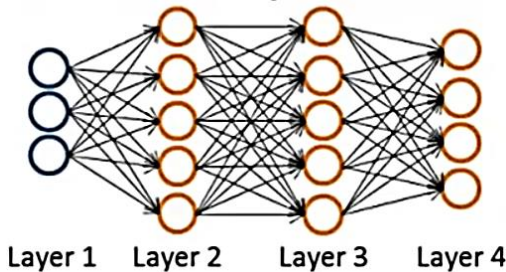


# 9. Neural Networks - Learning

⇒ **COST FUNCTION:**

## Neural Network (Classification)



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

→  $L =$  total no. of layers in network  $\underline{L = 4}$

→  $s_l =$  no. of units (not counting bias unit) in layer  $l$   $\underline{s_1 = 3, s_2 = 4, s_3 = 4, s_4 = 4}$

### Binary classification

$y = 0$  or  $1$  ←

1 output unit ←

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

$$s_L = 1, \quad \underline{K = 1}$$



### Multi-class classification (K classes)

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

### K output units

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

$$s_L = K \quad (K \geq 3)$$

Cost function in Neural Networks is just a generalization of Logistic regression:

With regularization term included.

## Cost function

Logistic regression:

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

We denote  $h_{\Theta}(x)_k$  as being a hypothesis that results in the  $k^{\text{th}}$  output.

## For Neural Networks:

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

We have added a few nested summations to account for our multiple output nodes. In the first part of the equation, before the square brackets, we have an additional nested summation that loops through the number of output nodes.

In the regularization part, after the square brackets, we must account for multiple theta matrices. The number of columns in our current theta matrix is equal to the number of nodes in our current layer (including the bias unit). The number of rows in our current theta matrix is equal to the number of nodes in the next layer (excluding the bias unit). As before with logistic regression, we square every term.

Note:

- the double sum simply adds up the logistic regression costs calculated for each cell in the output layer
- the triple sum simply adds up the squares of all the individual  $\Theta$ s in the entire network.
- the  $i$  in the triple sum does **not** refer to training example  $i$

---

**BACKPROPOGATION ALGORITHM:** to calculate the gradient of cost function to minimize it

## Gradient computation

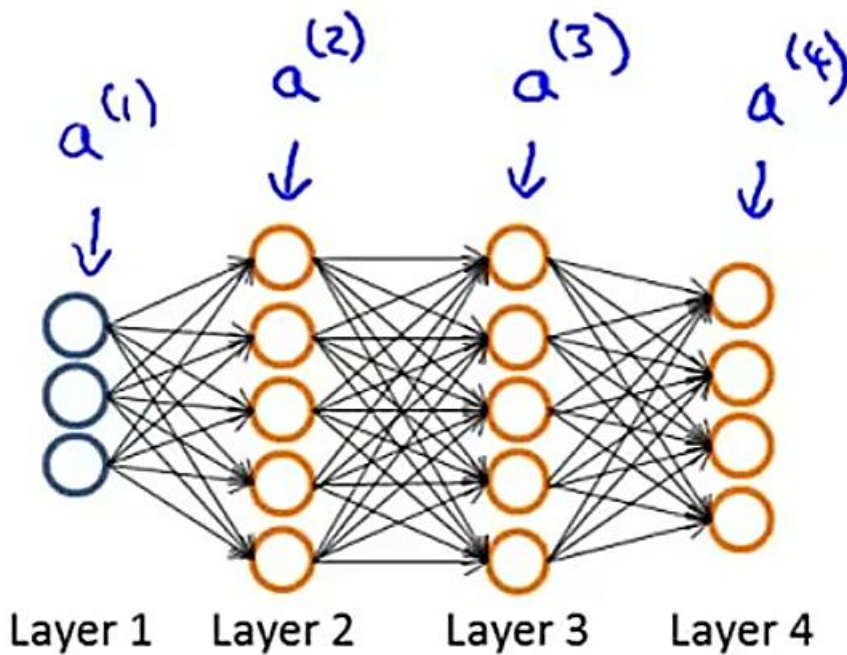
$$\min_{\Theta} J(\Theta)$$

Need code to compute:

$$\begin{aligned} &> - \frac{J(\Theta)}{\partial \Theta_{ij}^{(l)}} J(\Theta) \end{aligned}$$

