# Week 5: Neural Networks Learning

## *Cost Function*

- We are going to talk about the application of neural networks to classification problems. First we define $L$ as the total number of layers in a network and $s_j$ as the number of units in a layer. We will focus in 2 classification problems: Binary where the output layer = 1 and multiclass classification (K classes) where there are K outputs and $K \geq 3$.

**Neural Network (Classification)**

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

$L = $ total no. of layers in network     $L = 4$

$s_l = $ no. of units (not counting bias unit) in layer $l$     $S_1 = 3, \ S_2 = 5, \ S_4 = S_L = 4$

**Binary classification**

$y = 0$ or $1$ ←

$h_\Theta(x)$

1 output unit ←

$h_\Theta(x) \in \mathbb{R}$

$S_L = 1.$     $K = 1$ ←

**Multi-class classification** (K classes)

$y \in \mathbb{R}^K$  E.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ ←

pedestrian  car  motorcycle  truck

K output units

$h_\Theta(x) \in \mathbb{R}^k$

$S_L = K$     $(k \geq 3)$

Andrew Ng

- The cost function for a neural network is a generalization of the logistic regression regularized cost function. The differences is that first we have two sums where the first sum is over the all the neurons in one layer (sum from K=1 to K) and the regularized term iterates over every neuron except the bias factor neurons. In this case $h_\theta(x)_k$ is the hypothesis that results in the kth output

**Cost function**

Logistic regression:

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

Neural network:

$h_\Theta(x) \in \mathbb{R}^K$   $(h_\Theta(x))_i = i^{th}$ output

$$J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k)\right]$$
$$+ \frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_{ji}^{(l)})^2$$

$\Theta_{ji}^{(l)}$

$j = 0 \ldots S_\ell$

$\Theta_{io}^{(1)} x_o + \Theta_{ii}^{(1)} x_1 + \ldots$

$a_o$

$y_k \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

Andrew Ng

The double sum simply adds up the logistic regression costs calculated for each cell in the output layer. The triple sum adds up the squares of all the individual $\theta s$ in the entire network

## *Backpropagation*

- We are going to talk about an algorithm to minimize the cost function: the backpropagation algorithm. What we would like to do is find parameters $\theta$ that minimize the $J(\theta)$ cost function shown below. In order to use gradient descent or a more advanced optimization algorithm we need to write code that receives $\theta$ and computes $J(\theta)$ and the partial derivatives. We already have the $J(\theta)$ cost function defined so we will concentrate in this class on how to compute the partial derivatives

**Gradient computation**

$$\rightarrow J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right]$$
$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

$$\rightarrow \min_{\Theta} J(\Theta)$$

Need code to compute:

$$\rightarrow - J(\Theta)$$
$$\rightarrow - \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) \leftarrow$$

$\Theta_{ij}^{(l)} \in \mathbb{R}$

- 

- Let's take the example where we only have one training example, an $(x, y)$ pair. Given this one training example we apply forward propagation as shown below. $a^{(1)} = x$, then we move on the 2nd layer which is the sigmoid of the multiplication of the $\theta$ matrix * $a^{(1)}$, we then have $z^{(2)} = \theta^{(1)} * a^{(1)}$ and we then apply the sigmoid to get $a^{(2)} = g(z^{(2)})$ (remember to add the bias factor).

**Gradient computation**

Given one training example $(x, y)$:

Forward propagation:

$$a^{(1)} = x$$
$$\rightarrow z^{(2)} = \Theta^{(1)} a^{(1)}$$
$$\rightarrow a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$
$$\rightarrow z^{(3)} = \Theta^{(2)} a^{(2)}$$
$$\rightarrow a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$
$$\rightarrow z^{(4)} = \Theta^{(3)} a^{(3)}$$
$$\rightarrow a^{(4)} = h_\Theta(x) = g(z^{(4)})$$

Layer 1   Layer 2   Layer 3   Layer 4

$a^{(1)}$   $a^{(2)}$   $a^{(3)}$   $a^{(4)}$

We do the same calculations for the remaining layers and our last layer ends up being
$a^{(4)} = h_\theta(x) = g(\theta^{(3)} a^{(3)})$

- We will use backpropagation to compute the derivatives. It is called backpropagation because we start by computing our errors in the output layer and we work backwards by 'propagating' our errors to the previous layers. Our intuition here is that $\delta_j^{(l)} = error\ of\ node\ j\ in\ layer\ l$. For a neural network with 4 layers, like the one shown below, we can start with the output layer and $\delta_j^{(4)} = a_j^{(4)} - y_j$. Remember that $a_j^{(4)} = h_\theta(x)$ and the vectorized implementation is then $\delta^{(4)} = a^{(4)} - y$. Now we have the value for the derivative for the output layer.



## Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)} = $ "error" of node $j$ in layer $l$.

For each output unit (layer L = 4)

$\delta_j^{(4)} = a_j^{(4)} - y_j$   $(h_\theta(x))_j$;  $\delta^{(4)} = a^{(4)} - y$

Layer 1   Layer 2   Layer 3   Layer 4

$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} .* g'(z^{(3)})$   $a^{(3)} .* (1 - a^{(2)})$

$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} .* g'(z^{(2)})$   $a^{(2)} .* (1 - a^{(2)})$

(No $\delta^{(1)}$)   $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$   ( ignoring $\lambda$; if $\lambda = 0$)

Then we go backwards and compute the error/partial derivative for the previous layers. For $\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$ and we compute $g'$ as $a^{(3)} .* (1 - a^{(3)})$ .For $\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$ and $g'$ is $a^{(2)} .* (1 - a^{(2)})$. We can prove this mathematically (assuming there is no regularization $\lambda$ or $\lambda = 0$) but that is outside the scope of the class. Finally notice there is no $\delta^{(1)}$ for the input layer because those values are the features we observed in our training set and those do not have any error

- Generalizating for all kinds of training sets we have the following: given a training set of (x,y) value pairs
  - we set $\Delta_{ij}^{(l)} = 0\ (for\ all\ i,j,\ l)$. These are our accumulators
  - We then do a for loop $for\ i = 1.. m$
    - Set $a^{(1)} = x^{(i)}$
    - Perform forward propagation to compute $a^{(l)}$ for all layers
    - Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$
    - Do backpropagation to compute $\delta^{(L-1)}, \delta^{(L-2)}... \delta^{(2)}$
    - $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(L)} \delta_i^{(L+1)}$
  - Outside the loop we then have

- $D_{ij}^{(L)} := 1/m \; \Delta_{ij}^{(l)} + \lambda\theta_{ij}^{(L)}$ if j <> 0
- $D_{ij}^{(L)} := 1/m \; \Delta_{ij}^{(l)}$        if j = 0

## Backpropagation algorithm

→ Training set $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$

Set $\triangle_{ij}^{(l)} = 0$ (for all $l, i, j$).    (use to compute $\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta)$)

For $i = 1$ to $m$ ←   $(x^{(i)}, y^{(i)})$.

