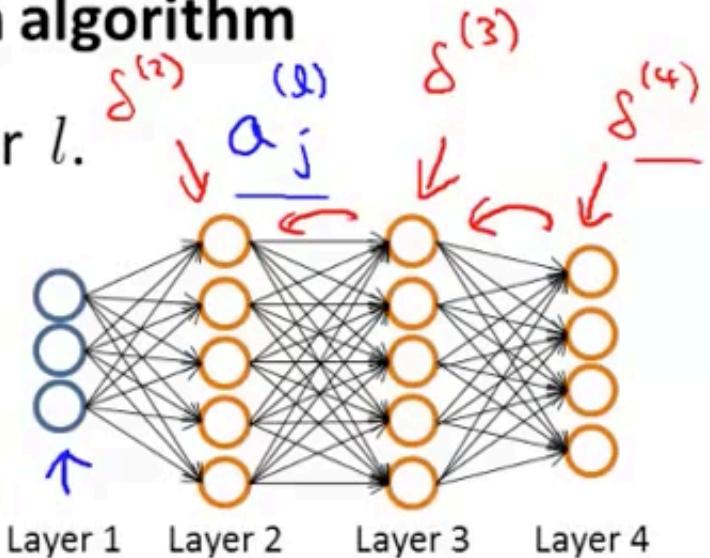


## Gradient computation: Backpropagation algorithm

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



For each output unit (layer  $L = 4$ )

$$\delta_j^{(4)} = \underline{a_j^{(4)}} - \underline{y_j} \quad (\underline{h_{\Theta}(x)})_j \quad \underline{\delta^{(4)}} = \underline{a^{(4)}} - \underline{y}$$

$$\rightarrow \delta^{(3)} = (\underline{\Theta^{(3)}})^T \underline{\delta^{(4)}} \cdot * g'(z^{(3)})$$

$$\rightarrow \delta^{(2)} = (\underline{\Theta^{(2)}})^T \underline{\delta^{(3)}} \cdot * g'(z^{(2)})$$

$$(No \quad \underline{\delta^{(1)}}) \quad \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

$$\frac{a^{(3)} - * (1 - a^{(3)})}{a^{(2)} - * (1 - a^{(2)})}$$

(ignoring  $\lambda$ ; if  
 $\lambda = 0$ ) ↵

## Backpropagation algorithm

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

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

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

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

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

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

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

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

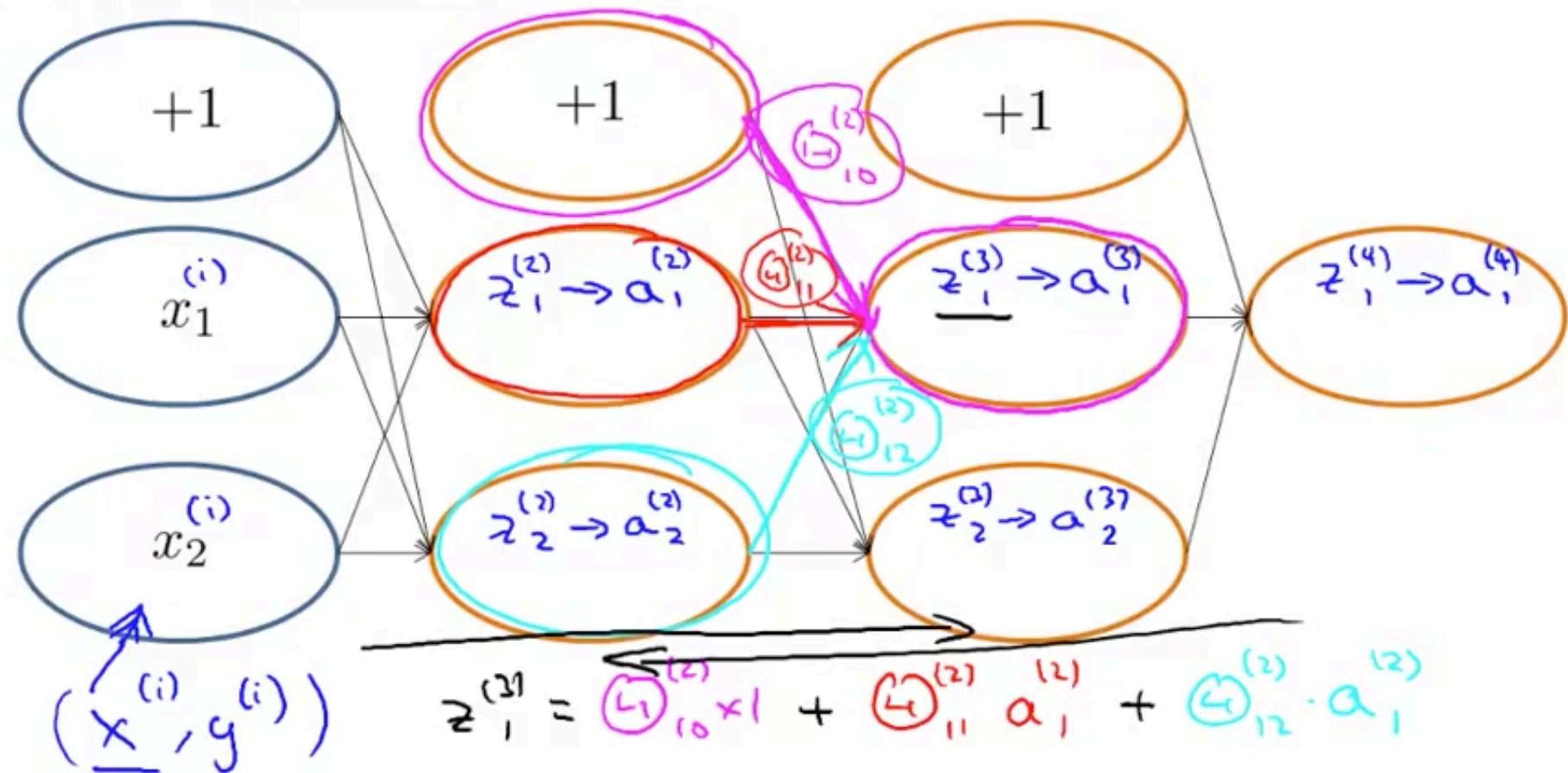
$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \delta^{(l+1)} (\underline{a}^{(l)})^T.$$

