

# Week : 6

## Part One

### 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  $x_1, x_2, x_3, \dots, x_{100}$
- $\rightarrow$  - Try getting additional features
- Try adding polynomial features ( $\underline{x_1^2}, \underline{x_2^2}, \underline{x_1 x_2}$ , etc.)
- Try decreasing  $\lambda$
- Try increasing  $\lambda$

Andrew Ng

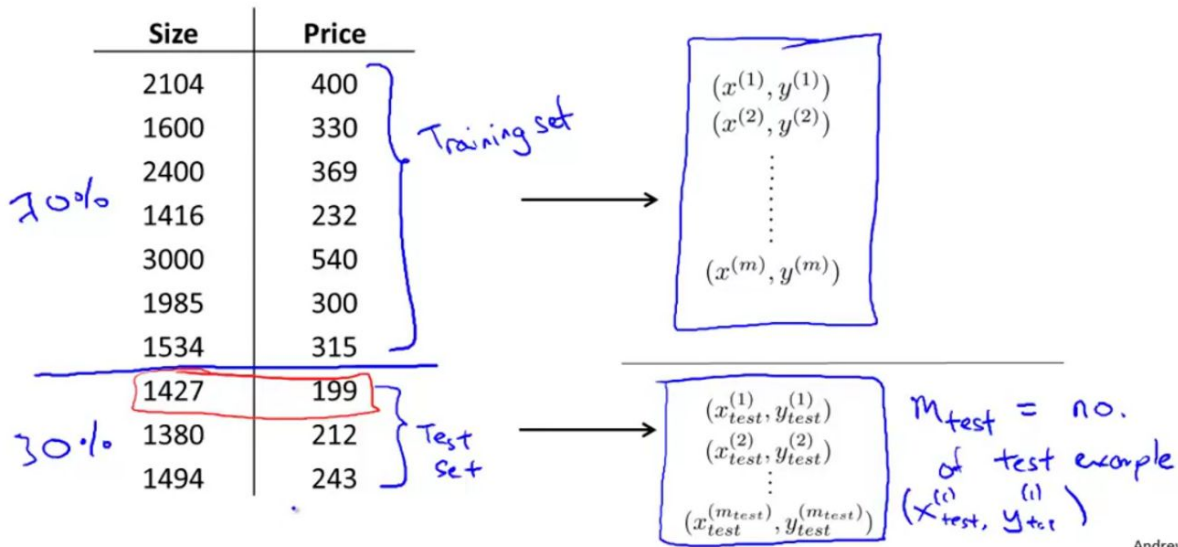
A hypothesis may have a low error for the training examples but still be inaccurate (because of overfitting). Thus, to evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: a **training set** and a **test set**. Typically, the training set consists of 70 % of your data and the test set is the remaining 30 %.

The new procedure using these two sets is then:

1. Learn  $\Theta$  and minimize  $J_{train}(\Theta)$  using the training set
2. Compute the test set error  $J_{test}(\Theta)$

## Evaluating your hypothesis

Dataset:



Andrew Ng

## Training/testing procedure for logistic regression

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

$m_{test}$

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

- Misclassification error (0/1 misclassification error):

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

0 otherwise

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

Andrew Ng

This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is:

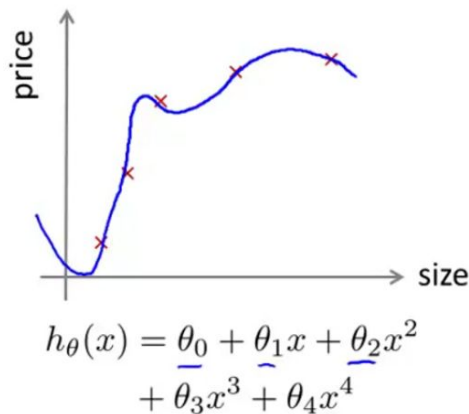
$$\text{Test Error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \text{err}(h_{\theta}(x_{test}^{(i)}), y_{test}^{(i)})$$

This gives us the proportion of the test data that was misclassified.

## What degree of polynomial should we choose?

Just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis. It could over fit and as a result your predictions on the test set would be poor.

