

## Week 4

### — Neural Networks: Representation (Non-linear hypotheses)

When we have  $n$  features and  $n$  is relatively large, the map features we use to add more polynomial features makes so many more features so we can end up overfitting the hypothesis and it requires so much computational effort to minimize the function.

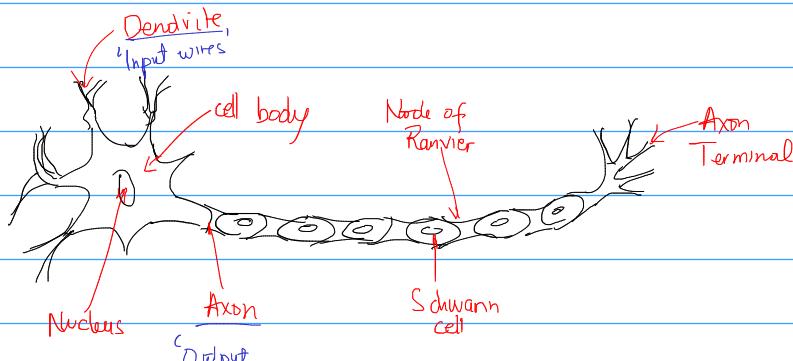
### — Neural Networks : Representation (Neurons and the brain)

Origins: Algorithms that try to mimic the brain

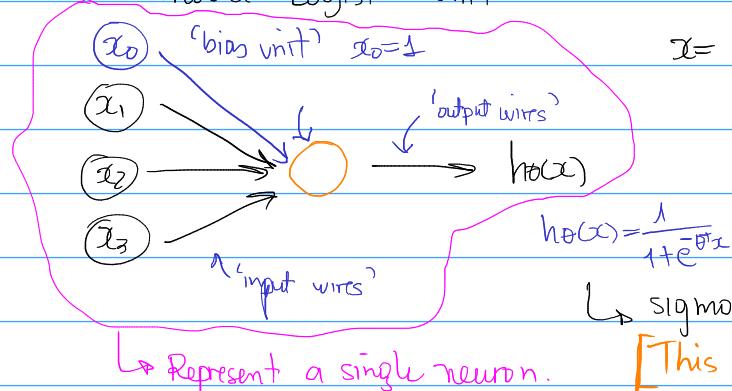
Was very widely used in 80s and early 90s; popularity diminished in late 90s

Recent resurgence: State of the art techniques for many applications.

### — Neural Networks : Representation (Model Representation I)



#### Neuron model: Logistic Unit

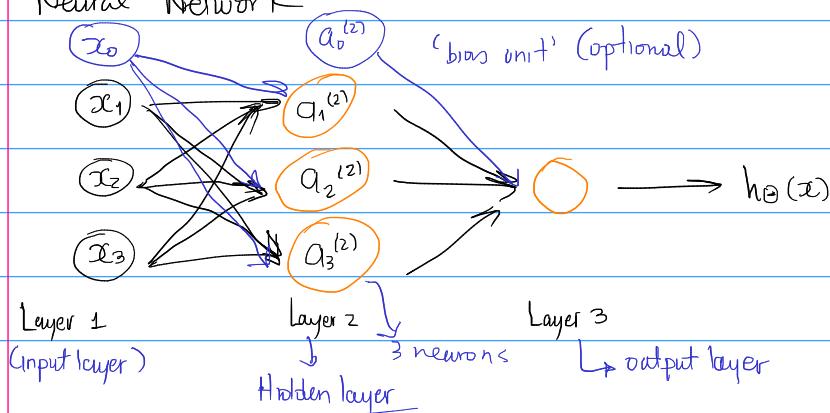


$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

↳ 'weights' = 'parameters'

