# Model Selection and Train/Validation/Test Sets

*Evaluating a Learning Algorithm*

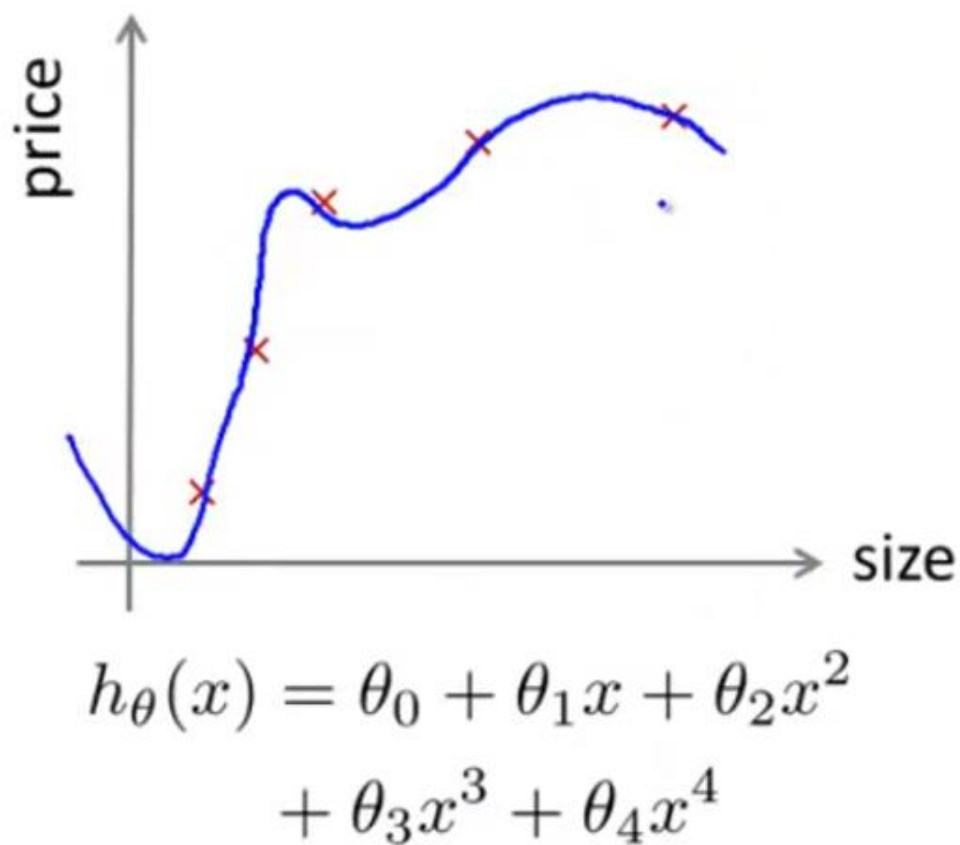Advice for Applying Machine Learning

# Introduction

- Suppose you are left to decide what degree of polynomial to fit to a data set.

- So that what features to include that gives you a learning algorithm.

- Or suppose you'd like to choose the regularization parameter lambda for learning algorithm

- These are called model selection problems.

# Introduction

- 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.

# Overfitting 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

# Model selection

1. $h_\theta(x) = \theta_0 + \theta_1 x$
2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$
   $\vdots$
10. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$

Andrew Ng

# Model selection

$d = $ degree of polynomial

1. $h_\theta(x) = \theta_0 + \theta_1 x$
2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$

    $\vdots$

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

Andrew Ng

# Model selection

$d =$ degree of polynomial

$d=1$   1. $\rightarrow h_\theta(x) = \theta_0 + \theta_1 x$

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

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

    $\vdots$     $\vdots$

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

Andrew Ng

# Model selection

$d = $ degree of polynomial

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

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

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

$\vdots$    $\vdots$

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

Andrew Ng

# Model selection

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

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

$\vdots \qquad \vdots \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \vdots$

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

Andrew Ng

**Model selection**

$d = $ degree of polynomial ↓

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

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

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

  ⋮    ⋮

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

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

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

Andrew Ng

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

$\vdots$

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

$\boxed{d = \text{degree of polynomial}}$

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

How well does the model generalize? Report test set $\boxed{\Theta_0, \Theta_1 \ldots}$
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 ($\underline{d}$ = 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 |

# 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 |

Training set (rows 2104/400 through 1985/300)

Cross validation set (CV) (rows 1534/315 and 1427/199)

test set (rows 1380/212 and 1494/243)

# 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 (rows 2104–1985)

20% — Cross validation set (CV) (rows 1534, 1427)

20% — test set (rows 1380, 1494)

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% } 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 example $(x_{cv}^{(i)}, y_{cv}^{(i)})$

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

$m_{test}$

# Train/validation/test error

Training error:

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

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:

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

# Model selection

1. $h_\theta(x) = \theta_0 + \theta_1 x$
2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$

$\vdots$

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

# Model selection

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

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

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

$\vdots$

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

# Model selection

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

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

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

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

$\vdots$

10. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$ $\longrightarrow$ $\theta^{(10)} \longrightarrow J_{cv}(\theta^{(4)})$

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

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

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

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

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

# Model selection

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

$d=2$    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.    $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$    $\longrightarrow$    $\theta^{(10)} \longrightarrow J_{cv}(\theta^{(4)})$

$d = 4$

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

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

- Consider the model selection procedure where we choose the degree of polynomial using a cross validation set. For the final model (with parameters θ), we might generally expect $J_{CV}(\theta)$ to be lower than $J_{test}(\theta)$
  - An extra parameter (d, the degree of the polynomial) has been fit to the cross validation set.
  - An extra parameter (d, the degree of the polynomial) has been fit to the test set.
  - The cross validation set is usually smaller than the test set.
  - The cross validation set is usually larger than the test set.