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

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

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

## Forward Propagation



## What is backpropagation doing?

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

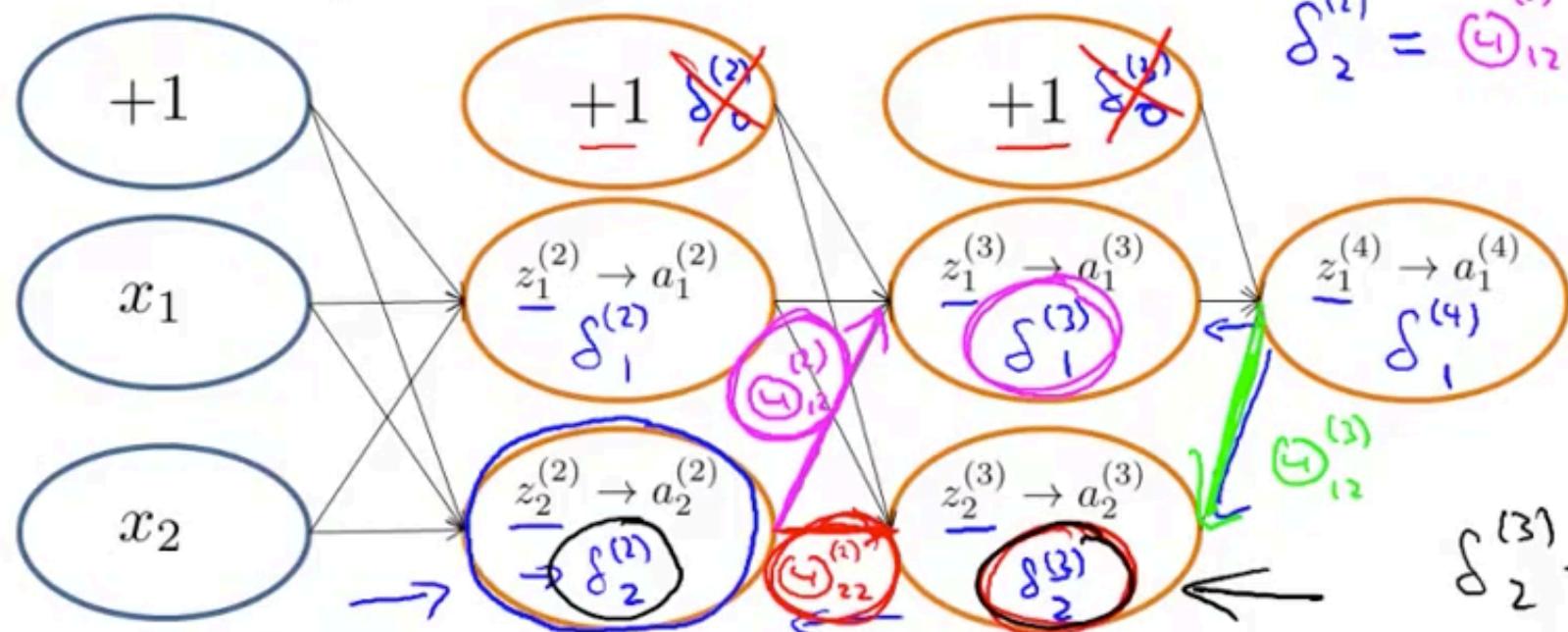
Focusing on a single example  $x^{(i)}, y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda = 0$ ),