### Overfitting example

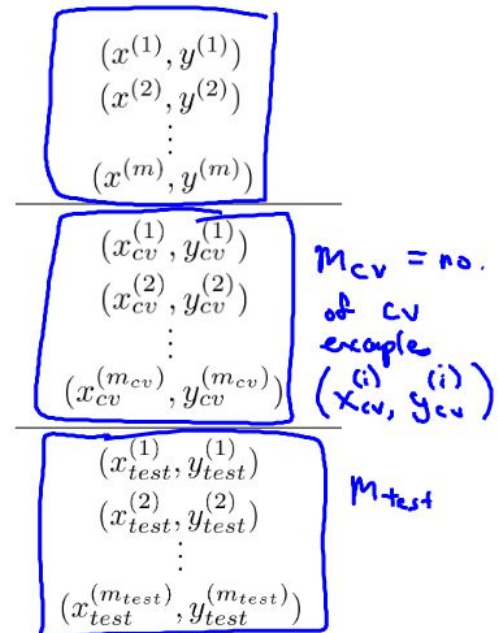


Once parameters  $\theta_0, \theta_1, \dots, \theta_4$  were fit to some set of data (training set), the error of the parameters as measured on that data (the training error  $J(\theta)$ ) is likely to be lower than the actual generalization error.

### Evaluating your hypothesis

Dataset:

	Size	Price	
60%	2104	400	Training set
	1600	330	
	2400	369	
	1416	232	
	3000	540	
	1985	300	
20%	1534	315	Cross validation set (CV)
	1427	199	
20%	1380	212	test set
	1494	243	



First we need to divide our data training set into three portions : Training set, cross validation set and test set.

- Training set: 60%
- Cross validation set: 20%
- Test set: 20%

We can now calculate three separate error values for the three different sets using the following method:

1. Optimize the parameters in  $\Theta$  using the training set for each polynomial degree.
2. Find the polynomial degree  $d$  with the least error using the cross validation set.
3. Estimate the generalization error using the test set with  $J_{test}(\Theta^{(d)})$ , ( $d$  = theta from polynomial with lower error);

This way, the degree of the polynomial  $d$  has not been trained using the test set.

### Train/validation/test error

Training error:

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

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:

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

## Model selection

$$\begin{array}{ll}
 \delta=1 & 1. \quad h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)}) \\
 \delta=2 & 2. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)}) \\
 \delta=3 & 3. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)}) \\
 \vdots & \vdots \\
 \delta=10 & 10. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})
 \end{array}$$

$\underline{\delta=4} \rightarrow$

Pick  $\theta_0 + \theta_1 x + \dots + \theta_4 x^4 \leftarrow$

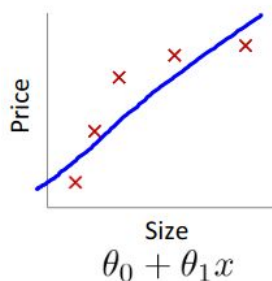
Estimate generalization error for test set  $\underline{J_{test}(\theta^{(4)})} \leftarrow$

Andrew Ng

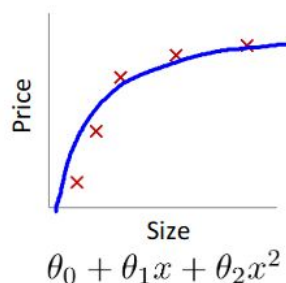
## Diagnosis bias vs variance :

We need to distinguish whether **bias** or **variance** is the problem contributing to bad predictions. High bias is underfitting and high variance is overfitting.

