# Notes : Week 6 part 1

## 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

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

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

- get more training examples.
- try smaller set of features.
- try getting additional features.
- try adding polynomial features $(x_1^2, x_2^2, x_1x_2, etc)$
- try decreasing $\lambda$
- try increasing $\lambda$

## Machine Learning Diagnostic:

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

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

# Notes: Week 6

## Advice for applying machine learning: evaluating a hypothesis

fails to generalize to new examples not in training set:

$$x_1 = \text{size of house}$$
$$x_2 = \text{no of bedrooms}$$
$$x_3 = \text{no of floors}$$
$$\vdots$$
$$x_{1000}$$

→ hard to imagine what this function even looks like.

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

## evaluating your hypothesis:

| size | price |
|------|-------|
| ~ | ~ |
| ~ | ~ |
| ~ | ~ |
| ~ | ~ |
| ~ | ~ |
| ~ | ~ |

70% } Training set →

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

30% } Test set →

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

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

## Training / testing procedure for [linear regression]

- Learn parameter $\theta$ from training data (minimizing training error $J(\theta)$) → 70%. take $\theta$ from traing set and plug it in here →

- Compute test set error.

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

# 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)

$$err(h_\Theta(x), y) = \begin{cases} 1 & \left( \begin{array}{l} \text{if } h_\Theta(x) \geq 0.5, \ y=0 \\ \text{or if } h_\Theta(x) < 0.5, \ y=1 \end{array} \right) \text{ error} \\ 0 \rightarrow \text{otherwise} \end{cases}$$

$\longrightarrow$ opposite to what we expect

$$\text{Test error} = \frac{1}{M_{test}} \sum_{i=1}^{M_{test}} err\left(h_\Theta(x_{test}^{(i)}), y_{test}^{(i)}\right)$$
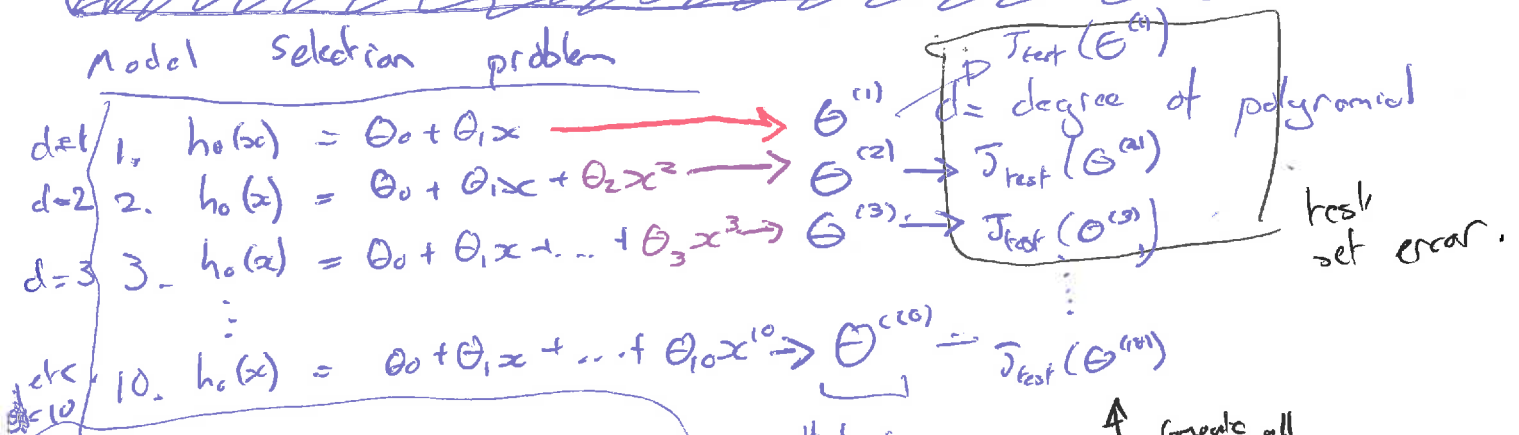
This is exactly the fraction of examples on the test set that the hypothesis has mislabelled.