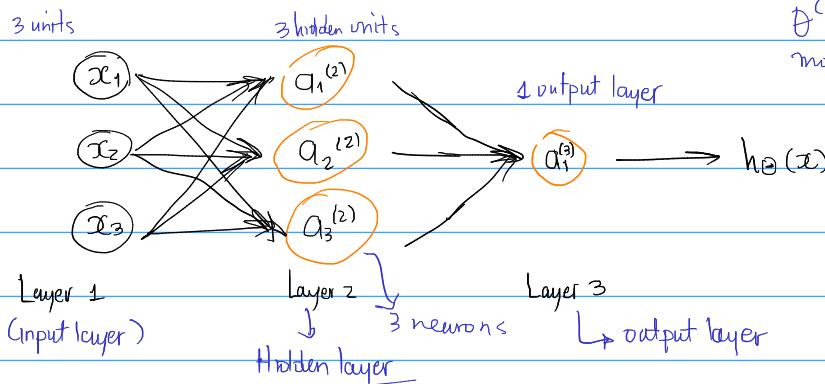
#### Neural Network



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

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

### Neural Network



\* Important \*

$\Theta^{(j)}$ : matrix of weights controlling function mapping from layer  $j=1$  to layer  $j+1=2$

(hidden layer (without bias) from it's layer)

$\Theta^{(j)} \in \mathbb{R}^{3 \times 4}$  input layer (with bias) from it's layer

$$\begin{matrix} (0) & (1) & (2) & (3) \\ (1) & \left[ \begin{array}{cccc} \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} & \theta_{13}^{(1)} \end{array} \right] & \left[ \begin{array}{c} x_0 \\ x_1 \\ x_2 \\ x_3 \end{array} \right] \\ (2) & \left[ \begin{array}{cccc} \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} & \theta_{23}^{(1)} \end{array} \right] & \left[ \begin{array}{c} x_0 \\ x_1 \\ x_2 \\ x_3 \end{array} \right] \\ (3) & \left[ \begin{array}{cccc} \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} & \theta_{33}^{(1)} \end{array} \right] & \left[ \begin{array}{c} x_0 \\ x_1 \\ x_2 \\ x_3 \end{array} \right] \end{matrix}$$

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

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

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

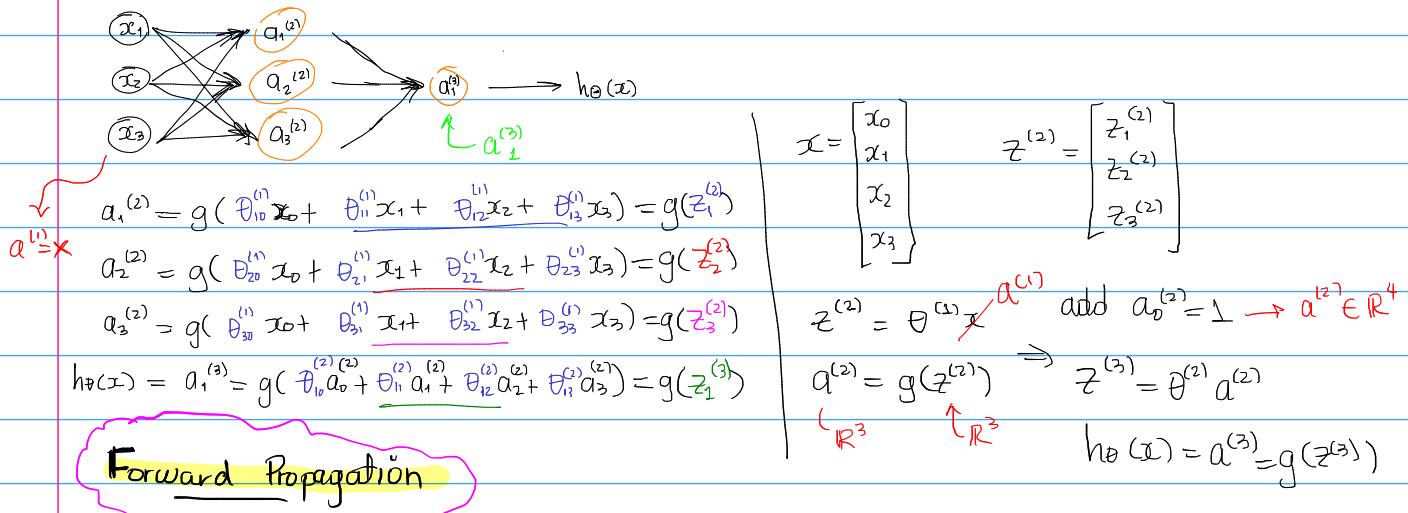
$$h_\theta(x) = a_1^{(3)} = g(\theta_{10}^{(2)}a_0^{(2)} + \theta_{11}^{(2)}a_1^{(2)} + \theta_{12}^{(2)}a_2^{(2)} + \theta_{13}^{(2)}a_3^{(2)})$$

$$\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

$$\left[ \begin{array}{cccc} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \end{array} \right], \left[ \begin{array}{c} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{array} \right]$$