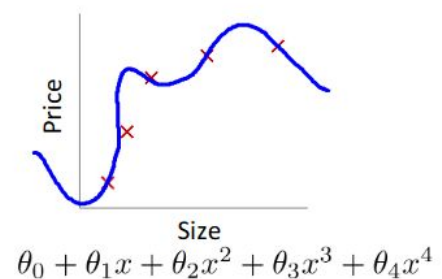
### Bias/variance



High bias  
(underfit)  
 $\delta=1$



"Just right"  
 $\delta=2$



High variance  
(overfit)  
 $\delta=4$

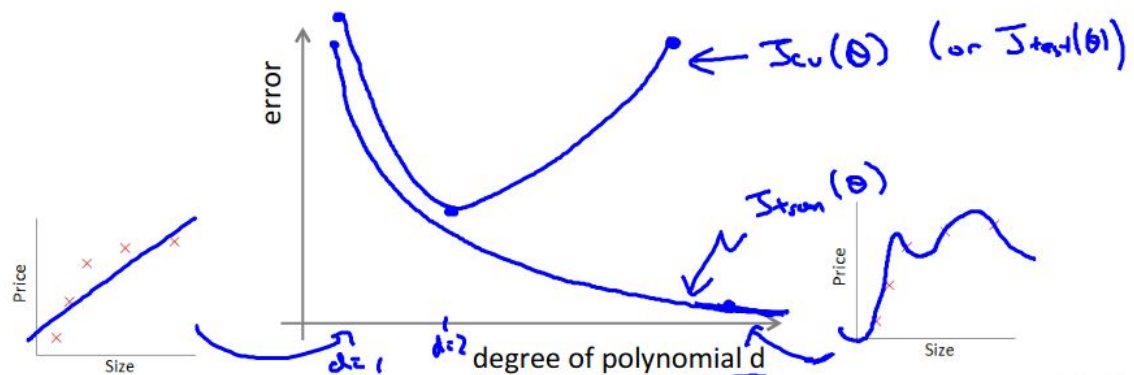
Andrew Ng



## Bias/variance

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

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

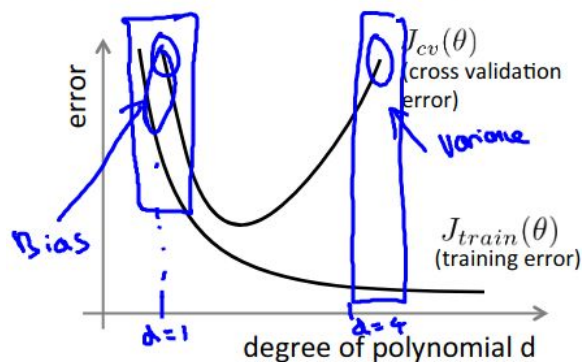


Andrew Ng

The training error will tend to decrease as we increase the degree  $d$  of the polynomial. At the same time, the cross validation error will tend to decrease as we increase  $d$  up to a point, and then it will increase as  $d$  is increased, forming a convex curve.

## 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):

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

Variance (overfit):

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

$\Rightarrow$

Andrew Ng

In an underfitting case, the equation will be very simple and the error difference of the training set and cross validation set will be nearly equal. But in an overfitting case, the

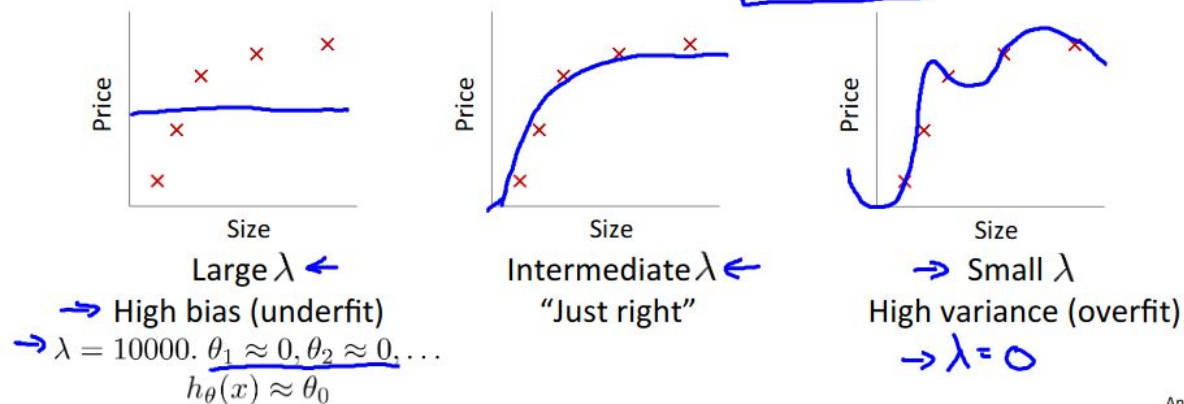
difference will be big. In overfitting case error of the training set will be low on the contrary for cross validation it will be high again.

## Regularization and Bias/Variance

### Linear regression with regularization

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

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



Andrew Ng

In the figure above, we see that as  $\lambda$  increases, our fit becomes more rigid. On the other hand, as  $\lambda$  approaches 0, we tend to over overfit the data. So how do we choose our parameter  $\lambda$  to get it 'just right' ? In order to choose the model and the regularization term  $\lambda$ , we need to:

1. Create a list of lambdas (i.e.  $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$ );
2. Create a set of models with different degrees or any other variants.
3. Iterate through the  $\lambda$  and for each  $\lambda$  go through all the models to learn some  $\Theta$ .
4. Compute the cross validation error using the learned  $\Theta$  (computed with  $\lambda$ ) on the  $J_{cv}(\Theta)$  **without** regularization or  $\lambda = 0$ .
5. Select the best combo that produces the lowest error on the cross validation set.