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

## Computing gradient:



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

$\Rightarrow x, y$  are vectors

1. First, we do the forward propagation:

Forward propagation:

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

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

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

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

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

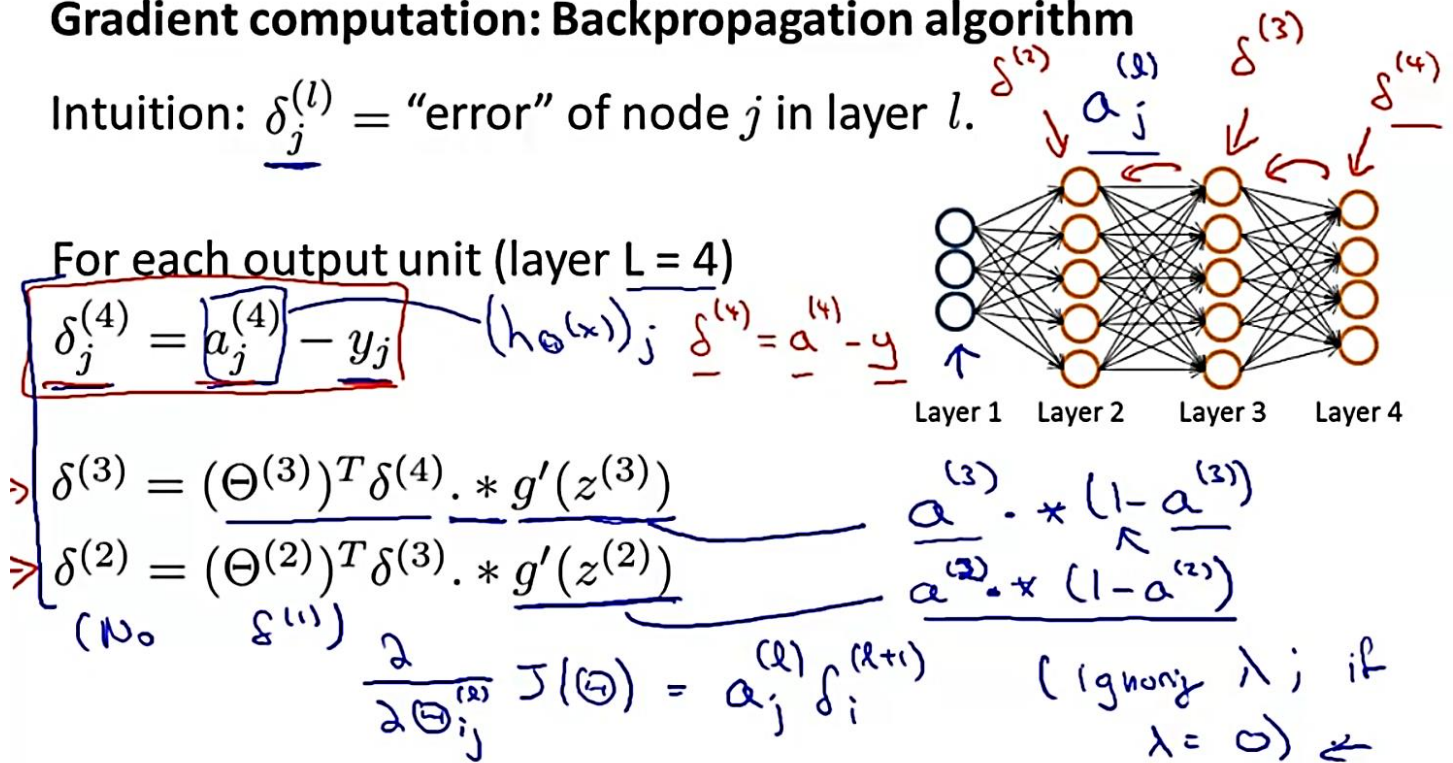
$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$

## 2. Next, we do back propagation:

### Gradient computation: Backpropagation algorithm

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



$j$  = a particular unit in layer  $L$

$g'$  =  $g$ -prime = derivative of  $g(z)$  wrt  $z$

⇒ we don't calculate  $\delta$  for 1<sup>st</sup> layer as it's the input layer and thus it has no errors.