Set $a^{(1)} = x^{(i)}$

Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \ldots, L$

Using $y^{(i)}$, compute $\delta^{(L)} = \boxed{a^{(L)}} - \boxed{y^{(i)}}$

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \ldots, \delta^{(2)}$   $\delta^{(1)}$ ✗

$\boxed{\triangle_{ij}^{(l)} := \triangle_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}}$ ←     $\triangle^{(\ell)} := \triangle^{(\ell)} + \delta^{(\ell+1)} (a^{(\ell)})^T$.

→ $\boxed{D_{ij}^{(l)}} := \frac{1}{m} \triangle_{ij}^{(l)} + \lambda\Theta_{ij}^{(l)}$ if $j \neq 0$      $\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

→ $\boxed{D_{ij}^{(l)}} := \frac{1}{m} \triangle_{ij}^{(l)}$      if $j = 0$

---

*Backpropagation Intuition*

- To get a better intuition and understanding of what backpropagation is doing let's start by reviewing how forward propagation works. In the neural network below we know how to get the values for each node: it is $z$ to which we apply the sigmoid function to get the activation value $a$ with the appropriate indices. Now let's concentrate on the 2nd neuron of the 2nd hidden layer and let's see how we compute the value there.



**Forward Propagation**

We know there is a $\theta$ vector in between the 2 layers and we know the other input is the appropriate $a$ value from the previous layer so for that neuron we compute its z value as

$z_1^{(3)} = \theta_{1,0}^{(2)} * 1 + \theta_{1,1}^{(2)} * a_1^{(2)} + \theta_{1,2}^{(2)} * a_2^{(2)}$ . We did these calculations in a forward way, backpropagation is doing something very similar but from right to left

- What is backpropagation doing? Let's start by examining the cost function for a single training example $(x^{(i)}, y^{(i)})$ with a 1 output unit and with no regularization ($\lambda = 0$ ). In this case the cost function is $cost(i) = y^{(i)} log(h_\theta(x^{(i)})) + (1 - y^{(i)}) log(1 - h_\theta(x^{(i)}))$ . For intuition purposes, we can also think of the cost as the square error which would give us $cost(i) = (h_\theta(x^{(i)} - y^{(i)})^2$

### What is backpropagation doing?

$$J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - (h_\Theta(x^{(i)})))\right]$$

$$+\frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_{ji}^{(l)})^2$$

$(x^{(i)}, y^{(i)})$

Focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$ ),

$$cost(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$$

(Think of $cost(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2$ )

I.e. how well is the network doing on example i?

- Now let's take a shot at see how backpropagation works. What we want to do is compute the $\delta$ values for each neuron in the network. We start with the output layer (working backwards) and that value is $\delta_1^{(4)} = a_1^{(4)} - y^{(i)}$ which is the cost error for $a_1^{(4)}$ .



$\delta_1^{(4)} = y^{(i)} - a_i^{(4)}$

$\delta_2^{(2)} = \theta_{12}^{(i)} \delta_1^{(3)} + \theta_{12}^{(2)} \delta_2^{(3)}$

**Forward Propagation**

$\delta_2^{(3)} = \theta_{12}^{(3)} \cdot \delta_1^{(4)}$ .

$\delta_j^{(l)} = $ "error" of cost for $a_j^{(l)}$ (unit $j$ in layer $l$).

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} cost(i)$ (for $j \geq 0$), where $cost(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$

Andrew Ng

Now to compute a $\delta$ value for one of the nodes in the hidden layers remember there is a $\theta$ matrix in between layers. We are going to use that $\theta$ vector and the 'previous' $\delta$ values to calculate the cost error for that node. Let's take the 2nd node in the first hidden layer

so $\delta_2^{(2)} = \theta_{1,2}^{(2)} * \delta_1^{(3)} + \theta_{2,2}^{(2)} * \delta_2^{(3)}$. Similarly, we can compute the error cost for the 2nd node in the 3rd layer as follows: $\delta_2^{(3)} = \theta_{1,2}^{(3)} * \delta_1^{(4)}$

- We need to be concerned about an implementation detail: some of the parameters we use in advanced optimization functions expect vectors as inputs and outputs. However, neural networks have matrices: $\theta^{(1)} \ldots \theta^{(n)}$ and $D^{(1)} \ldots D^{(n)}$ so we need to 'unroll' our matrices into vectors so they can be used with these advanced optimization algorithms

**Advanced optimization**

```
function [jVal, gradient] = costFunction(theta)
    ...
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network (L=4):
$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (`Theta1`, `Theta2`, `Theta3`)
$D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (`D1`, `D2`, `D3`)
"Unroll" into vectors

- Let's see how we do this using an example. We have a neural network with 10 nodes per layer and 1 output unit. The dimensions of the $\theta$ matrices and the $D$ matrices are shown below. We 'unroll' the theta matrices into vectors using the octave command thetaVec = [ Theta1(:); Theta2(:); Theta3(:)]; and we bring back the matrices from a vector by using the following commands:
    - *Theta1 = reshape(thetaVec(1:110), 10, 11);*
    - *Theta2 = reshape(thetaVec(111:220), 10,11);*
    - *Theta3 = reshape(thetaVec(221:231), 10, 11);*

**Example**

$s_1 = 10, s_2 = 10, s_3 = 1$

$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$

$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$

```
thetaVec = [ Theta1(:); Theta2(:); Theta3(:)];
DVec = [D1(:); D2(:); D3(:)];

Theta1 = reshape(thetaVec(1:110),10,11);
Theta2 = reshape(thetaVec(111:220),10,11);
Theta3 = reshape(thetaVec(221:231),1,11);
```

- Here's how we are going to use these commands in a real situation:
    - We have the initial theta matrices
    - We unroll them into $thetaVec$ so it can be passed to our cost function
    - Inside the cost function I get the individual theta matrices using reshape
    - I use forward and back propagation to compute the $D$ matrices and $J(\theta)$
    - I unroll the $D$ matrices into a $gradientVec$ which is what is returned by the cost function

**Learning Algorithm**

→ Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.

→ Unroll to get `initialTheta` to pass to

→ `fminunc(@costFunction, initialTheta, options)`