6. Using the best combo  $\Theta$  and  $\lambda$ , apply it on  $J_{Test}(\Theta)$  to see if it has a good generalization of the problem

### Choosing the regularization parameter $\lambda$

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

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

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

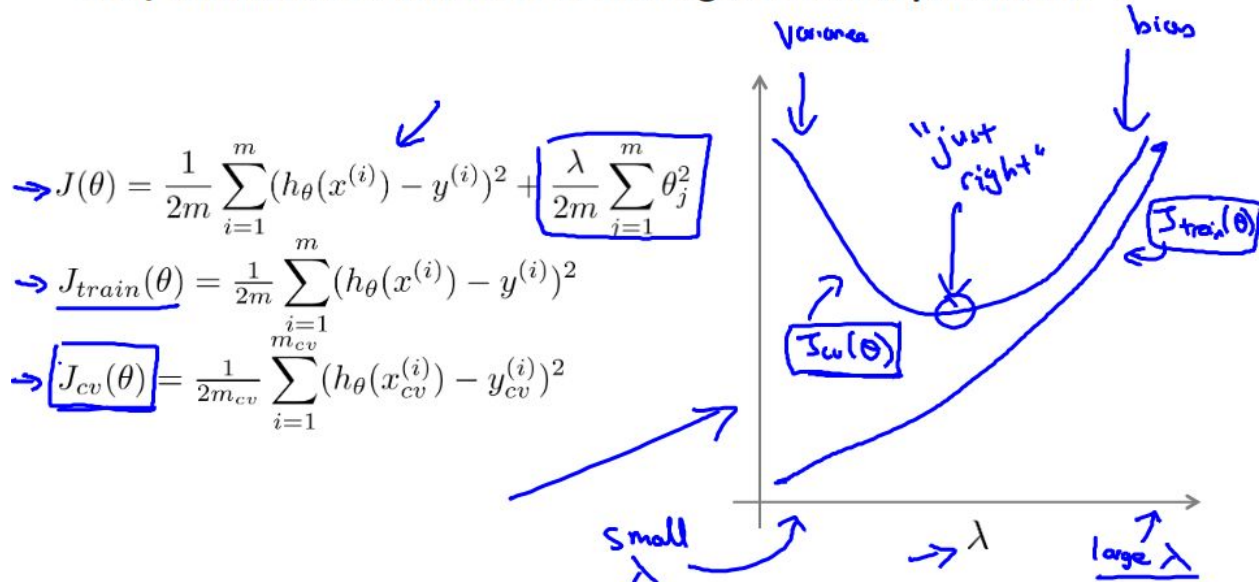
$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

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

$J(\theta)$   
 $J_{train}$   
 $J_{cv}$   
 $J_{test}$

Andrew Ng

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



Andrew Ng

At first I was confused to see the graph as there is no lambda parameter in the train and cv function but they change with value lambda. It is because we get the min theta from the first equation and then apply it on the next two equations. As a result when the value



of lambda is low it's impact on the first equation is not so big and our train set will fit well but cross validation set will not. Similarly when lambda is too large, it's contribution to the first equation is too big. As a result the next two equations will show a large error without the lambda part.

## Learning curves

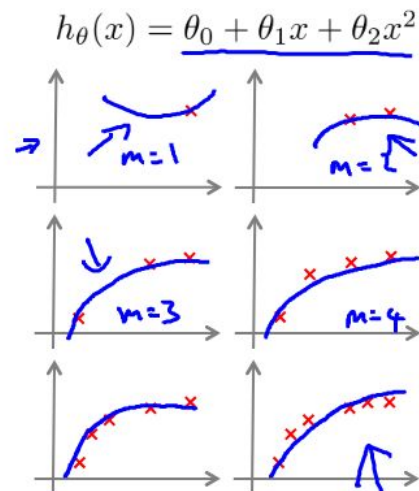
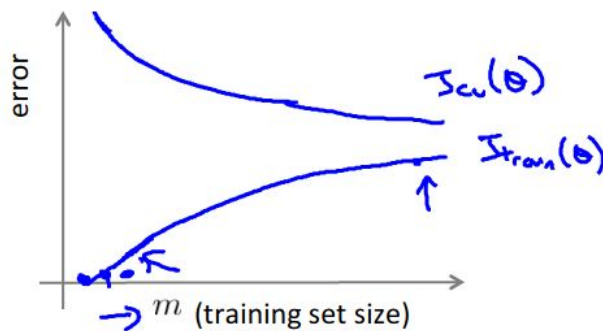
Training an algorithm on a very few number of data points (such as 1, 2 or 3) will easily have 0 errors because we can always find a quadratic curve that touches exactly those number of points. Hence:

- As the training set gets larger, the error for a quadratic function increases.
- The error value will plateau out after a certain m, or training set size.

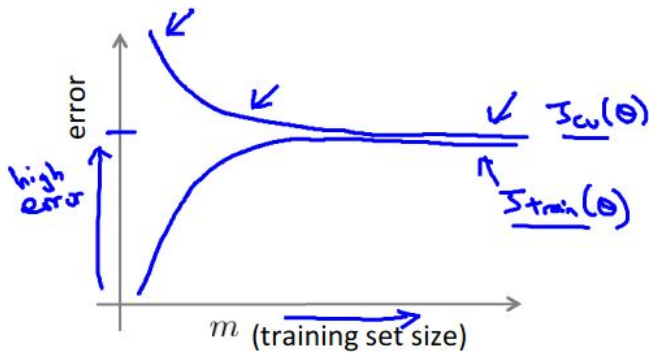
### Learning curves

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

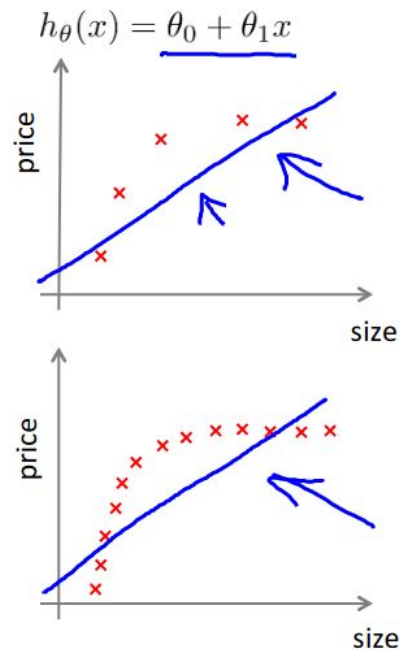
$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



## High bias

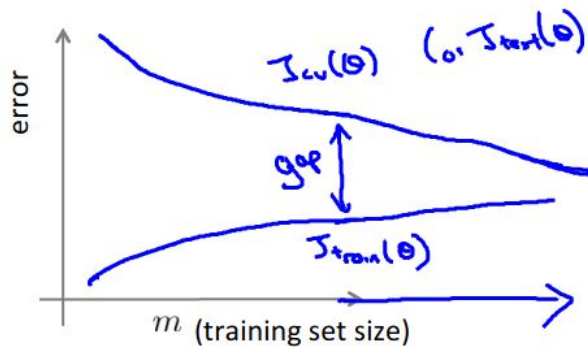


If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.



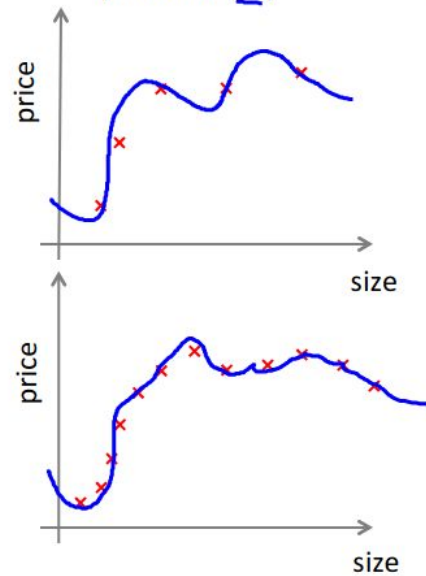
Andrew Ng

## High variance



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

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



Andrew Ng

# Summary of Part One:

## Debugging a learning algorithm:

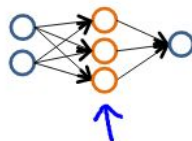
Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- 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 ( $x_1^2, x_2^2, x_1x_2$ , etc) → fixes high bias.
- Try decreasing  $\lambda$  → fixes high bias
- Try increasing  $\lambda$  → fixes high variance