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

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

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

## Forward Propagation



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

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

$$\delta_1^{(4)} = y^{(i)} - a_1^{(4)}$$

$$\delta_2^{(2)} = (w_{12}^{(1)}) \delta_1^{(3)} + (w_{22}^{(1)}) \cdot \delta_2^{(1)}$$

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

## Advanced optimization

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

Annotations:

- A red arrow points from the gradient parameter to the text  $\mathbb{R}^{n+1}$ .
- A red arrow points from the theta parameter to the text  $\mathbb{R}^{n+1}$  (vectors).
- A red arrow points from the options parameter to the Neural Network ( $L=4$ ) section.

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

## Example

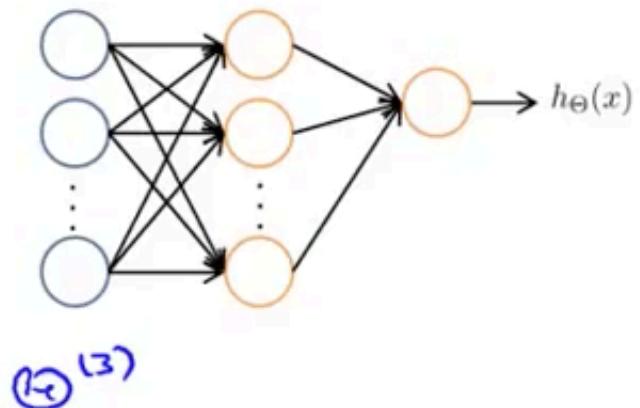
$$s_1 = 10, s_2 = 10, s_3 = 1$$

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

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

$\rightarrow$  thetaVec = [ Theta1(:); Theta2(:); Theta3(:) ];  
 $\rightarrow$  DVec = [ D1(:); D2(:); D3(:) ];

$\rightarrow$  Theta1 = reshape(thetaVec(1:110), 10, 11);  
 $\rightarrow$  Theta2 = reshape(thetaVec(111:220), 10, 11);  
 $\hookrightarrow$  Theta3 = reshape(thetaVec(221:231), 1, 11);



Octave-3.2.4

```
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1
```

```
>> Theta2 = 2*ones(10,11)
```

```
Theta2 =
```

```
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2
```

```
>> Theta3 = 3*ones(1,11)
```

Octave-3.2.4

```
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2
```

```
>> Theta3 = 3*ones(1,11)
```

```
Theta3 =
```

```
3 3 3 3 3 3 3 3 3 3 3
```

```
>> thetaVec = [ Theta1(:); Theta2(:); Theta3(:) ];
```

```
>> size(thetaVec)
```

```
ans =
```

```
231 1
```

```
>>
```



5:44 PM  
10/16/2011

Octave-3.2.4

```
1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1
```

```
>> reshape(thetaVec(111:220), 10,11)
```

```
ans =
```

```
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2
```

```
>> reshape(thetaVec(221:231), 1,11)
```

```
ans =
```

```
3 3 3 3 3 3 3 3 3 3
```

```
>>
```