```
function [jval, gradientVec] = costFunction(thetaVec)
```
→ From `thetaVec`, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$    reshape

→ Use forward prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.

Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get `gradientVec`.

*Gradient Checking*

- Backprop is susceptible to some subtle bugs that are hard to spot and that you might not know about. Your algorithm could be working but you might end up with a bad answer. Because of this, the professor recommends to always implement gradient checking when using backprop
- Concretely this is how we check the gradients. We take a point $\theta$ in the $J(\theta)$ curve and create 2 points: $\theta + \varepsilon$ and $\theta - \varepsilon$ and draw a line between the 2; that line is our approximation to the derivative of our cost function. Mathematically, the slope of that line is the height divided by the width The corresponding $J(\theta)$ points for those new theta values are $J(\theta + \varepsilon)$ and $J(\theta - \varepsilon)$. To get the value of slope we then take the height which is $J(\theta + \varepsilon) - J(\theta - \varepsilon)$ divided by the width which is $2\varepsilon$. My equation for my approximation for the derivative at that $\theta$ points is then $(J(\theta + \varepsilon) - J(\theta - \varepsilon))/2\varepsilon$. The $\varepsilon$ value is usually very small; the professor uses $10^{-4}$. Some people might have seen an alternate equation, the one-sided difference equation which is $(J(\theta + \varepsilon) - J(\theta))/\varepsilon$; ours is called the two-sided difference equation and it gives us a slightly better result so we use that equation. Below we also see the octave code to use for this equation

**Numerical estimation of gradients**



$J(\Theta+\varepsilon)$
$J(\Theta)$
$J(\Theta-\varepsilon)$
$\dfrac{J(\Theta+\varepsilon)-J(\Theta-\varepsilon)}{\Theta}, \Theta \in \mathbb{R}$

$2\varepsilon$

$\Theta-\varepsilon \quad \Theta \quad \Theta+\varepsilon$

$\dfrac{d}{d\Theta}J(\Theta) \approx \dfrac{J(\Theta+\varepsilon)-J(\Theta-\varepsilon)}{2\varepsilon}$ ← $\dfrac{J(\Theta+\varepsilon)-J(\Theta)}{\varepsilon}$ ✗

$\varepsilon = 10^{-4}$ ←

Implement: `gradApprox = (J(theta + EPSILON) - J(theta - EPSILON))`
`/ (2*EPSILON)`

- The previous example was for the case when θ was a rolled number. Now let's look at the case when θ is a vector of parameters. In this case we use the same logic as in the previous slide, the difference is that here we have a vector so we do it in a one by one basis. For example, to approximate the partial derivative for the cost function with respect to $\theta_1$ we use the formula
$J(\theta) \approx J((\theta_1 + \varepsilon, \theta_2, \theta_3, \theta_4...\theta_n) - J(\theta_1 - \varepsilon, \theta_2, \theta_3, \theta_4...\theta_n) / 2\varepsilon$ and we just repeat the formula for each θ value

**Parameter vector** $\theta$

$\rightarrow \theta \in \mathbb{R}^n$  (E.g. $\theta$ is "unrolled" version of $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ )

$\rightarrow \theta = [\theta_1, \theta_2, \theta_3, \ldots, \theta_n]$

$\rightarrow \dfrac{\partial}{\partial \theta_1}J(\theta) \approx \dfrac{J(\theta_1+\varepsilon,\theta_2,\theta_3,\ldots,\theta_n)-J(\theta_1-\varepsilon,\theta_2,\theta_3,\ldots,\theta_n)}{2\varepsilon}$

$\rightarrow \dfrac{\partial}{\partial \theta_2}J(\theta) \approx \dfrac{J(\theta_1,\theta_2+\varepsilon,\theta_3,\ldots,\theta_n)-J(\theta_1,\theta_2-\varepsilon,\theta_3,\ldots,\theta_n)}{2\varepsilon}$

$\vdots$

$\rightarrow \dfrac{\partial}{\partial \theta_n}J(\theta) \approx \dfrac{J(\theta_1,\theta_2,\theta_3,\ldots,\theta_n+\varepsilon)-J(\theta_1,\theta_2,\theta_3,\ldots,\theta_n-\varepsilon)}{2\varepsilon}$

- What we actually implement in octave is shown below. At the end of the for loop we have a *gradApprox* value. We can compare that value with the DVec that we got from backprop and if the difference between the two is just a few decimal points I can be confident my backprop algorithm is working correctly

```
for i = 1:n, ←
    [ thetaPlus = theta;
    [ thetaPlus(i) = thetaPlus(i) + EPSILON;
    [ thetaMinus = theta;
    [ thetaMinus(i) = thetaMinus(i) - EPSILON;
    [ gradApprox(i) = (J(thetaPlus) - J(thetaMinus))
                    /(2*EPSILON);
end;
```

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_i + \varepsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta_i - \varepsilon$$

$$\frac{\partial}{\partial \theta_i} J(\theta).$$

Check that $\texttt{gradApprox} \approx \texttt{DVec}$

From backprop.

- When implementing this we need to do the following steps:
    - Implement backprop to compute DVec
    - Implement gradient check to compute gradApprox
    - Make sure they give similar values
    - Turn off gradient checking before using backprop code for learning. We do this because backprop is very fast but gradient checking is very slow. If we leave it 'on' our code will be very slow

### Implementation Note:
- Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- Implement numerical gradient check to compute gradApprox.
- Make sure they give similar values.
- Turn off gradient checking. Using backprop code for learning.

↳ DVec

### Important:
- Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of costFunction(…))your code will be very slow.

- The octave code that implements gradient checking is shown below for reference.