Andrew Ng

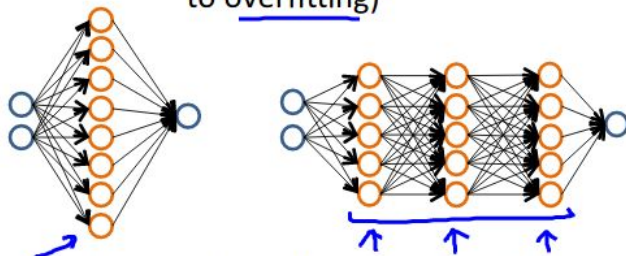
## Neural networks and overfitting

→ “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.

$J_{\text{reg}}(\theta)$  ↑

Andrew Ng

## Diagnosing Neural Networks

- A neural network with fewer parameters is **prone to underfitting**. It is also **computationally cheaper**.
- A large neural network with more parameters is **prone to overfitting**. It is also **computationally expensive**. In this case you can use regularization (increase  $\lambda$ ) to address the overfitting.

Using a single hidden layer is a good starting default. You can train your neural network on a number of hidden layers using your cross validation set. You can then select the one that performs best.

### Model Complexity Effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly consistently.
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

## Part Two : ML System Design

### Decide an email is spam or not:

To decide if an email is spam or not, we have some way. One way is to collect so many emails and select which ones of those are spam and which are not. From the spam set we will collect a set of words and will use them to differentiate new email.

---

#### Building a spam classifier

Supervised learning.  $x$  = features of email.  $y$  = spam (1) or not spam (0).

Features  $x$ : Choose 100 words indicative of spam/not spam.

E.g. deal, buy, discount, andrew, now, ...

$$X = \begin{bmatrix} 0 & \text{andrew} \\ 1 & \text{buy} \\ 1 & \text{deal} \\ 0 & \text{discount} \\ \vdots & \vdots \\ 1 & \text{now} \\ \vdots & \vdots \end{bmatrix} \quad x \in \mathbb{R}^{100}$$

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise} \end{cases}$$

From: cheapsales@buystufffromme.com  
To: ang@cs.stanford.edu  
Subject: Buy now!

Deal of the week! Buy now!

Note: In practice, take most frequently occurring  $n$  words (10,000 to 50,000) in training set, rather than manually pick 100 words.

Andrew Ng

Besides this, we have some other way to do this :

#### Building a spam classifier

How to spend your time to make it have low error?

- Collect lots of data
  - E.g. "honeypot" project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should "discount" and "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. m0rtgage, med1cine, w4tches.)



## Error analysis :

### Recommended approach

- Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross-validation data.
- Plot learning curves to decide if more data, more features, etc. are likely to help.
- Error analysis: Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

### Error Analysis

$m_{CV} = 500$  examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

- (i) What type of email it is *pharma, replica, steal passwords, ...*
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma:	12	→ Deliberate misspellings:	5
Replica/fake:	4		(m0rgage, med1cine, etc.)
→ Steal passwords:	53	→ Unusual email routing:	16
Other:	31	→ Unusual (spamming) punctuation:	32

Andrew Ng

It is very important to get error results as a single, numerical value. Otherwise it is difficult to assess your algorithm's performance. For example if we use stemming, which is the process of treating the same word with different forms (fail/failing/failed) as one word (fail), and get a 3% error rate instead of 5%, then we should definitely add it to our model. However, if we try to distinguish between uppercase and lowercase letters and end up getting a 3.2% error rate instead of 3%, then we should avoid using this new feature. Hence, we should try new things, get a numerical value for our error rate, and based on our result decide whether we want to keep the new feature or not.

## The importance of numerical evaluation

Should discount/discounts/discounted/discounting be treated as the same word?

Can use “stemming” software (E.g. “Porter stemmer”)  
universe/university.

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

Need numerical evaluation (e.g., cross validation error) of algorithm’s performance with and without stemming.

Without stemming: 5% error With stemming: 3% error

Distinguish upper vs. lower case (Mom/mom): 3.2%

## Error matrix for skewed classes :

If we have some skewed or biased training sets then error rate can’t help us to understand how much correct our algorithm is.

### Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

→ skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

→ 0.5% error

→ 99.2% accuracy (0.8% error)  
→ 99.5% accuracy (0.5% error)

To understand our algorithm’s accuracy we have two formulas by which we can measure accuracy of our algorithm more correctly.

