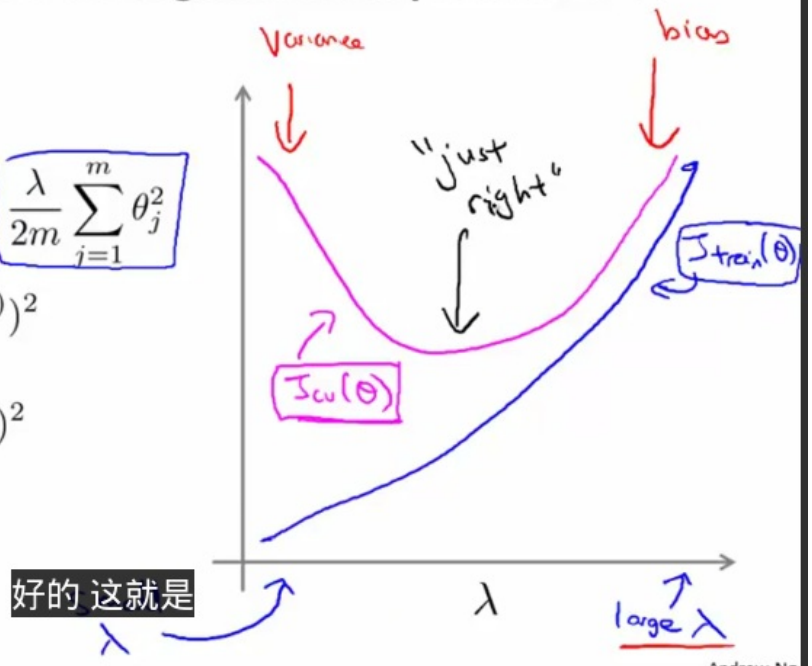
## Choosing the regularization parameter $\lambda$

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 $\lambda = 0 \leftarrow$    $\longrightarrow$   $\min_\theta J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. Try $\lambda = 0.01$   $\longrightarrow$   $\min_\theta J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. Try $\lambda = 0.02$    $\longrightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
4. Try $\lambda = 0.04$
5. Try $\lambda = 0.08$    $\longrightarrow \theta^{(5)}$   $J_{cv}(\theta^{(5)})$

   $\vdots$       $\vdots$

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

      $\uparrow$ $10.24$

来测出它对测试集的

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

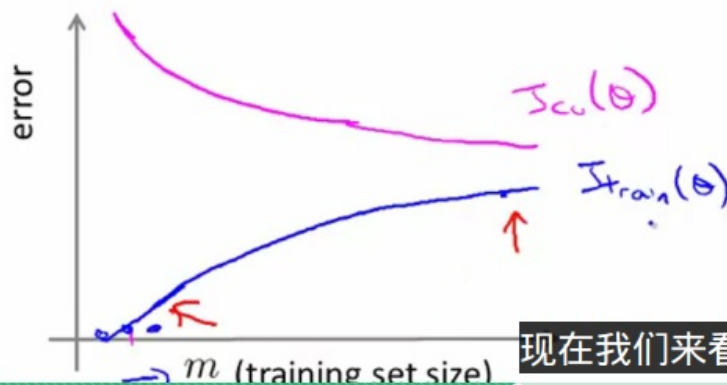## Bias/variance as a function of the regularization parameter $\lambda$

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

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

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

Variance

bias

"just right"