```
1   epsilon = 1e-4;
2   for i = 1:n,
3     thetaPlus = theta;
4     thetaPlus(i) += epsilon;
5     thetaMinus = theta;
6     thetaMinus(i) -= epsilon;
7     gradApprox(i) = (J(thetaPlus) - J(thetaMinus))/(2*epsilon)
8   end;
9
```

*Random Initialization*

- One last detail before putting everything together is to select the initial value of theta. In previous occasions we were fine with initializing theta to all zeros. However, that does not work for neural networks

**Initial value of** $\Theta$

For gradient descent and advanced optimization method, need initial value for $\Theta$.
```
optTheta = fminunc(@costFunction,
                   initialTheta, options)
```

Consider gradient descent
Set `initialTheta = zeros(n,1)` ?

- Why is this the case? Because what ends up happening is that the parameters corresponding to inputs into two hidden units are identical and we end up with not only say $a_1^{(2)} = a_2^{(2)}$ but also with $\delta_1^{(2)} = \delta_2^{(2)}$. In practice this fact means our neural network is not computing interesting calculations and therefore we end up with a not very interesting hypothesis $h_\theta(x)$. This problem is also known as symmetric weights

**Zero initialization**



$$\Rightarrow \Theta_{ij}^{(l)} = 0 \text{ for all } i, j, l.$$

$\rightarrow h_\Theta(x)$

$a_1^{(2)} = a_2^{(2)}$, Also $\delta_1^{(2)} = \delta_2^{(2)}$.

$\dfrac{\partial}{\partial \Theta_{01}^{(1)}} J(\Theta) = \dfrac{\partial}{\partial \Theta_{02}^{(1)}} J(\Theta)$   $\Theta_{01}^{(1)} = \Theta_{02}^{(1)}$

After each update, parameters corresponding to inputs going into each of two hidden units are identical.

$a_1^{(2)} = a_2^{(2)}$

---

- In order to get around this problem we initialize our parameters with random initialization. What we do is we will initialize our theta values to a random value between $-\varepsilon \text{ and } \varepsilon$. This epsilon value is different from the epsilon value we were using in gradient checking.

**Random initialization: Symmetry breaking**

$\rightarrow$ Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
(i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

E.g.   $\rightarrow$ random 10×11 matrix (betw. 0 and 1)

$\rightarrow$ Theta1 = rand(10,11)*(2*INIT_EPSILON) - INIT_EPSILON;   $[-\epsilon, \epsilon]$

$\rightarrow$ Theta2 = rand(1,11)*(2*INIT_EPSILON) - INIT_EPSILON;

The octave code to do this is shown above. We start a matrix of random numbers between 0 and 1 and then multiply it by $(2 * init\ epsilon) - initi\ epsilon$ and we do the same for all thetas and that's how we end up with all thetas with values between $[-\varepsilon, \varepsilon]$

- Summarizing, to train a neural network we should:
  - Initialize the weights to small values close to 0, $[-\varepsilon, \varepsilon]$
  - Implement back propagation
  - Do gradient checking
  - Use gradient descent (or another advanced optimization method) to train to minimize $J(\theta)$ as a function of $\theta$
- Below is some sample octave code we can use to experiment with this method

