

# Machine Learning Part II

Neural Network  
ML Diagnostic

# Motivation

- When we have too many features, and we need a non-linear hypothesis, the linear (polynomial) regression and logistic regression is not a good algorithm.

# Model Representation

- The function maps inputs to output is called **activation function**.
  - also uses sigmoid function
- The parameters are also called **weights**.
- Input layer, hidden layers, output layer

$a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

$\Theta^{(j)}$  = matrix of weights controlling function mapping from layer  $j$  to layer  $j + 1$

If we had one hidden layer, it would look like:

$$[x_0 x_1 x_2 x_3] \rightarrow [a_1^{(2)} a_2^{(2)} a_3^{(2)}] \rightarrow h_{\theta}(x)$$

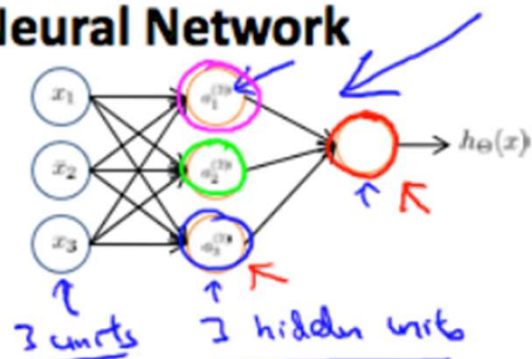
The values for each of the "activation" nodes is obtained as follows:

$$\begin{aligned}a_1^{(2)} &= g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3) \\a_2^{(2)} &= g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3) \\a_3^{(2)} &= g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3) \\h_{\theta}(x) = a_1^{(3)} &= g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})\end{aligned}$$

This is saying that we compute our activation nodes by using a 3×4 matrix of parameters. We apply each row of the parameters to our inputs to obtain the value for one activation node. Our hypothesis output is the logistic function applied to the sum of the values of our activation nodes, which have been multiplied by yet another parameter matrix  $\Theta^{(2)}$  containing the weights for our second layer of nodes.

Each layer gets its own matrix of weights,  $\Theta^{(j)}$ .

## Neural Network



$\rightarrow a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

$\rightarrow \Theta^{(j)}$  = matrix of weights controlling function mapping from layer  $j$  to layer  $j + 1$

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

$$h_{\Theta}(x)$$

$$\rightarrow a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$\rightarrow a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$\rightarrow a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$\rightarrow h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

$$\Theta^{(2)}$$



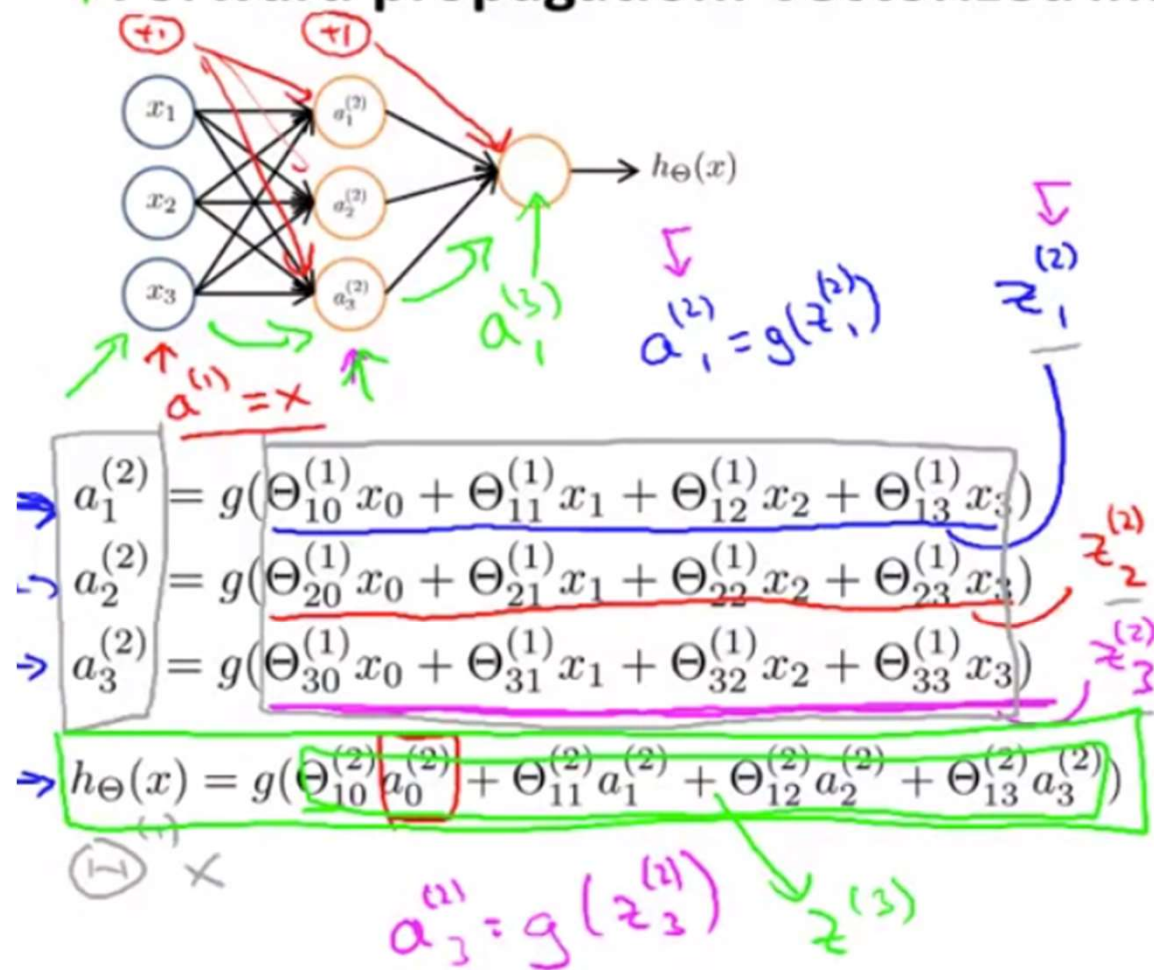
$\rightarrow$  If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j + 1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