- First, we calculate  $\delta$  for all units of o/p layer
- then for all other layers in backwards order
- we don't calculate  $\delta$  for 1<sup>st</sup> layer
- then using all  $\delta$ 's we calculate  $\Delta$  for all layers
- then using  $\Delta$  we calculate  $D$  for all layers .  $D$ =derivative of  $J(\Theta)$  wrt  $\Theta$  of layer  $l$ )



## Backpropagation algorithm

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

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

(use 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  $l = 2, 3, \dots, L$

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

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

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

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

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

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

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

$j=0 \Rightarrow$  corresponds to bias unit in layer  $l$

## Summary:

### Back propagation Algorithm

Given training set  $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$

- Set  $\Delta_{i,j}^{(l)} := 0$  for all  $(l,i,j)$ , (hence you end up having a matrix full of zeros)

For training example  $t = 1$  to  $m$ :

- Set  $a^{(1)} := x^{(t)}$
- Perform forward propagation to compute  $a^{(l)}$  for  $l=2,3,\dots,L$
- Using  $y^{(t)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(t)}$

"error values" for the last layer are simply the differences of our actual results in the last layer and the correct outputs in  $y$

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  using  $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) .* a^{(l)} .* (1 - a^{(l)})$

The delta values of layer  $l$  are calculated by multiplying the delta values in the next layer with the theta matrix of layer  $l$ . We then element-wise multiply that with a function called  $g'$ , or  $g$ -prime, which is the derivative of the activation function  $g$  evaluated with the input values given by  $z^{(l)}$ .

The  $g$ -prime derivative terms can also be written out as:

$$g'(z^{(l)}) = a^{(l)} .* (1 - a^{(l)})$$

5.  $\Delta_{i,j}^{(l)} := \Delta_{i,j}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  or with vectorization,  $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

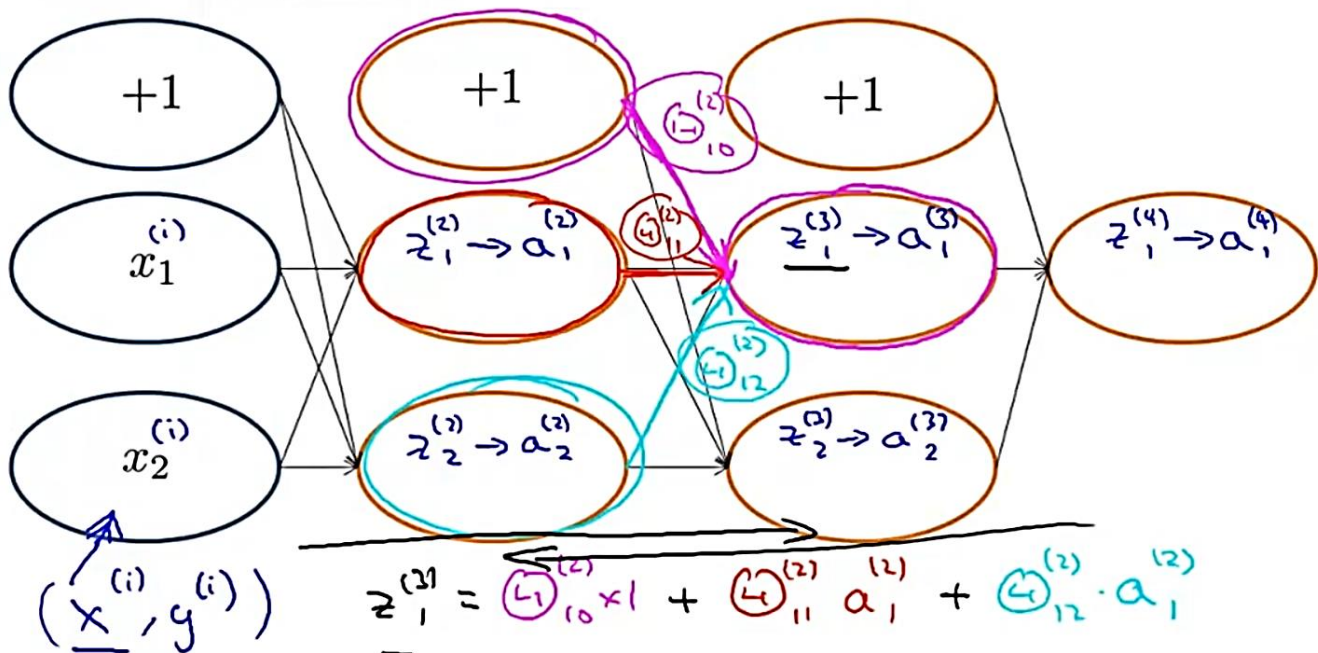
Hence we update our new  $\Delta$  matrix.