```
1  If the dimensions of Theta1 is 10x11, Theta2 is 10x11 and Theta3 is 1x11.
2
3  Theta1 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
4  Theta2 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
5  Theta3 = rand(1,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
6
```

*Putting it all together*

- One of the first steps in training a neural network is to pick a network architecture. You have to make a choice on the number of input & output units, the number of hidden layers and the number of units per hidden layer. In general we can say
  - No of input units is defined by the number of features, i.e. dimension of features
  - No of output units is the number of classes
  - A reasonable default number is 1 hidden layer. If we have to have > 1 hidden layer then have the same number of units in each hidden layer as shown in the sample architectures below



- Once we have a network architecture the next step is to train the neural network. There are 6 steps to train the neural network:
  - Randomly initialize weights with values close to, but not quite, zero
  - Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
  - Implement code to compute the cost function
  - Implement code to compute the partial derivative of the cost function
    - This code takes the form of a for loop where we perform forward and back propagation for a given $(x^{(i)}, y^{(i)})$
    - We get all the $a^{(l)}$ and all the $\delta^{(l)}$ and we add them to our $D$ accumulators
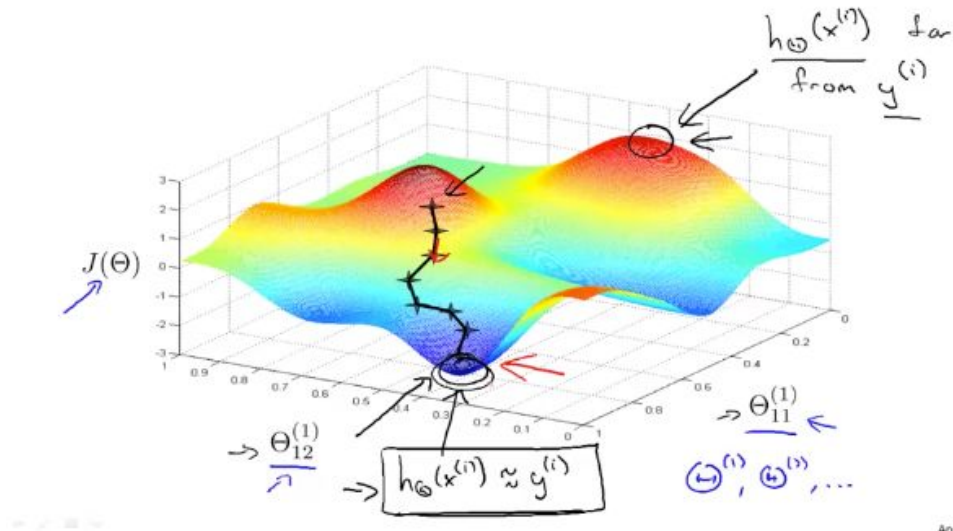  - Outside the for loop we compute the partial derivatives for $J(\theta)$

- The next 2 steps are
  - Use gradient checking to make sure back prop is correct. We do this by comparing the partial derivatives of $J(\theta)$ vs the numerical estimate of gradient of $J(\theta)$
    - We have to disable the gradient checking code because it is slow
  - Use gradient descent (or another advanced optimization method) to minimize $J(\theta)$

**Training a neural network**

→5. Use gradient checking to compare $\frac{\partial}{\partial\Theta^{(l)}_{ik}}J(\Theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.
   →Then disable gradient checking code.
→6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters $\Theta$

$\frac{\partial}{\partial\Theta^{(l)}_{jk}} J(\Theta)$

$J(\Theta)$  —  non-convex.

   Finally, $J(\theta)$ can end up being non convex and this, as you might remember from previous classes, might cause gradient descent not to converge on the global minimum. In practice this is not a big problem
- As a refresher, gradient descent with backpropagation is trying to start from one point and then it is computing the direction or path so we can end up at a point where the difference between the hypothesis $h_\theta(x^{(i)}) \approx y^{(i)}$. In the example below that would be the point circled in black. In contrast, a point atop of one of the red mountains is one where $h_\theta(x^{(i)})$ is far different from $y^{(i)}$

$$h_\Theta(x^{(i)}) \text{ far from } y^{(i)}$$

$J(\Theta)$

$$\Theta^{(1)}_{12}$$

$$h_\Theta(x^{(i)}) \approx y^{(i)}$$

$$\Theta^{(1)}_{11}$$

$$\Theta^{(1)}, \Theta^{(2)}, \ldots$$

# Week 6: Neural Networks Learning

*Advice for Applying Machine Learning*

- Now we are going to learn some practical advice on how to apply the learning algorithms we have learned so far.
- Suppose we have a learning algorithm like regularized linear regression and we apply it a new set of data and we find out there are very large errors on predictions, what can we do then? There are a few choices (shown below)

**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?

- Get more training examples
- Try smaller sets of features    $x_1, x_2, x_3, \ldots, x_{100}$
- Try getting additional features
- Try adding polynomial features $(x_1^2, x_2^2, x_1x_2, \text{etc.})$
- Try decreasing $\lambda$
- Try increasing $\lambda$

The problem is that usually people do not have a method to pick from this menu and go by gut feeling. To make it even worse, some of the choices above could be very long projects and we might end up wasting months and not get any better predictions. So how do we choose from the choices above? Turns out there is a set of techniques that can help us eliminate at least half of the choices; this method is….

- These techniques are known as Machine Learning Diagnostic. A diagnostic is a test that we can run to understand what is/is not working with an algorithm and we can gain guidance on how to improve its performance.

## 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.

*Evaluating a Hypothesis*

- One way to evaluate a hypothesis is to plot it against the training set as we have been doing. However, that only works for problems where the number of features is small. It is difficult to plot if we have hundreds of features

### Evaluating your hypothesis



Fails to generalize to new examples not in training set.

$$x_1 = \text{size of house}$$
$$x_2 = \text{no. of bedrooms}$$
$$x_3 = \text{no. of floors}$$
$$x_4 = \text{age of house}$$
$$x_5 = \text{average income in neighborhood}$$
$$x_6 = \text{kitchen size}$$
$$\vdots$$
$$x_{100}$$

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

Andrew Ng

- The standard way of evaluating a hypothesis is to take the data we have and assign a certain percentage as a training set (say 70%) and the remaining 30% is our test set.
  $m_{test} = number\ of\ test\ examples$ .

## 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 — 70%

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

Test set — 30%

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

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

Andrew N

It is also desirable to randomize the order of the training set so we can use a random 70% as the training set and a random 30% as the test set. In the example above we just used the first 70% but it's better if we make sure it is a random 70% set

- The procedure then for linear regression is to:
  - Learn θ from training data
  - Compute test set error which is expressed as $J(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m} (h_\theta(x^{(i)}_{test}) - y^{(i)}_{test})^2$

## Training/testing procedure for linear regression

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

- Compute test set error:

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

- The same process for logistic regression is similar: we learn θ from training set and then compute the error of the test set using the logistic regression function but expressed with $m_{test}$ as follows $J_{test}(\theta) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y^{(i)}_{test} log(h_\theta(x^{(i)}_{test})) - (1 - y^{(i)}_{test}) log(h_\theta(x^{(i)}_{test}))$ .

In this case sometimes we can use another test metric called the misclassification error (also known as the 0/1 misclassification error) and is defined below. Basically it is showing the number of examples on my test set that the hypothesis got wrong

## Training/testing procedure for logistic regression

- Learn parameter $\theta$ from training data $\quad M_{test}$
- 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):

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

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

Finally, to get the average test error we do a summation of the individual errors scaled by 1/mtest as shown above