$$s_{j+1} \times (s_j + 1)$$

Andrew N

Example: If layer 1 has 2 input nodes and layer 2 has 4 activation nodes. Dimension of  $\Theta^{(1)}$  is going to be  $4 \times 3$  where  $s_j = 2$  and  $s_{j+1} = 4$ , so  $s_{j+1} \times (s_j + 1) = 4 \times 3$ .

## Forward propagation: Vectorized implementation



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

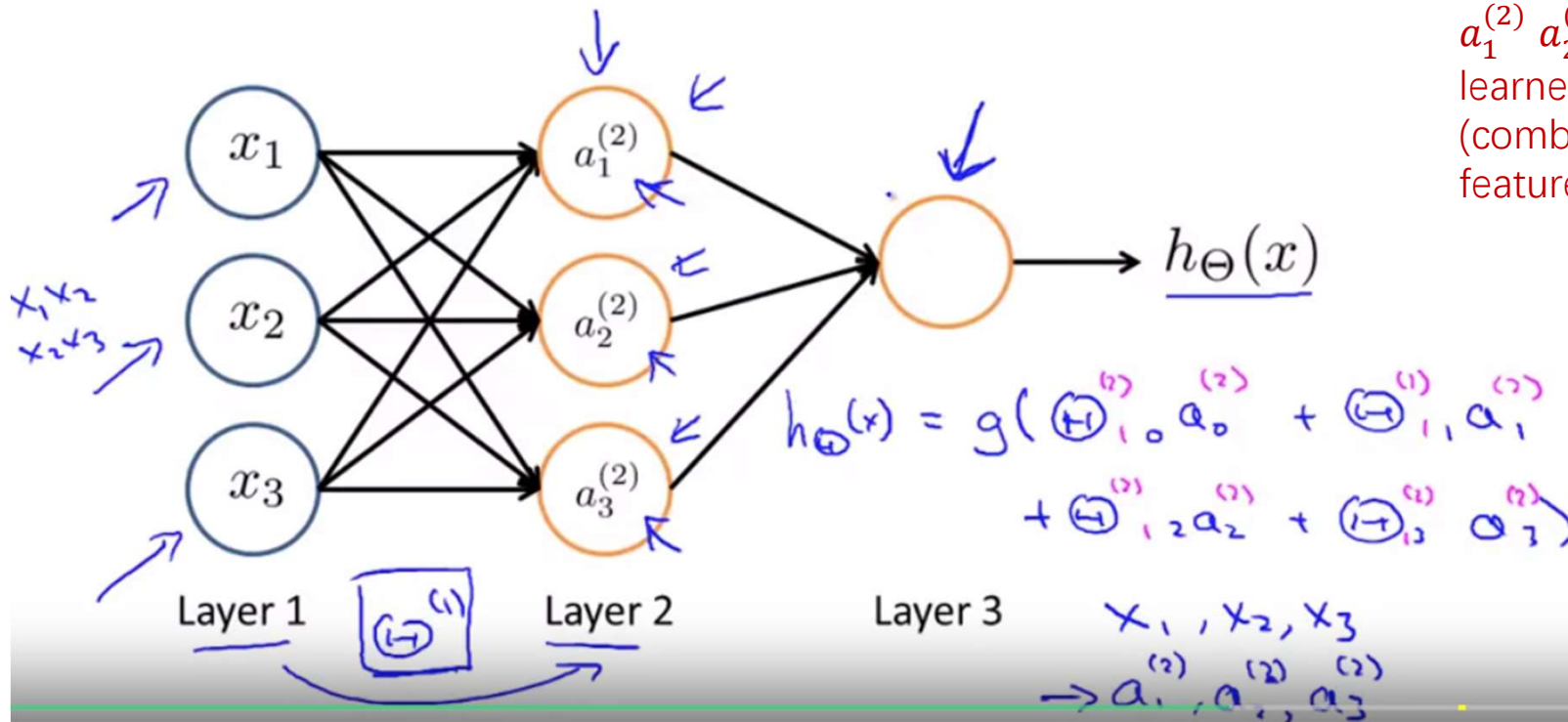
Add  $a_0^{(2)} = 1$   $\rightarrow a^{(2)} \in \mathbb{R}^4$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$$