$J_{cv}(\theta)$

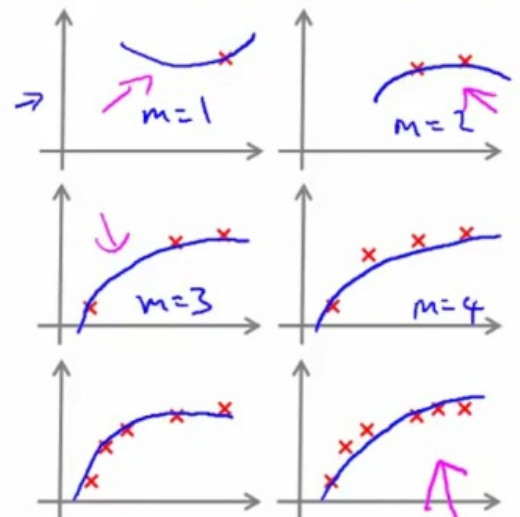$J_{train}(\theta)$

好的 这就是

$\lambda$

large $\lambda$

Andrew Ng

学习曲线：

## Learning curves

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

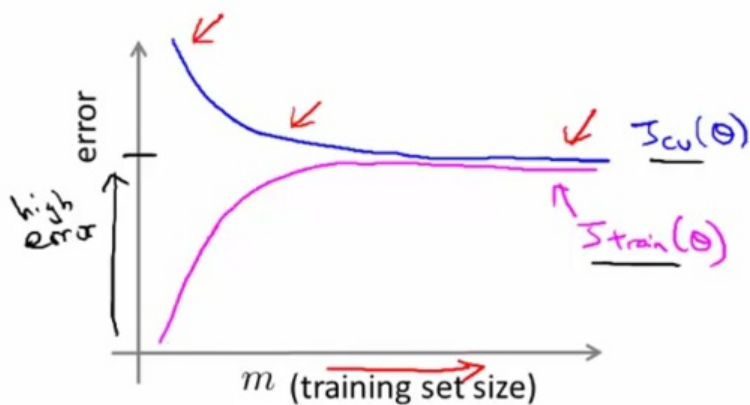$\Rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}}\sum_{i=1}^{m_{cv}}(h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$
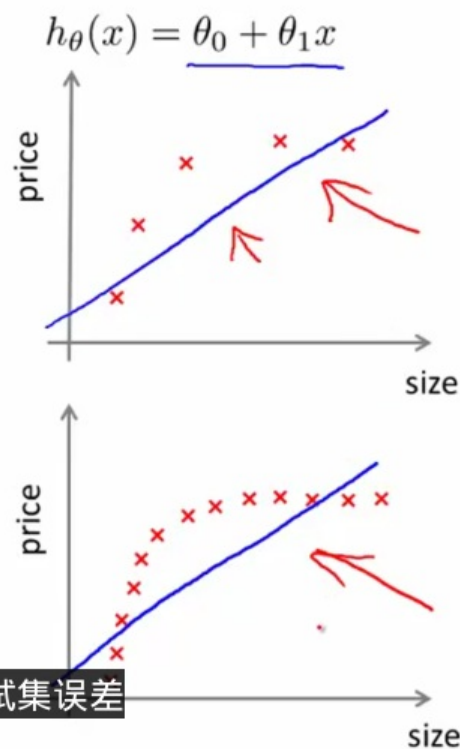
$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

m=1

M=2

m=3

M=4

error

$J_{cv}(\theta)$

$J_{train}(\theta)$

$m$ (training set size)

现在我们来看看

高偏差状态（欠拟合）：

## High bias



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

error — high error

$J_{cv}(\theta)$

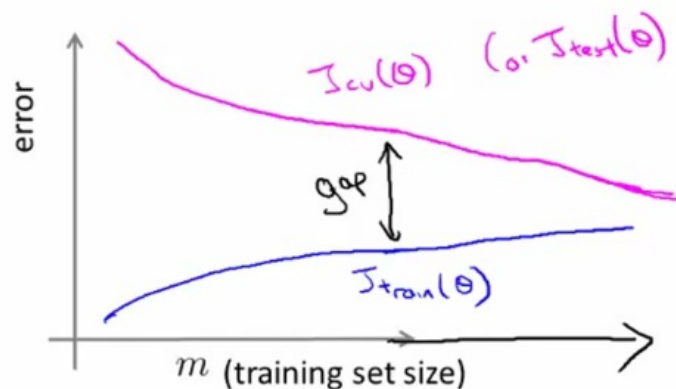$J_{train}(\theta)$

$m$ (training set size)

交叉验证集误差或测试集误差

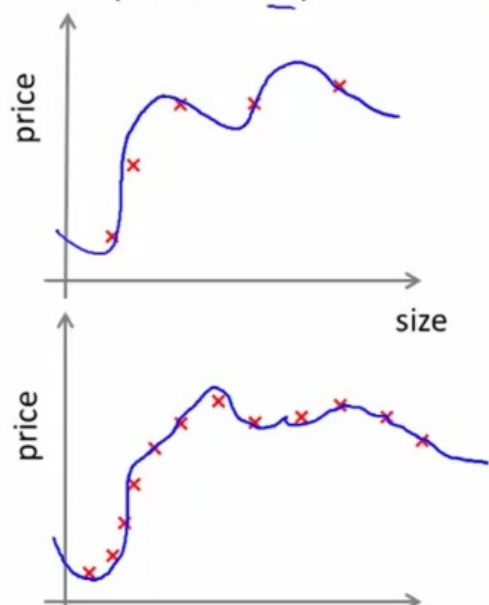If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

price / size

过拟合：

## High variance



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

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.
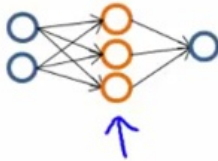
price / size

**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 $(x_1^2, x_2^2, x_1x_2, \text{etc}) \rightarrow$ fixes high bias.
- Try decreasing $\lambda$ $\rightarrow$ fixes high bias
- Try increasing $\lambda$ $\rightarrow$ fixes high variance

## Neural networks and overfitting

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

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

Computationally cheaper

Computationally more expensive.

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

$J_{c_0}(\theta)$