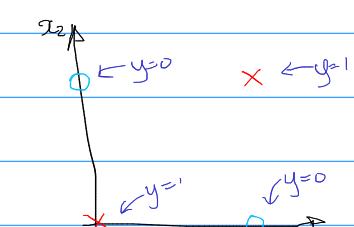
If network has  $S_j$  units in layer  $j$ ,  $S_{j+1}$  units in layer  $j+1$ , then  $\Theta^{(j)}$  will be of dimension  $S_{j+1} \times (S_j + 1)$

### — Neural Network: Representation (Model Representation II)



The are several different network architecture

## — Neural Network: representation (Examples and intuitions 1)

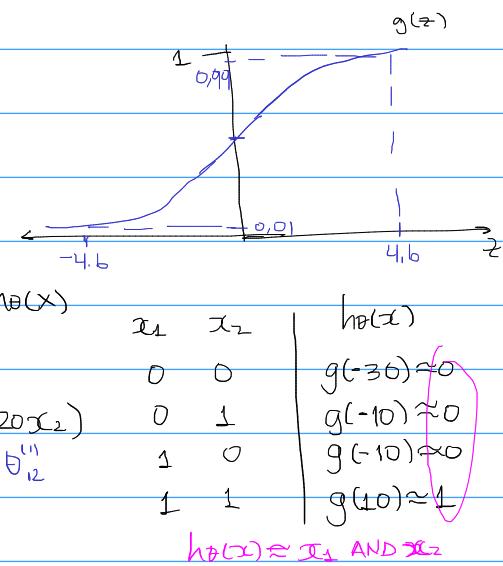


$y = x_1 \text{ XOR } x_2$   
 $x_1 \text{ XNOR } x_2$   
 $\text{NOT}(x_1 \text{ XOR } x_2)$

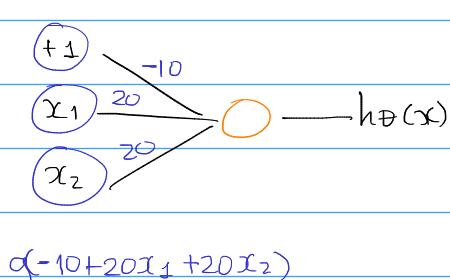
Simple Example: AND  
 $x_1, x_2 \in \{0, 1\}$   
 $y = x_1 \text{ AND } x_2$

$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

$$\theta^{(1)}_{10} \quad \theta^{(1)}_{11} \quad \theta^{(1)}_{12}$$



Example: OR function

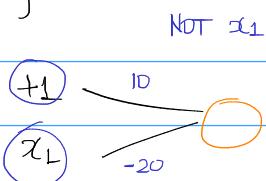


$x_1$	$x_2$	$h_{\theta}(x)$
0	0	0
0	1	1
1	0	1
1	1	1

$h_{\theta}(x) \approx x_1 \text{ OR } x_2$

## — Neural Networks: Representation (Examples and Intuitions II)

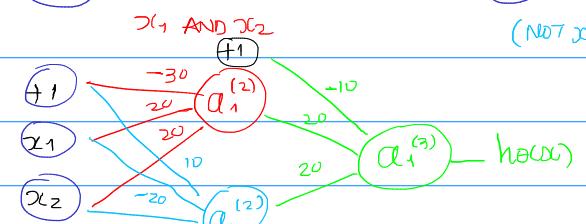
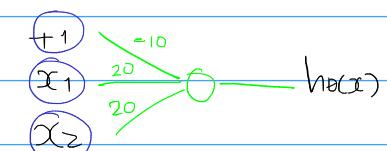
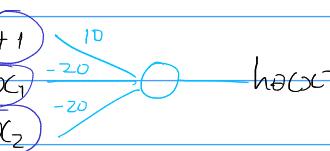
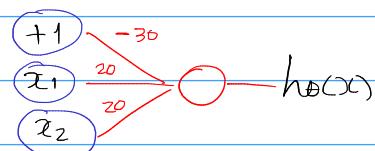
Negation