*Model Selection and Train/Validation/Test Sets*

- We start by generalizing an overfitting example. We have seen overfitting many times before but the more general principle is that once my parameters $\theta$ fit some set of data (the training set), the error of these parameters as measured to that data is likely to be lower than the actual generalization error.

### 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.

In other words, just because the error is low for the training set it does not mean it is going to be the same when applied to other data. In fact, it will probably be higher for the general case

- A related question is how do we select the model to use. That is, what hypothesis should we choose? One with a quadratic equation? One with a 5th order polynomial? One with a 10th order polynomial? What we do here is we introduce a variable,

*d = degree of polynomial* and then we calculate the cost function $J_{test}(\theta)$ for each model choice and choose the one with the lowest cost. Let's assume that it is the 5th order polynomial. Here the question is still the same: how well does this generalize? The answer is, as with the previous slide, not very well and is likely to be an optimistic estimate of the generalization error. Why? For the same reason as before: we choose the *d* that best fits the test set.

**Model selection**

→ $d = $ degree of polynomial ↓

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

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

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

⋮

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

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

How well does the model generalize? Report test set $\boxed{\Theta_0, \Theta_1 \ldots}$

error $\underline{J_{test}(\theta^{(5)})}$. $\Theta^{(5)}$

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.

- What we do to get over this problem is to divide our data set in 3 ways: 60% is the training set, 20% is the cross validation (cv) test and the remaining 20% is the usual test set. We then have $m_{cv}$ = *number of cv examples* and an example from this data set is $(x_{cv}^{(i)}, y_{xv}^{(i)})$

**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 $(x_{cv}^{(i)}, y_{cv}^{(i)})$

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

$M_{test}$

- The cost function for each of the 3 sections is shown below. It's really the same equation but just uses the training set from each section

## Train/validation/test error

Training error:

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

Cross Validation error:

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

Test error:

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

- Our model selection process changes: we still minimize $J(\theta)$ for each of our model selections and that gives us a $\theta$ vector for each choice. We then take that parameter vector and apply it to the cross validation set to get the error for each possible choice. Once we are done calculating the error we pick the $h_\theta(x)$ that has the lowest error. In the example below we pick the 4th-order polynomial and use it to estimate the generalization error for the test set (we are now using the 3rd section of our data set). With this method our $d$ parameter has not been trained using the test set, it was trained using the cross validation set

### Model selection

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

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

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

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

$\vdots \qquad\qquad \vdots$

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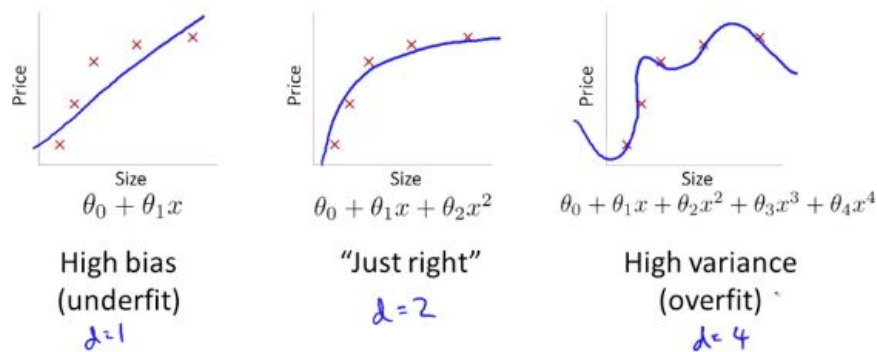
$d = 4$

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

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

*Diagnosing Bias vs Variance*

- To refresh our memory, high bias is underfitting when we have a hypothesis that is too simple (with a low order polynomial) and high variance is the opposite case: overfitting, where we have hypothesis that fits the data set too well (with a somewhat high order polynomial) but does not generalize. Both cases are shown below

## Bias/variance



$\theta_0 + \theta_1 x$

High bias
(underfit)
$d=1$

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

"Just right"
$d=2$

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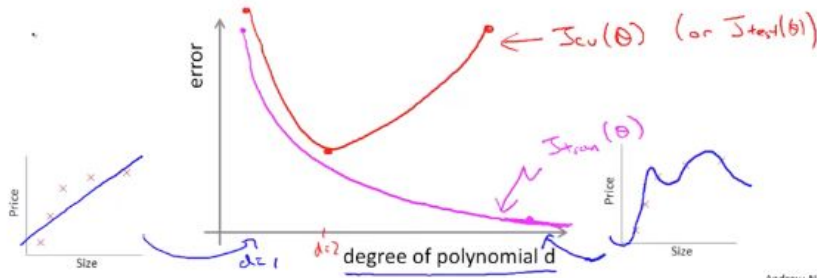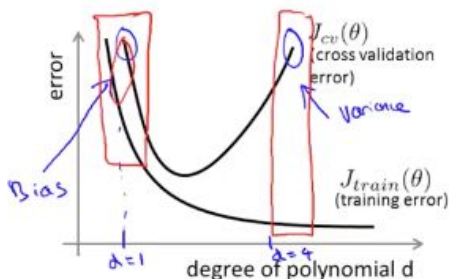
High variance
(overfit)
$d=4$

- We can get a graphical understanding of bias and variance by taking our $J_{train}(\theta)$ and our $J_{cv}(\theta)$ cost functions and graphing them against the degree of polynomial $d$. We end up with a chart as below. The $J_{train}(\theta)$ is like that because our error is high when we have a simple hypothesis (underfitting) but it is close to 0 (or perhaps 0) because a higher order polynomial allows to fit (and overfit) our hypothesis better with the training set.

### Bias/variance

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

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



$\leftarrow J_{cv}(\theta) \quad (or\ J_{test}(\theta))$

$J_{train}(\theta)$

degree of polynomial d

The cross validation error is a bit different because it starts similar (high error going down as the polynomial order increases) and then it reaches a point where the error starts to increase again (the point where overfitting starts)
- This graph then is a tool to help us identify if we have a bias or variance problem. If we have a bias problem then $J_{train}(\theta)$ is going to be high and $J_{train}(\theta) \approx J_{cv}(\theta)$. If we have a variance problem $J_{train}(\theta)$ is going to be very low and and $J_{cv}(\theta) \gg J_{train}(\theta)$

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

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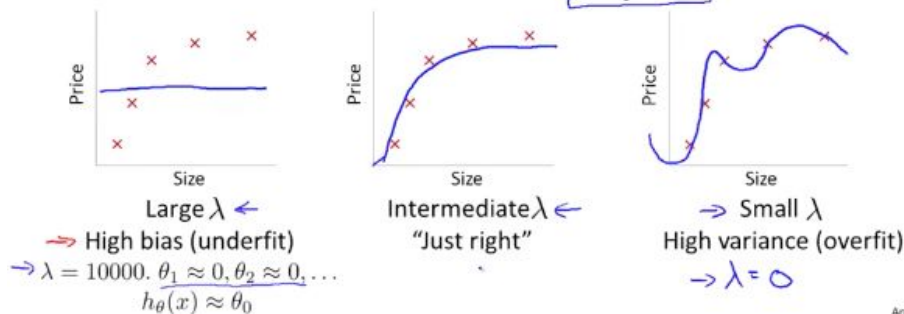$$J_{cv}(\theta) \gg J_{train}(\theta)$$

*Andrew*

## *Regularization & Bias/Variance*

- Now we take a look to how regularization impacts bias and variance. Considering the equation below: if we have a big $\lambda$ value then all our $\theta's$ are 0 and we end up with $h_\theta(x) = \theta_0$. If we have a small $\lambda$ value (say $\lambda = 0$) then we end up with original $h_\theta(x)$ which is overfitting. An intermediate $\lambda$ value is the one that gets it 'just right'

### Linear regression with regularization

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

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



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

Intermediate $\lambda$ ←
"Just right"

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

*Andrew*

- How can we then automatically select a good $\lambda$ value? We start by changing our cost function definition as shown below. Our original cost function includes the regularization parameter. In this case, $J_{train}(\theta)$ and $J_{cv}(\theta)$ and $J_{test}(\theta)$ will NOT have a regularization param, they will just be the sum of the square of the errors * 1/(2m)

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

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

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

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