- $D_{i,j}^{(l)} := \frac{1}{m} (\Delta_{i,j}^{(l)} + \lambda \Theta_{i,j}^{(l)})$ , if  $j \neq 0$ .
- $D_{i,j}^{(l)} := \frac{1}{m} \Delta_{i,j}^{(l)}$  if  $j = 0$

The capital-delta matrix  $D$  is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative. Thus we get  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

## BACKPROPPOGATION: INTUTION:

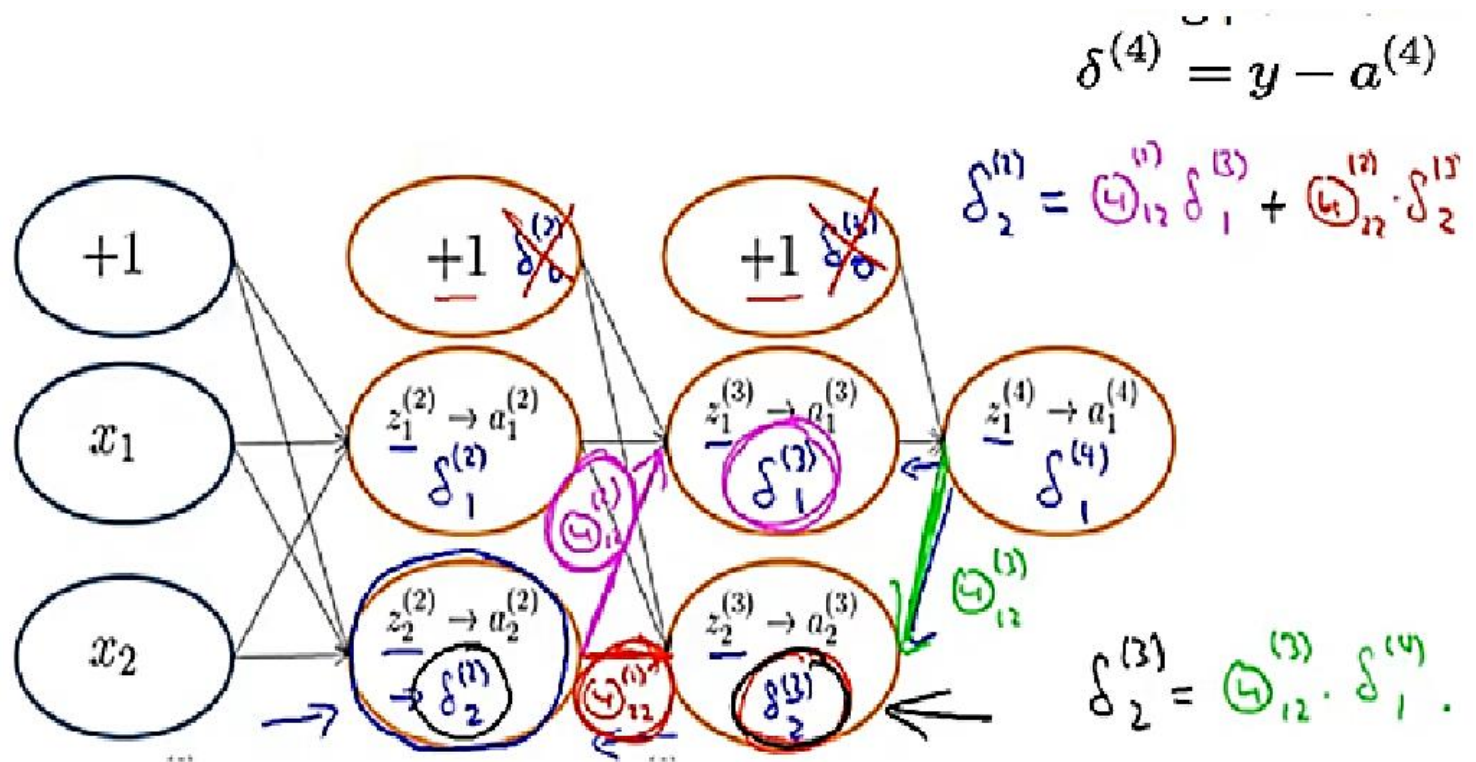
### Forward Propagation



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

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$  (for  $j \geq 0$ ), where

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



## MATRIX vs VECTORS DURING IMPLEMENTATION:

Matrices are useful when doing forward and backward propagation

Vectors are useful when using advanced optimization algo like `fminunc()`

`fminunc` assume the  $\Theta$  passed as argument is a vector and the gradient which the cost fn returns is also a vector