$x_1$	$h_{\theta}(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$h_{\theta}(x) = g(10 - 20x_1)$$

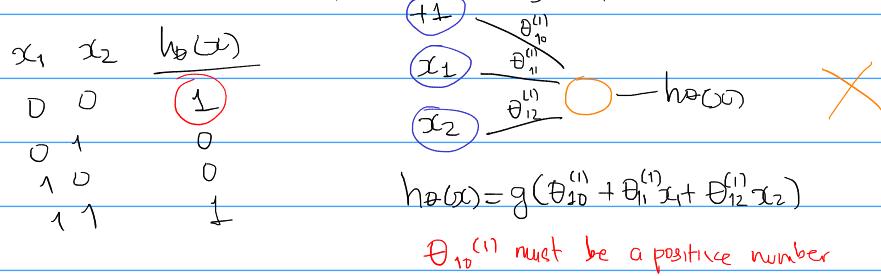
## — Putting it together —



$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

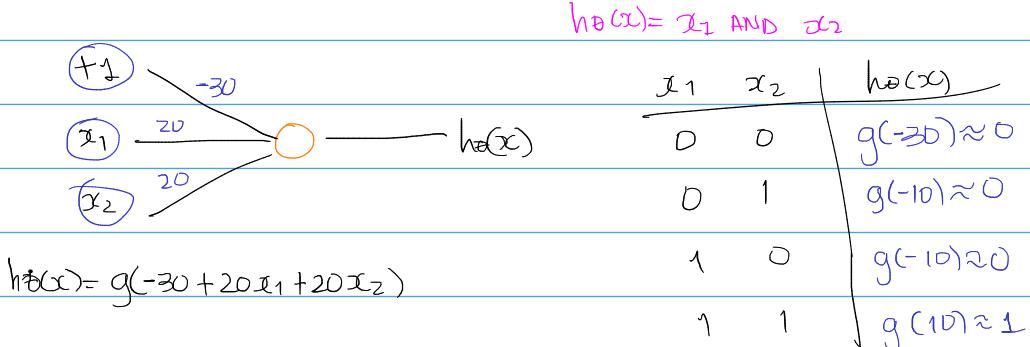
## Quiz

- ① - The activation values of the hidden units in a neural network, with the sigmoid activation function applied at every layer, are always in the range (0,1)

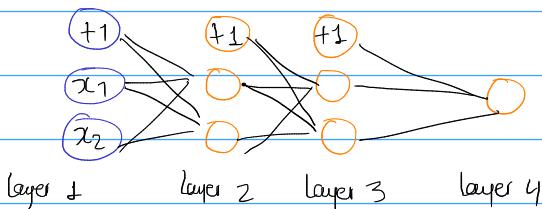


- Any logical function over binary-valued (0 or 1) inputs  $x_1$  and  $x_2$  can be (approximately) represented using some neural network.

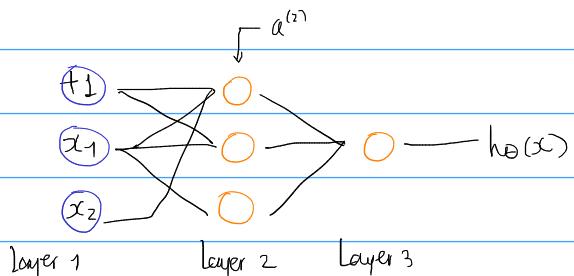
②



$$\text{③ } a_1^{(2)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)})$$



④



Compute the activations of hidden layer  $a^{(2)} \in \mathbb{R}^3$

$$\theta^{(1)} \in \mathbb{R}^{3 \times 3} \rightarrow \text{length}(input\ vector) + 1$$