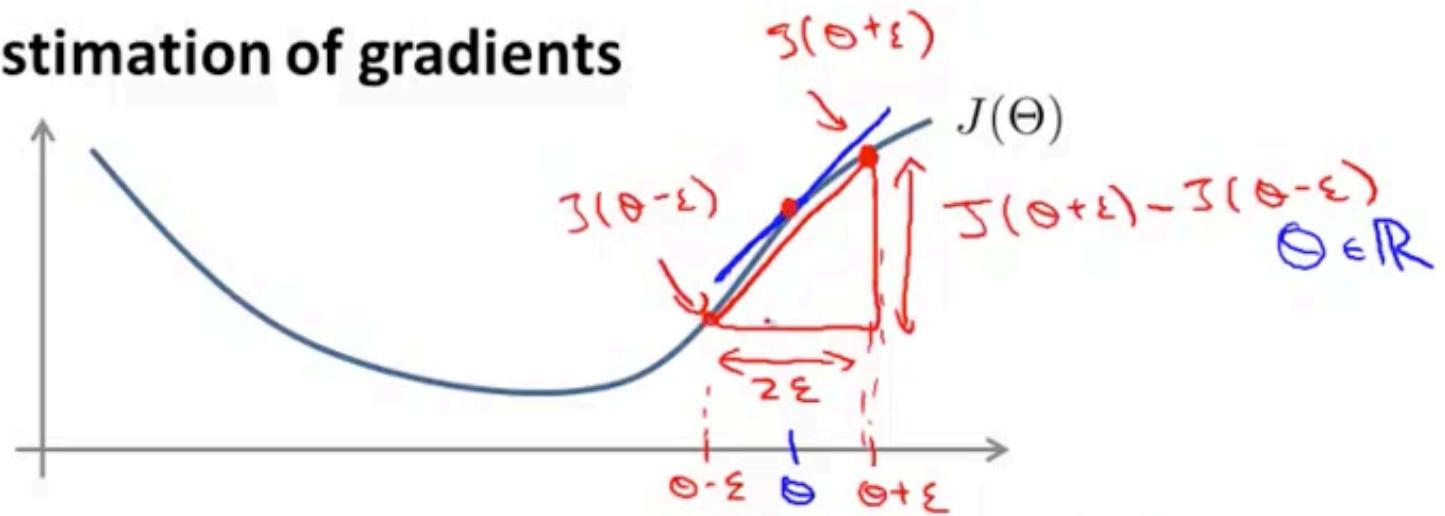
5:45 PM  
10/16/2011

## 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.
```

## Numerical estimation of gradients



$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \leftarrow \quad \cancel{\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}}$$

$\epsilon = 10^{-4} \quad \leftarrow$

Implement: gradApprox =  $(J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON})) / (2 * \text{EPSILON})$

## Parameter vector $\theta$

- $\theta \in \mathbb{R}^n$  (E.g.  $\theta$  is “unrolled” version of  $\underline{\Theta^{(1)}}, \underline{\Theta^{(2)}}, \underline{\Theta^{(3)}}$ )
- $\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$
- $\frac{\partial}{\partial \theta_1} \underline{J(\theta)} \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$
- $\frac{\partial}{\partial \theta_2} \underline{J(\theta)} \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$
- ⋮
- $\frac{\partial}{\partial \theta_n} \underline{J(\theta)} \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$

```

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;

```

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

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

Check that gradApprox  $\approx$  DVec  $\leftarrow$   
 From backprop.

## 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.

↓ ↓ ↓

gradApprox.

→ 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.

## Initial value of $\Theta$

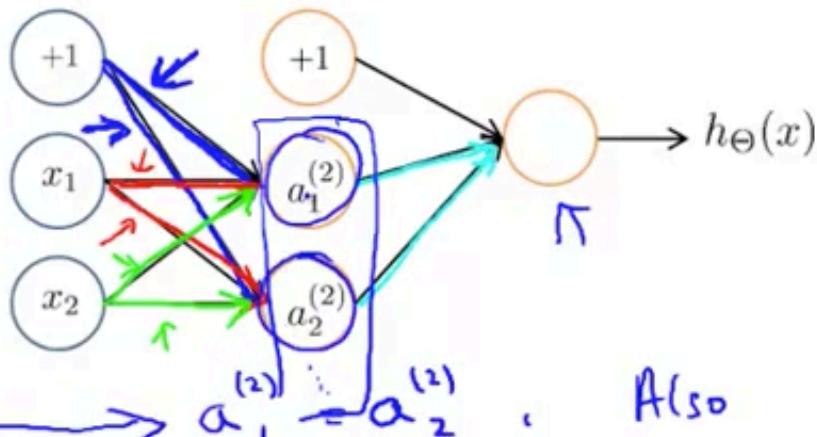
For gradient descent and advanced optimization method, need initial value for  $\Theta$ .