---



## Advanced optimization

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

Handwritten notes:  $\mathbb{R}^{n+1}$  (under gradient),  $\mathbb{R}^{n+1}$  (vectors) (under initialTheta)

But, original  $\Theta$  and gradient are matrices: so we need to unroll them into vectors

Neural Network ( $L=4$ ):

→  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  - matrices (Theta1, Theta2, Theta3)

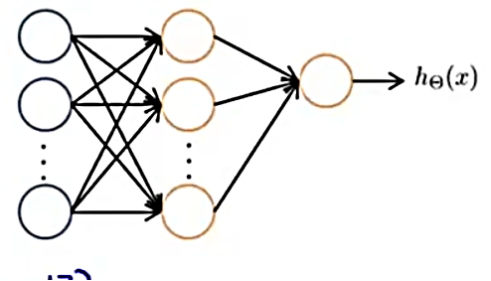
→  $D^{(1)}, D^{(2)}, D^{(3)}$  - matrices (D1, D2, D3)

"Unroll" into vectors

## Example: Binary Classification

### Example

$s_1 = 10, s_2 = 10, s_3 = 1$   
 $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$   
 $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$



To unroll into vectors:

Handwritten notes:  $\Theta^{(1)}$ ,  $\Theta^{(2)}$ ,  $\Theta^{(3)}$

```
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);
```