- Neural Networks is actually a multi-units multi-layer logistic regression algorithm

### Neural Network learning its own features

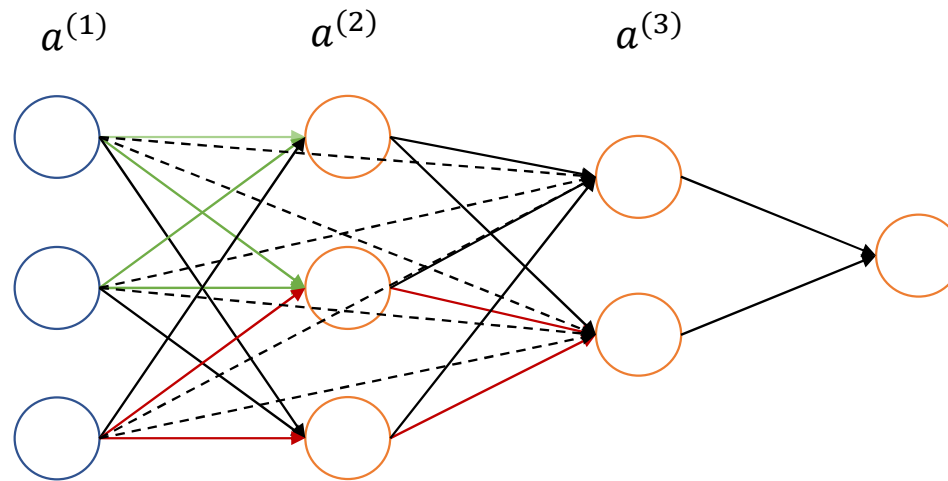


$a_1^{(2)} a_2^{(2)} a_3^{(2)}$  is the learned features (combination of input features)

- $a_1^{(2)} a_2^{(2)} a_3^{(2)}$  and  $h_{\Theta}(x)$  form a logistic regression algorithm unit;
- similarly,  $x_1 x_2 x_3$  and  $a_1^{(2)}$  form a a logistic regression algorithm unit



- Q:  $a^{(1)}$  connects to  $a^{(2)}$ ,  $a^{(2)}$  connects to  $a^{(3)}$ , can we connect  $a^{(1)}$  to  $a^{(3)}$  ?



Patterns in neurons: bi-fan, bi-parallel and **feedforward loop**

Green links form the bi-fan motif

Red links form the bi-parallel motif

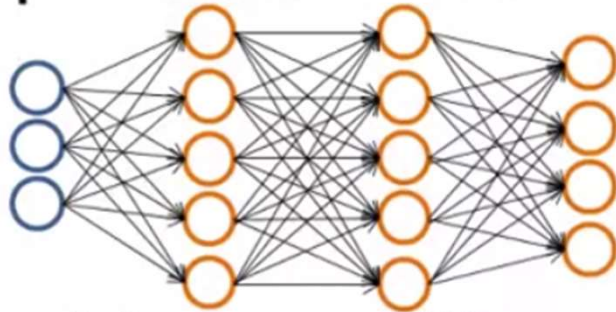
As feedforward loop is another motif in neurons, is it possible that we can make the dotted link between, say the input layer  $a^{(1)}$  and the 2<sup>nd</sup> hidden layer  $a^{(3)}$ . etc..