## Precision/Recall

$y = 1$  in presence of rare class that we want to detect

	Actual class	
Predicted class	1	0
	True positive False positive	False negative True negative

$y = 0$   
recall = 0

→ Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

$$\frac{\text{True positives}}{\text{\#predicted positive}} = \frac{\text{True positive}}{\text{True pos + False pos}}$$

→ Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positives}}{\text{\#actual positives}} = \frac{\text{True positives}}{\text{True pos + False neg}}$$

## Trading of Precision and Recall :

Sometimes we need to tell our client if their test is negative or positive correctly. In that case we can get help from precision and recall. When we want to be confident then we will choose high precision and lower recall. Again when we want to avoid false negative result we will choose higher recall and lower precision.

### Trading off precision and recall

→ Logistic regression:  $0 \leq h_{\theta}(x) \leq 1$

Predict 1 if  $h_{\theta}(x) \geq 0.5$  ~~0.5~~ ~~0.7~~ ~~0.9~~ ~~0.3~~ ←

Predict 0 if  $h_{\theta}(x) < 0.5$  ~~0.5~~ ~~0.7~~ ~~0.9~~ ~~0.3~~

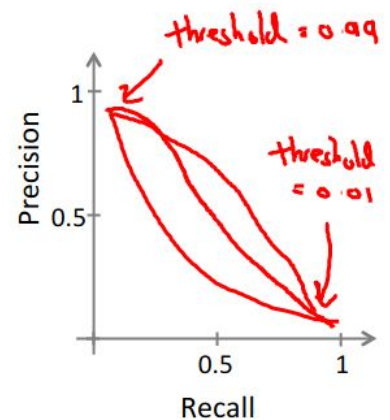
→ Suppose we want to predict  $y = 1$  (cancer) only if very confident.

→ Higher precision, lower recall

→ Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ Higher recall, lower precision.

$$\begin{aligned} \rightarrow \text{precision} &= \frac{\text{true positives}}{\text{no. of predicted positive}} \\ \rightarrow \text{recall} &= \frac{\text{true positives}}{\text{no. of actual positive}} \end{aligned}$$



More generally: Predict 1 if  $h_{\theta}(x) \geq \text{threshold}$  ←

We need to choose one algorithm over others. In that case, F score will help us. We will apply the f score formula on cross validation sets and pick that algorithm with maximum f score

## F<sub>1</sub> Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)
→ Algorithm 1	<u>0.5</u>	<u>0.4</u>
→ Algorithm 2	<u>0.7</u>	<u>0.1</u>
Algorithm 3	<u>0.02</u>	1.0

Average:  ~~$\frac{P+R}{2}$~~   
 F<sub>1</sub> Score:  $2 \frac{PR}{P+R}$

Predict y=1 all the time  
 P=0 or R=0 ⇒ F-score = 0  
 P=1 and R=1 ⇒ F-score = 1

Andrew Ng

## Data for Machine learning :

With small training set different algorithms act significantly differently. But with large number of training set different algorithms can predict with nearly same training error.

### Designing a high accuracy learning system

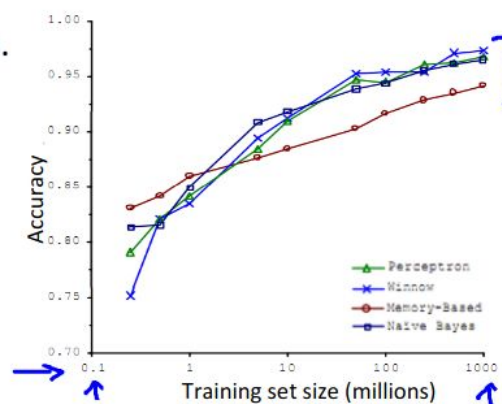
E.g. Classify between confusable words.

{to, two, too}, {then, than}

→ For breakfast I ate two eggs.

Algorithms

- - Perceptron (Logistic regression)
- - Winnow
- - Memory-based
- - Naïve Bayes



“It’s not who has the best algorithm that wins.

It’s who has the most data.”

[Banko and Brill, 2001]



## Large data rationale

- Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms. ←

→  $J_{\text{train}}(\theta)$  will be small.

Use a very large training set (unlikely to overfit) low variance ←

→  $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

→  $J_{\text{test}}(\theta)$  will be small