# Notes: Week 6

## Advice for applying machine learning: Model selection and training / Validation / test sets

### Model Selection problem

$$\text{def} \begin{cases} 1. & h_\theta(x) = \theta_0 + \theta_1 x \\ d=2 \quad 2. & h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \\ d=3 \quad 3. & h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \\ \vdots \\ \text{etc} \quad 10. & h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \end{cases}$$

$\rightarrow \Theta^{(1)}$

$\rightarrow \Theta^{(2)} \rightarrow J_{test}(\Theta^{(2)})$

$\rightarrow \Theta^{(3)} \rightarrow J_{test}(\Theta^{(3)})$

$\Theta^{(10)} \rightarrow J_{test}(\Theta^{(10)})$

$$J_{test}(\Theta^{(1)})$$
$d =$ degree of polynomial

test set error.

thetas from different hypotheses.

Compute all $J_{test}(\Theta)$'s. See which model has the lowest test set error

so for this example choose:

$\theta_0 + \dots \theta_5 x^5$

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

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 the test set.

$\hookrightarrow$ just like in the extra practise exercise

in homework exercise 4. The NN overfit our training data ( the handwritten characters), however this highly optimized NN would not be all that useful for say a new dataset of handwritten characters
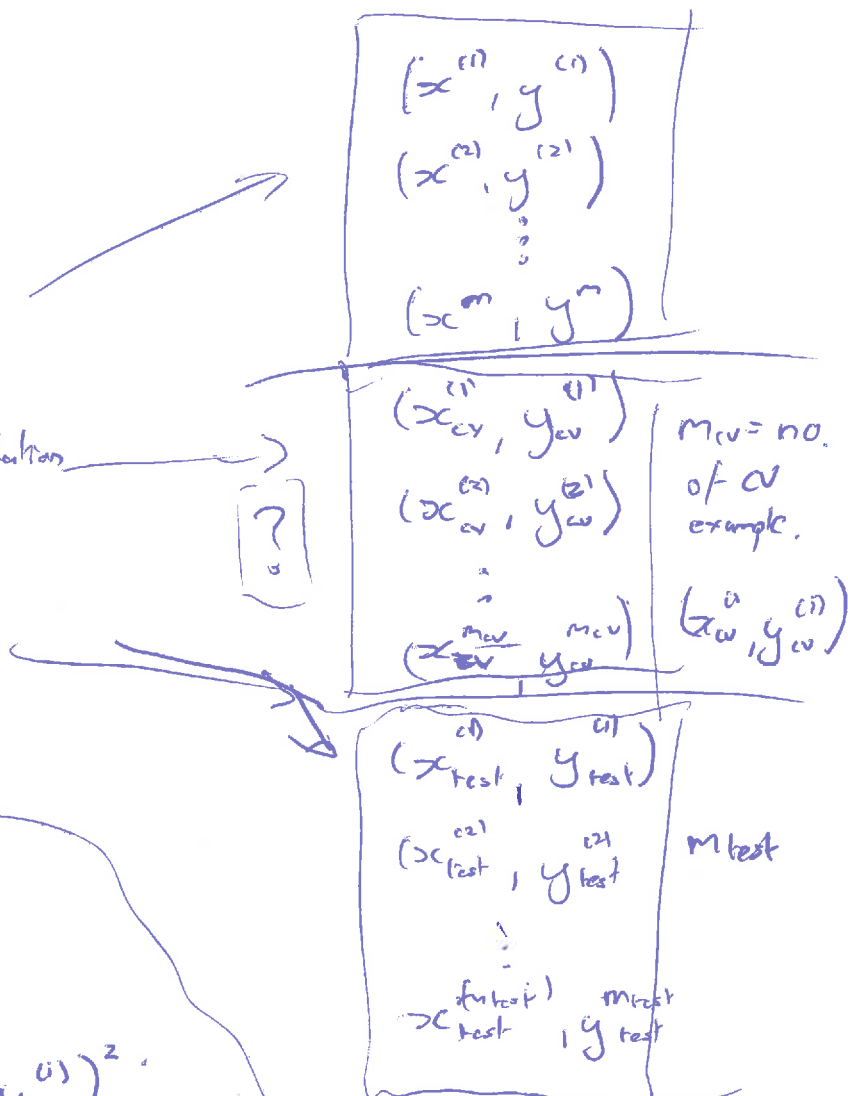
Advice for applying machine learning: Model selection and training / Validation / test sets