## 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)`

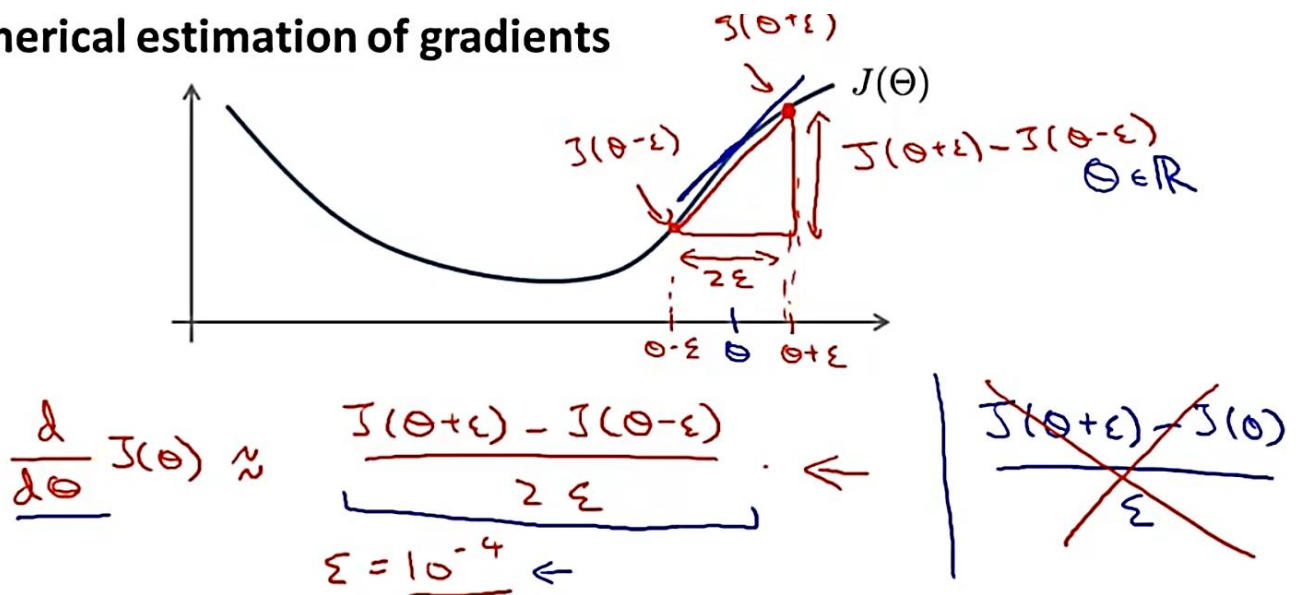
→ From `thetaVec`, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  *reshape*

→ 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



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

## Parameter vector $\theta$

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

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

$$\begin{aligned} \rightarrow \frac{\partial}{\partial \theta_1} J(\theta) &\approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon} \\ \rightarrow \frac{\partial}{\partial \theta_2} J(\theta) &\approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon} \\ &\vdots \\ \rightarrow \frac{\partial}{\partial \theta_n} J(\theta) &\approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon} \end{aligned}$$

## Octave code:

```
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;
```

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

$\frac{2}{2\theta_i} J(\theta)$

Check that  $\text{gradApprox} \approx \text{DVec}$  ←

$\uparrow$   
 From backprop.

## Implementation Note:

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

## Important:

$\begin{matrix} \downarrow & \downarrow & \downarrow \\ f^{(1)} & f^{(2)} & f^{(3)} \end{matrix} \rightarrow \text{DVec}$

## RANDOM INITIALIZATION:

### Initial value of $\Theta$

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

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

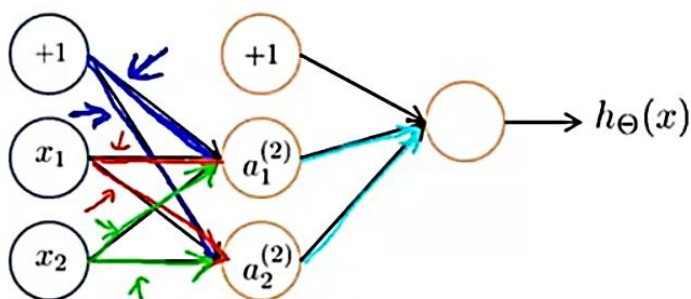
Consider gradient descent

Set initialTheta = zeros(n,1) ?

$\Theta = 0$  doesn't work in Neural Networks:

When we backpropagate, all nodes will update to the same value repeatedly. Instead we can randomly initialize our weights.