# Multiclass Classification

- one-hot vector for each class

## Multiple output units: One-vs-all.



Want  $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ , etc.  
when pedestrian      when car      when motorcycle

Training set:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

$\rightarrow y^{(i)}$  one of  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$   
pedestrian   car   motorcycle   truck

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

~~Previously~~

~~$y \in \{1, 2, 3, 4\}$~~

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

# Cost Function

- $L$  = total number of layers in the network
- $s_l$  = number of units (not counting bias unit) in layer  $l$
- $K$  = number of output units/classes

Recall that in neural networks, we may have many output nodes. We denote  $h_{\Theta}(x)_k$  as being a hypothesis that results in the  $k^{th}$  output. Our cost function for neural networks is going to be a generalization of the one we used for logistic regression. Recall that the cost function for regularized logistic regression was:

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

For neural networks, it is going to be slightly more complicated:

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

Double sum adds up the logistic regression costs calculated for each cell in the output layer

Triple sum adds up the squares of all  $\Theta$ s in the entire network.

# Backpropagation

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ).

(used to compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ )

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

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

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

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

→ Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$   ~~$\delta^{(1)}$~~

→  $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  ←

$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$ .

→  $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$  if  $j \neq 0$

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

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

# Implementation: Unrolling parameters

```
→ 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);
```

Share

## Learning Algorithm

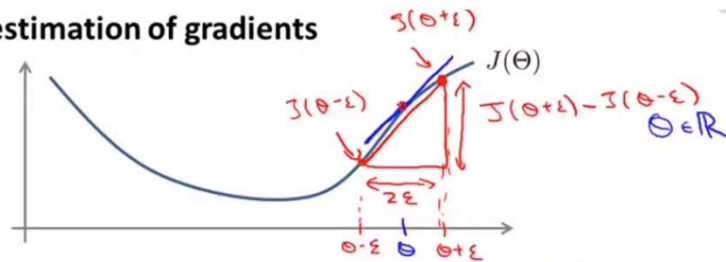
- Have initial parameters  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ .
- Unroll to get `initialTheta` to pass to
- `fminunc(@costFunction, initialTheta, options)`

```
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

Numerical estimation of gradients



$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$\epsilon = 10^{-4}$

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