## evaluating your hypothesis
Dataset:

| Size | Price | |
|---|---|---|
| 60% | ~ ~ ~ ~ ~ | } Training set. |
| 20% | ~ ~ | } Cross Validation Set (CV) |
| 20% | ~ ~ ~ | } test set. |

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

$(x_{cv}^{(1)}, y_{cv}^{(1)})$     $m_{cv} = $ no.
$(x_{cv}^{(2)}, y_{cv}^{(2)})$     of CV
$\vdots$                            example.
$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$     $(x_{cv}^{(i)}, y_{cv}^{(i)})$

$\boxed{?}$

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

## Train / Validation / Test error

### Training error:

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

### Cross - Validation error

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

### Test error:

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

# Model Selection

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

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

    $\boxed{e.g}$     ⟵ $\boxed{C.V.}$

      Pick: $J_{cv}(\theta^{(4)})$   e.g.

d=10 10. $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \longrightarrow \theta^{10} \to J_{cv}(\theta^{10})$    lets assume that $\theta^4$ had the lowest error.

$\not{A}$ pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4$ ⟵    $d = 4$

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

# Notes : Week 6

> **Advice for applying machine learning : Diagnosing bias vs. Variance**

bias = underfitting
variance = overfitting

## bias / variance



size
high bias
(underfit)
$d=1$
$\Theta_0 + \Theta_1 x$

size
"Just right"
$d=2$.
$\Theta_0 + \Theta_1 x + \Theta_2 x^2$

High variance
(overfit)
$d=4$
$\Theta_0 + \Theta_1 x + \Theta_2 x^2$
$+\Theta_3 x^3 + \Theta_4 x^4$

## Bias / variance

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

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



$\leftarrow J_{cv}(\Theta)$

$J_{train}(\Theta)$

$d=1$    $d=2$    degree of polynomial $d$

# Diagnosing bias vs. variance

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



high bias (under fitting)

error

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

high variance (over fitting)

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

degree of polynomial d

Bias (underfit)

$J_{train}(\theta)$ will be high

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

Variance (overfit)

$J_{train}(\theta)$ will be low

$J_{cv}(\theta) \gg J_{train}$
& much greater than.

## Advice for applying machine learning: Regularization and bias/variance

### Linear regression with regularization

Model: $h_\theta(x) = \theta_0 + \theta_1 x \cdots \theta_4 x^4$

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

Large $\lambda$
High bias (underfit)

Intermediate $\lambda$
Just right

Small $\lambda$
High variance (overfit)

$\lambda = 10000 \quad \theta_1 = 0, \theta_2 = 0$

$\qquad h_\theta(a) \approx \theta_0$

$\boxed{\lambda = 0}$

### Choosing the regularization parameter $\lambda$

Model: $h_\theta(x) = \theta_0 + \theta_1 x \cdots \theta_4 x^4$

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

1) try $\lambda = 0$     $\to \min_\theta J(\theta) \to \theta^{(1)} \to J_{cv}(\theta^{(1)})$

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

3)    $\lambda = 0.02$       $\theta^{(3)}$

4)    $\lambda = 0.04$

5)    $\lambda = 0.08$    $\theta^{(5)} \to J_{cv}(\theta^{(5)})$

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

$\downarrow$ multiples of 2.

pick whichever $\lambda$ gives the lowest error of the cross-validation set.

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

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

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

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

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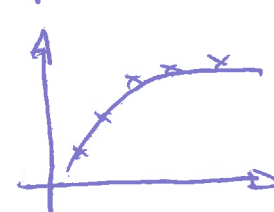
high variance

high bias

"just right"

$J_{cv}(\theta)$

$J_{train}(\theta)$

Small $\lambda$ (overfitting)

$\lambda$

Large $\lambda$ (underfitting)

# Advice for applying machine learning : learning curves

## learning curves

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

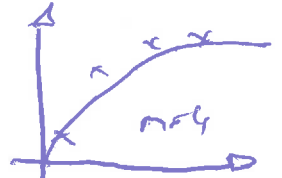$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} \left( h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)} \right)^2$$
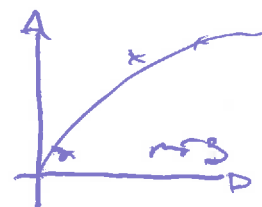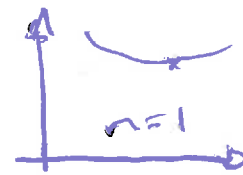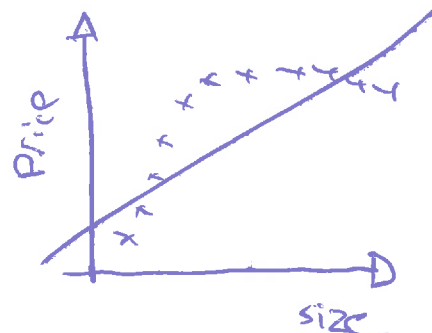