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

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

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

- Once we have our cost functions for all 3 sections of the data set what we do is to create a list of possible $\lambda$ values. Professor recommends that you start with a small value (say 0.01) and double it until you get to around 10. After having the list what we do is we try to minimize the original cost function (the one that has the regularization parameter) using all the $\lambda$ values from the list and we end up with a list of possible parameters $\theta$. We take each member of the list and apply it to our $J_{cv}(\theta)$ equation and we then select the cost function that minimizes the cross validation error( $\theta^{(5)}$ in the example below). We take $\theta^{(5)}$ and apply it to $J_{test}(\theta)$ so we end up with $J_{test}(\theta^{(5)})$

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

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$
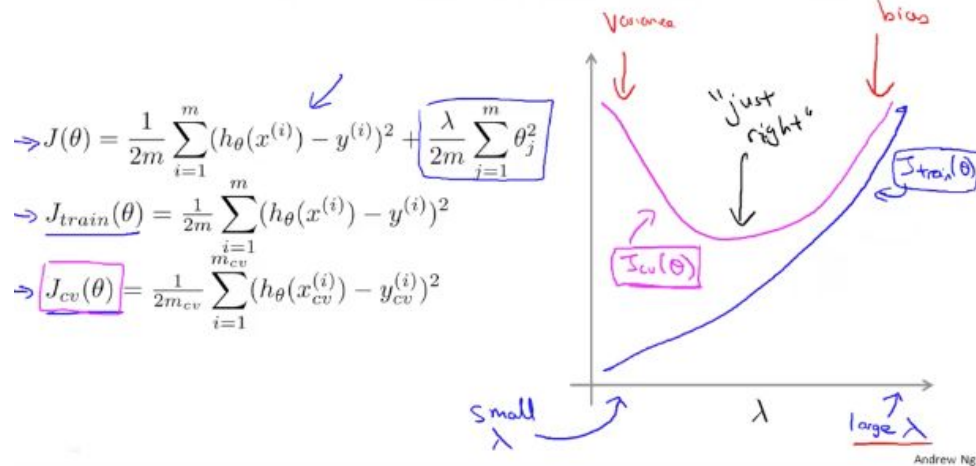
1. Try $\lambda = 0 \leftarrow$ $\longrightarrow$ $\min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. Try $\lambda = 0.01$ $\longrightarrow$ $\min J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. Try $\lambda = 0.02$ $\longrightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
4. Try $\lambda = 0.04$
5. Try $\lambda = 0.08$ $\longrightarrow \theta^{(5)} \quad J_{cv}(\theta^{(5)})$

   $\vdots$

12. Try $\lambda = 10$ $\longrightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
    $\uparrow$ 10.24
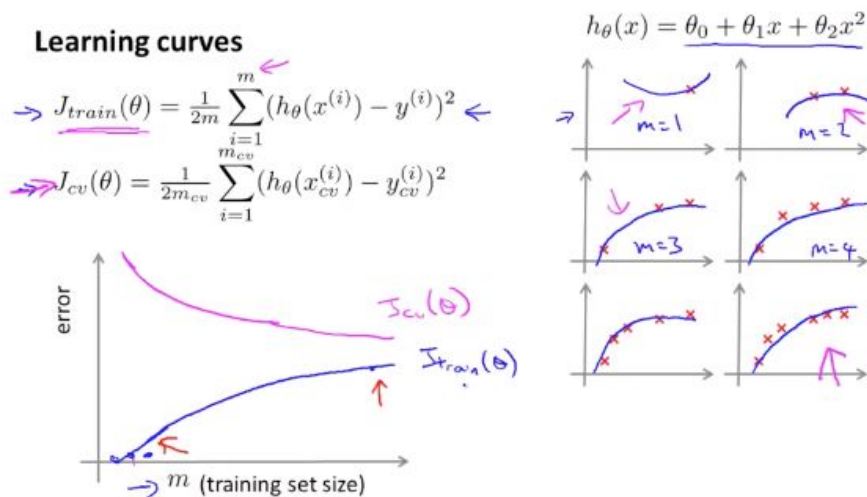    Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$

- We can then try to graph $J_{train}(\theta)$ and $J_{cv}(\theta)$ against $\lambda$ and we end up with a chart as shown below. If $\lambda$ is big it will penalize all my $\theta's$ and I end up with $h_\theta(x) = \theta_0$ which is then again a simple hypothesis with high bias so my chart below has the bias regime on the right. Conversely, if $\lambda$ is low then I end up with the original high order polynomial and am therefore in the high variance (overfit) area. Just as in the previous chart there is a portion where the $\lambda$ value is 'just right'

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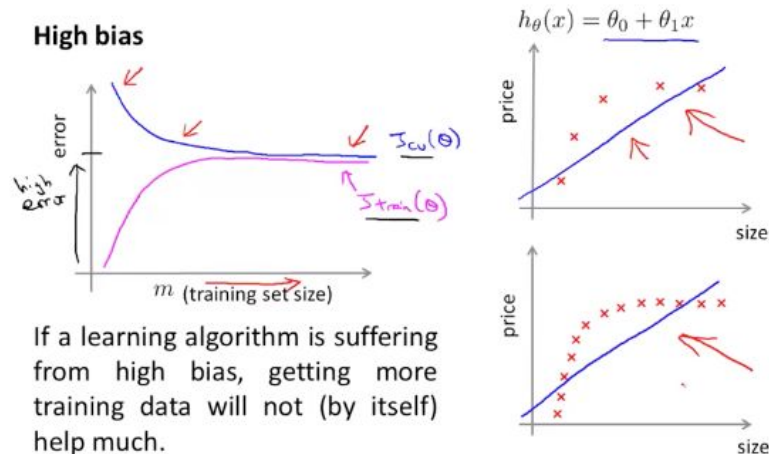


$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$

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

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

*Andrew Ng*

*Learning Curves*

- Learning curves are useful to check if the algorithm is working correctly or if you want to improve its performance. The learning curves for $J_{train}(\theta)$ and $J_{cv}(\theta)$ are shown below. The intuition is that if we have a small set of training examples (say 1-3) we will be able to fit $h_\theta(x)$ to the data so our error is 0 or very close to 0. However, as m starts to increase it gets harder to fit the dataset and our error increases.

### Learning curves



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

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

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

*Andrew N*

Cross validation is different. We trained on a small number of examples so we don't expect our generalization to be very good therefore our error will be high. However, as we get more training examples we expect our cross validation error to decrease.

- Now let's look at how the learning curves change when we have high bias or high variance. Let's start with high bias. Suppose we have a hypothesis like the one below which is a line. What we see in this case is that we start with a small amount of data but

as we increase the amount of data we end up with kind of the same line. This means that initially I will have a small error for $J_{train}(\theta)$ and the error increases and kind of plateaus as we get a bigger m. For $J_{cv}(\theta)$ is kind of the opposite: it starts with a high error because we have started training with a small $m$ and as the $m$ gets bigger we get kind of the best fit for the data and our error decrease also plateaus. As m gets bigger $J_{train}(\theta)$ and $J_{cv}(\theta)$ kind of converge. What this also tells us is if you are underfitting (high bias) your data you will not get any significant improvement from getting from training data (potentially saving us days/months of work of data collection that will not work)



**High bias**

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

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

- Now let's take a look at high variance (overfitting). If we are overfitting we have a high order polynomial. In this case, graphing our error against $m$ will show us that
  - $J_{train}(\theta)$ will be initially low and will remain low as the training size increases. Makes sense since we are overfitting for this training set and the error will then be low.
  - $J_{cv}(\theta)$ will be initially high which also makes sense because overfitting for one training set does not generalize to other data sets. As we get bigger training sizes the cross validation error goes down.



**High variance**

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

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

In a case like this getting more training data will help because it will help us to see if $J_{cv}(\theta)$ continues to decrease. Finally, we can identify if the algorithm suffers from high variance if we see a relatively big gap between the $J_{cv}(\theta)$ and $J_{train}(\theta)$

*Deciding what to do next Revisited*

- Now we can go back to original problem and provide cases where the below options are going to be helpful
    - Getting more training examples is going to fix *high variance*
    - Trying a smaller set of features is going to fix *high variance*
    - Getting additional features is going to fix *high bias*
    - Adding more polynomial features is going to fix *high bias*
    - Decreasing $\lambda$ is going to fix *high bias*
    - Increasing $\lambda$ is going to fix *high variance*

**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

- Now let's apply what we learned to neural networks. A small neural network is prone to underfitting but is computationally cheap. A large neural network is prone to overfitting and computational expensive. We can use regularization to address overfitting

**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.

- We are going to try to build a spam classifier. We note some characteristics of spam emails (words that have numbers instead of letters). If an email is classified as spam we have a positive id (1), if non-spam we have a negative id (0).

**Building a spam classifier**

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

Deal of the week! Buy now!
Rolex w4tchs - $100
Medlcine (any kind) - $50
Also low cost M0rgages
available.

Spam (1)

From: Alfred Ng
To: ang@cs.stanford.edu
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans
for Xmas. When do you get off
work. Meet Dec 22?
Alf

Non-spam (0)

- Now we need to create the X vector of features. In this case X will be 100 words indicative of spam or not spam. We order them alphabetically and add a 1 or 0 in our X vector if the word appears on an email message. We are then saying that $x_j = 1$ *if word j appears in email*, $x_j = 0$ *otherwise*. In practice, we don't chose 100 words, we take the most frequently occurring $n$ words (10,000 to 50,000)