### Zero initialization



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

$$a_1^{(2)} = a_2^{(2)} \quad \text{Also} \quad \delta_1^{(2)} = \delta_2^{(2)}$$

$$\frac{\partial \Theta_{0,1}^{(1)}}{\partial \Theta_{0,1}^{(1)}} J(\Theta) = \frac{\partial \Theta_{0,1}^{(1)}}{\partial \Theta_{0,1}^{(1)}} J(\Theta)$$

$$\underline{\Theta_{0,1}^{(1)}} = \underline{\Theta_{0,2}^{(1)}}$$

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

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

## 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.

```
Theta1 = rand(10,11) * (2*INIT_EPSILON) - INIT_EPSILON;
```

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

This  $\epsilon$  is different from the one used in gradient checking.

Doing this will give a good variation in values of  $\Theta$  and the  $J(\Theta)$  will be best minimized.

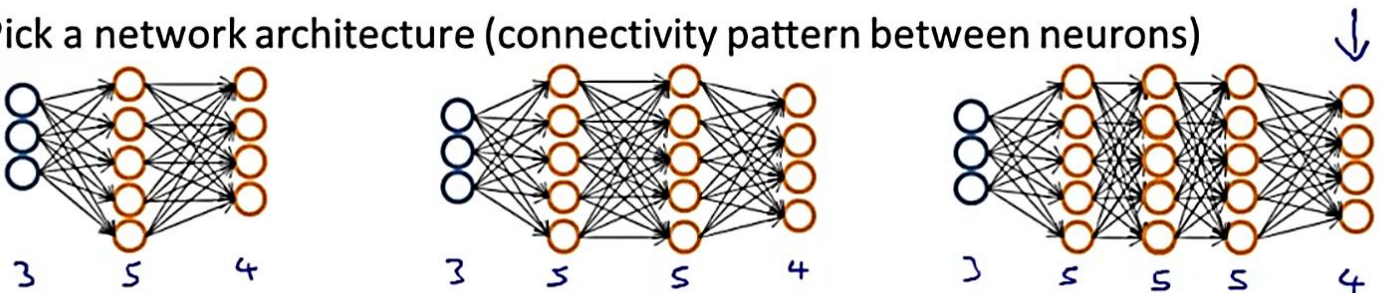
---

## PUTTING IT TOGETHER:

First, pick a network architecture; choose the layout of your neural network, including how many hidden units in each layer and how many layers in total you want to have.

### Training a neural network

Pick a network architecture (connectivity pattern between neurons)



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

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)



Number of hidden units per layer  $\rightarrow$  usually more the better (must balance with cost of computation as it increases with more hidden units)

Defaults: 1 hidden layer. If you have more than 1 hidden layer, then it is recommended that you have the same number of units in every hidden layer.

## Training a neural network

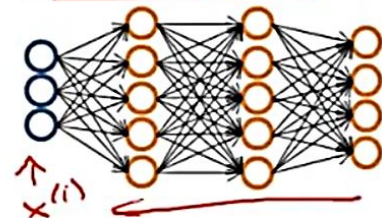
- $\rightarrow$  1. Randomly initialize weights
  - $\rightarrow$  2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
  - $\rightarrow$  3. Implement code to compute cost function  $J(\Theta)$
  - $\rightarrow$  4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
- $\rightarrow$  for  $i = 1:m$  {  $(x^{(1)}, y^{(1)})$   $(x^{(2)}, y^{(2)})$ , ...,  $(x^{(m)}, y^{(m)})$  }

$\rightarrow$  Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$

(Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ ).

$$\Delta^{(2)} := \Delta^{(2)} + \delta^{(n)} (a^{(2)})^T$$

compute  $\frac{\partial}{\partial \Theta_{jk}^{(2)}} J(\Theta)$ .



## Training a neural network

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

$$\frac{\partial}{\partial \Theta_{jk}^{(2)}} J(\Theta)$$

Ideally, you want  $h_{\Theta}(x^{(i)}) \approx y^{(i)}$ . This will minimize our cost function. However, keep in mind that  $J(\Theta)$  is not convex and thus we can end up in a local minimum instead.