## High bias (underfitting)
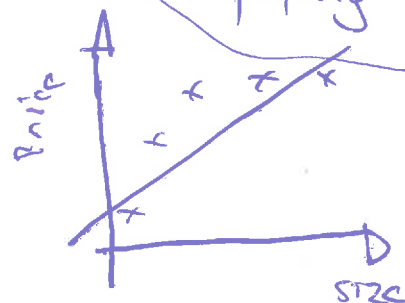


If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

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



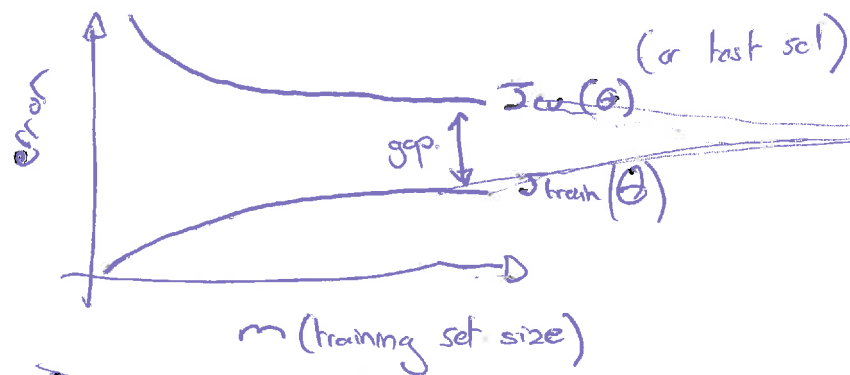harder and harder to find a quadratic function that fits the data perfectly as $d \to 0 \ldots \infty +$



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

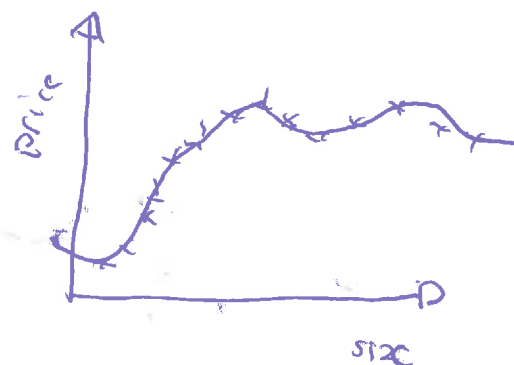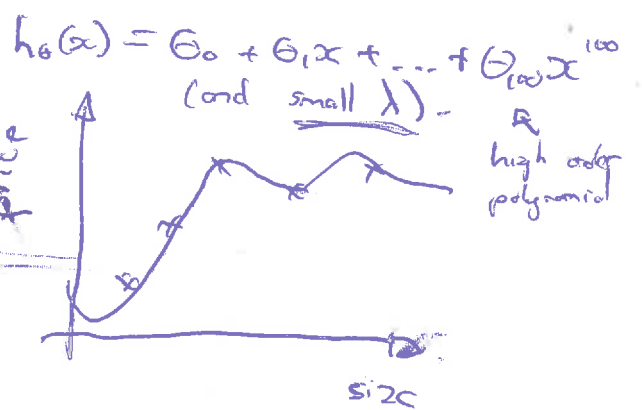straight line

more data
⇓
still underfitting (high bias)

# High variance



error

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

gap

$J_{train}(\theta)$

m (training set size)

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

high order polynomial

Price

size

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

Price

size

## Advice for applying machine learning: deciding what to try next (revisited)

### debugging a learning algorithm

- get more examples → fixes high variance
- try smaller sets of features → fixes high variance
- try getting additional features → fixing high bias
- try adding polynomial features → ~~adding polynomial features~~ fixes high bias
- try decreasing $\lambda$ → fixes high bias (underfitting)
- " increasing $\lambda$ → fixes high variance (overfitting)

### 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 ($\lambda$) to address overfitting