```
optTheta = fminunc(@costFunction,  
                    initialTheta, options)
```

Consider gradient descent

Set initialTheta = zeros(n,1) ?

## Zero initialization



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

$$\text{Also } \delta_1^{(2)} = \delta_2^{(2)}.$$

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

$$\underline{\Theta_{01}^{(1)}} = \underline{\Theta_{02}^{(1)}}$$

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

$$\underline{\underline{\Theta_{01}^{(1)}}} = \underline{\underline{\Theta_{02}^{(1)}}}$$

## Random initialization: Symmetry breaking

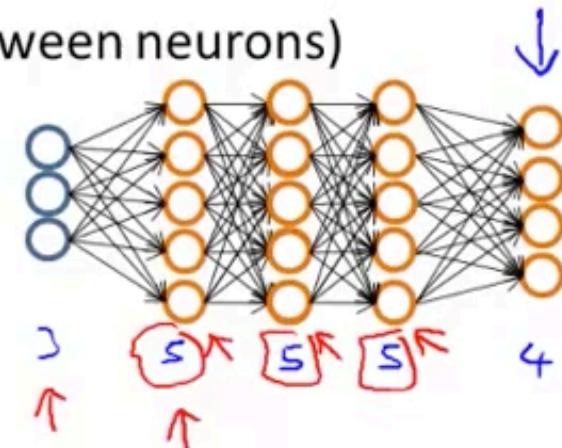
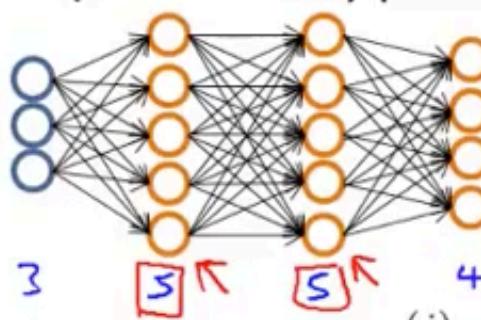
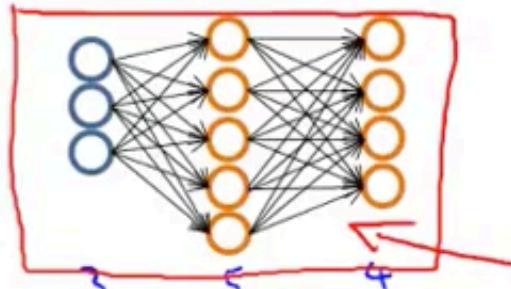
- Initialize each  $\Theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$   
(i.e.  $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$ )

E.g.

- $\Theta_{\text{theta1}} = \boxed{\text{rand}(10, 11) * (2 * \text{INIT\_EPSILON})}$   
-  $\text{INIT\_EPSILON}$ ;  $[-\epsilon, \epsilon]$
- $\Theta_{\text{theta2}} = \boxed{\text{rand}(1, 11) * (2 * \text{INIT\_EPSILON})}$   
-  $\text{INIT\_EPSILON}$ ;
- ↑  
Random 10x11 matrix (betw. 0  
and 1)

## Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features  $x^{(i)}$

→ No. output units: Number of classes

[Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

$$y \in \{1, 2, 3, \dots, 10\}$$

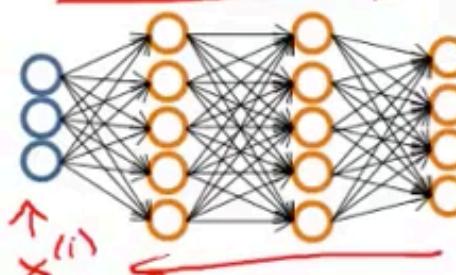
~~$y = 5$~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \leftarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1 \\ 0 \end{bmatrix}$$

Andrew Ng

## Training a neural network

- 1. Randomly initialize weights
- 2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
- 3. Implement code to compute cost function  $J(\Theta)$
- 4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
- **for**  $i = 1:m$  {  $(x^{(1)}, y^{(1)})$     $(x^{(2)}, y^{(2)})$  , ...  $(x^{(m)}, y^{(m)})$  }
  - Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$
  - (Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ .)
  - $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (\alpha^{(l)})^T$
  - ...  
}
  - - -  
compute  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ .



## Training a neural network

- 5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$  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  $\Theta$

$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta) \quad \uparrow$$

$J(\Theta)$  — non-convex.