```
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;
```

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

Check that `gradApprox`  $\approx$  `DVec`

From backprop.



# Random Initialization

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

E.g.

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

*Handwritten notes:*  
A blue arrow points from the text "Random 10x11 matrix (betw. 0 and 1)" to the `rand(10,11)` function call in the code above.

→ `Theta2 = rand(1,11) * (2*INIT_EPSILON) - INIT_EPSILON;`

NOT initialize theta by same value

# Machine Learning Diagnostic



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

## The test set error

1. For linear regression:  $J_{test}(\Theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\Theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$
2. For classification ~ Misclassification error (aka 0/1 misclassification error):

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

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

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

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

# Train/ Validation/ Test Sets

- Using **validation** set to choose hyperparameters (**Model selection**), such as the degree of polynomial, the number of hidden units, etc
  - 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);

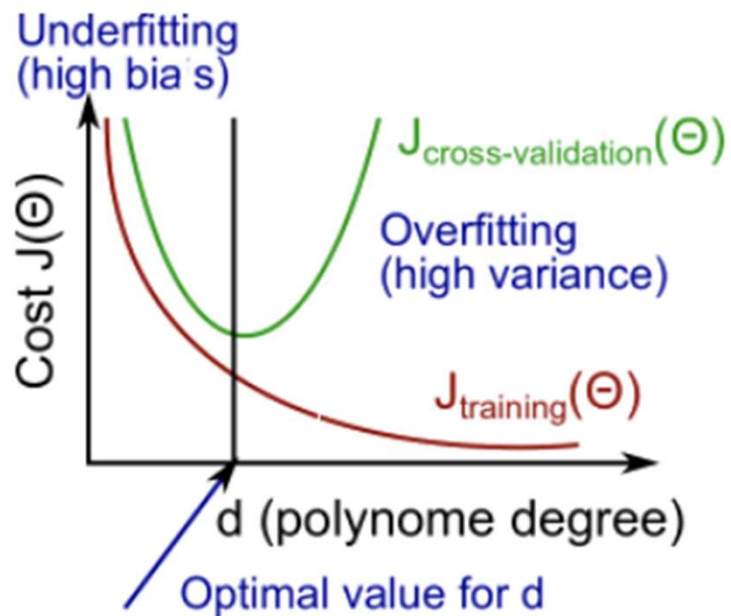
This way, the degree of the polynomial  $d$  has not been trained using the test set.

# Bias vs. Variance

**High bias (underfitting):** both  $J_{train}(\Theta)$  and  $J_{CV}(\Theta)$  will be high. Also,  $J_{CV}(\Theta) \approx J_{train}(\Theta)$ .

**High variance (overfitting):**  $J_{train}(\Theta)$  will be low and  $J_{CV}(\Theta)$  will be much greater than  $J_{train}(\Theta)$ .

The is summarized in the figure below:



# Choosing regularization parameter

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  $\lambda = 0$   $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
  2. Try  $\lambda = 0.01$   $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
  3. Try  $\lambda = 0.02$   $\rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
  4. Try  $\lambda = 0.04$
  5. Try  $\lambda = 0.08$   $\rightarrow \theta^{(5)}$
  - $\vdots$
  12. Try  $\lambda = 10$   $\rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
- $\uparrow$  10.24 Pick (say)  $\theta^{(5)}$ . Test error:  $J_{test}(\theta^{(5)})$

# Learning Curve (error – training set size)

Experiencing high variance:

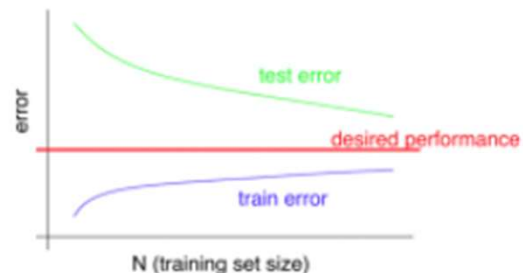
**Low training set size:**  $J_{train}(\Theta)$  will be low and  $J_{CV}(\Theta)$  will be high.

**Large training set size:**  $J_{train}(\Theta)$  increases with training set size and  $J_{CV}(\Theta)$  continues to decrease without leveling off. Also,  $J_{train}(\Theta) < J_{CV}(\Theta)$  but the difference between them remains significant.

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

## More on Bias vs. Variance

Typical **learning curve** for high variance(at fixed model complexity):



- **Getting more training examples:** Fixes high variance
- **Trying smaller sets of features:** Fixes high variance
- **Adding features:** Fixes high bias
- **Adding polynomial features:** Fixes high bias
- **Decreasing  $\lambda$ :** Fixes high bias
- **Increasing  $\lambda$ :** Fixes high variance.

## 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.

# Ex5 Notes

- Optional 3.4

Note that set **lambda = 0** when test learned model.

because lambda is introduced to deal with overfitting and train parameters of the model.

- Optional 3.5 (should be all example from validation set)

```
for i = 1:m
```

```
    for t = 1:50
```

```
        k = randperm(length(X));
```

```
        [theta] = trainLinearReg(X(k(1:i), :), y(k(1:i)), lambda);
```

```
        error_train(i) = error_train(i) + linearRegCostFunction(X(k(1:i), :), y(k(1:i)), theta, 0);
```

```
        error_val(i) = error_val(i) + linearRegCostFunction(Xval, yval, theta, 0);
```

```
    end
```

```
    error_train(i) = error_train(i) / 50;
```

```
    error_val(i) = error_val(i) / 50;
```

```
end
```



# Machine Learning Design

The recommended approach to solving machine learning problems is to:

- Start with a simple algorithm, implement it quickly, and test it early on your cross validation data.
- Plot learning curves to decide if more data, more features, etc. are likely to help.
- Manually examine the errors on examples in the cross validation set and try to spot a trend where most of the errors were made.

Also, use a **single numerical value** to evaluate the performance of algorithm, so we can quickly try out different strategies.

|                 |   | Actual class      |                   |
|-----------------|---|-------------------|-------------------|
|                 |   | 1                 | 0                 |
| Predicted class | 1 | True<br>Positive  | False<br>Positive |
|                 | 0 | False<br>Negative | True<br>Negative  |

$$F_1 \text{ Score: } 2 \frac{PR}{P+R}$$

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

$$\text{Recall} = \frac{\text{True positives}}{\# \text{ actual positives}} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$