**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, discont, andrew, now, ...

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

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

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.

- Then we try to answer the question of what is the best use of my time if I want to have a classifier with low error? Some options are
  - Collect lots of data, i.e. honeypot projects
  - Develop sophisticated features based on email header routing info
  - Develop sophisticated features for for message body
  - Develop sophisticated algorithms to detect misspellings

**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.)

*Building a Spam Classifier: Error Analysis*

- When faced with a new machine learning problem we can start with the approach detailed below: start with a simple algo, implement it and test it on the cv data. Plot learning curves to decide what to do next and, finally, examine the examples where the algo made errors to see if we can spot a trend in what examples we are making errors on. In this class we are going to concentrate in Error Analysis, the last step.

**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.

- Here's an example of how we use Error Analysis: in a cross validation set of 500 examples, 100 were misclassified. We manually examine the errors and categorize them by type of email and see if we have any clues on what features could have helped in classifying those samples correctly. The manual classification allows us to understand the type of email that is giving us trouble along with the feature(s) that is also fooling the algorithm. In the example below is 'steal passwords' and 'unusual punctuation'. This is why it pays to use a simple algorithm: with it we can quickly see what type of emails and/or features are giving us trouble.

**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

- Next, we will be faced with some decisions where error analysis will not help. In such cases, it is important to get our error result as a single, numerical value. Otherwise, it is hard to assess the algo's performance. We should try new things, get their error rate and based on that result decide if we want to use 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 Metrics for Skewed Classes*

- Skewed classes are problems where we have way more examples from one class than from the other class as shown below with the cancer classification example. In this case we had an algorithm that achieved 99% correctness/accuracy but a non-learning algorithm can do better: 0.5% error. It then becomes harder to evaluate if we did a useful change to the algorithm or we just replaced it with a non-learning one? It is then harder to just use classification accuracy because we can get high classification accuracy with an algorithm that does not look like such a good choice. We then need a different error/evaluation metric. One such metric is called Precision/Recall.

**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.

$\longrightarrow$ skewed classes.

$\rightarrow$ 0.5% error

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

→ 99.2% accuracy (0.8% error)

→ 99.5% accuracy (0.5% error)

- To define precision recall we start with an example where $y = 1$ (i.e. a tumor is malignant). In this case we can compute a table like the one below using the following cases:
    - When we predicted y = 1 or y =0
    - When the actual y = 1 or 0

If we predicted y=1 and y was actually =1 we have a true positive.
If we predicted y=0 and y was actually =0 we have a true negative.
If we predicted y=1 but y was actually = 0 we have a false positive.
If we predicted y=0 but y was actually = 1 we have a false negative.

**Precision/Recall**

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



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

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

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

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

We now can compute
- $Precision = True\ Positive/(True\ Positive + False\ Positive)$ and
- $Recall = True\ Positive / (True\ Positive + False\ Negative)$

Finally, if we have an algorithm where y =0 all the time then $Recall = 0$ because it will not have any True Positives. Now we can see Precision/Recall is helping us to see if an algorithm is 'cheating' and gives us high confidence that we have a good algorithm.

*Trading off Precision and Recall*

- There is a trade off to be made between precision and recall. To illustrate it let's continue with our cancer example. We start with a prediction threshold of 0.5. First we consider that we want to predict y=1 only if we are very confident, we then change our threshold to something higher, say 0.7. In this case we have a classifier with higher precision but lower recall. Now let's consider the opposite case: we want to avoid missing too many cases of cancer so we lower the threshold to say 0.3. In this case we decrease our precision but increase our recall.



A couple of different precision curves is shown on the chart above. We generalize this trade off as predict 1 if $h_\theta(x) \geq threshold$.

- Next question is how do we compare precision/recall numbers? The problem with precision/recall is that now we have 2 numbers to compare whereas before we insisted in the importance of having a single number to evaluate. The solution is to use the $F_1 score$ defined as $2\frac{PR}{P+R}$. This score gives us a good way to evaluate our algorithm and takes care of cases where we have a not very good classifier (i.e. one that always predict y=1 or y=0). It is tempting to use the average but the $F_1 score$ is better

- How much data should we train on? Let's start with an influential study by Banko and Brill that showed that the 4 algorithms shown below benefited greatly from having more data available to them. Other studies have shown the same result. This fact has led to the saying 'It's not who has the best algorithm that wins, it's who has the most data'.



- Now let's see the assumptions under which a massive set of data will help. The rationale says that feature x has sufficient information to predict y accurately. An example supporting this rationale is the 'fill in the blank" example below. A counterexample is the housing price prediction with only size. Professor Ng uses the following way to test this rationale : given the feature x, would a human expert confidently predict y?



- Let's see when having a lot of data helps. We have a learning algorithm with many parameters -> ensures low bias which means $J_{train}(\theta)$ will be small. Now assume we use a huge training set. With such a large training set (unlikely to overfit) we will have low variance and, therefore, $J_{train}(\theta) \approx J_{test}(\theta)$ and $J_{test}(\theta)$ will be small

**Large data rationale**

$\Rightarrow$ 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.

$$\rightarrow J_{train}(\Theta) \text{ will be small.}$$

Use a very large training set (unlikely to overfit)

$$\rightarrow J_{train}(\Theta) \approx J_{test}(\Theta)$$

$$\rightarrow J_{test}(\Theta) \text{ will be small}$$