$$\begin{matrix} (1) & \theta_{10}^{(1)} & \theta_{11}^{(1)} & \theta_{12}^{(1)} \\ (2) & \theta_{20}^{(1)} & \theta_{21}^{(1)} & \theta_{22}^{(1)} \\ (3) & \theta_{30}^{(1)} & \theta_{31}^{(1)} & \theta_{32}^{(1)} \end{matrix} \cdot \begin{matrix} (1) & +1 \\ (2) & x_1 \\ (3) & x_2 \end{matrix}$$

$$3 \times 3 = 3 \times 3$$

⑤

$$a^{(2)} = \theta^{(1)} x = \begin{bmatrix} 2 \times 3 & 3 \times 1 \\ 1 & -1.5 & 3.7 \\ 1 & 5.1 & 2.3 \end{bmatrix} \cdot \begin{bmatrix} +1 \\ x_1 \\ x_2 \end{bmatrix}$$

$$= \begin{bmatrix} 1 - 1.5x_1 + 3.7x_2 \\ 1 + 5.1x_1 + 2.3x_2 \end{bmatrix} = a^{(2)}$$

$$a^{(3)} = \theta^{(2)} a^{(2)} = \begin{bmatrix} 1 & 0.6 & -0.8 \end{bmatrix} \cdot \begin{bmatrix} 1 - 1.5x_1 + 3.7x_2 \\ 1 + 5.1x_1 + 2.3x_2 \end{bmatrix}$$

$$\theta^{(2)} \cdot a^{(2)} = \begin{bmatrix} \theta_{10}^{(2)} + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2 \\ \theta_{20}^{(2)} + \theta_{21}^{(2)}x_1 + \theta_{22}^{(2)}x_2 \\ \theta_{30}^{(2)} + \theta_{31}^{(2)}x_1 + \theta_{32}^{(2)}x_2 \end{bmatrix}$$

$$a^{(3)} = 1 + 0.6(1 - 1.5x_1 + 3.7x_2) - 0.8(1 + 5.1x_1 + 2.3x_2) = h_{\theta}(x) \quad 1$$

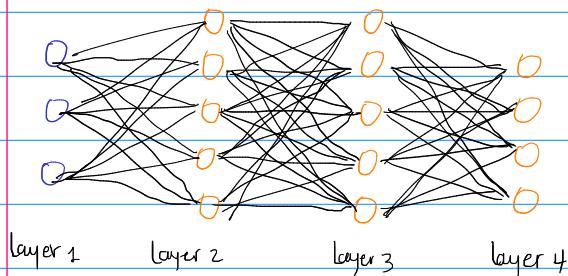
$$A^{(2)} = \theta^{(1)} x = \begin{bmatrix} 2 \times 3 \\ 1 & 5.1 & 2.3 \\ 1 & -1.5 & 3.7 \end{bmatrix} \begin{bmatrix} 3 \times 1 \\ +1 \\ x_1 \\ x_2 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 + 5.1x_1 + 2.3x_2 \\ 1 - 1.5x_1 + 3.7x_2 \end{bmatrix}$$

$$A^{(3)} = \theta^{(2)} A^{(2)} = \begin{bmatrix} 1 & -0.8 & 0.6 \\ 1 \times 3 \end{bmatrix} \begin{bmatrix} 3 \times 1 \\ +1 \\ 1 + 5.1x_1 + 2.3x_2 \\ 1 - 1.5x_1 + 3.7x_2 \end{bmatrix} = 1 - 0.8(1 + 5.1x_1 + 2.3x_2) + 0.6(1 - 1.5x_1 + 3.7x_2) = h_\theta(x)_2$$

$$\Rightarrow h_\theta(x)_1 = h_\theta(x)_2$$

## Week 5

### — Neural Networks: Learning (Cost function)



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

$L = \text{Total number of layers in Network}$

$S_l = \text{no. of units (not counting bias unit) in layer } l$ .

Binary Classification      Multi-class classification ( $K$  classes)

$$y = 0 \text{ or } 1$$

1 output unit

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

$$S_L = 1, K = 1$$

$$y \in \mathbb{R}^K \quad \text{e.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$K$  output units      Pedestrian Car motorcycle truck

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

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

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

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

### — Neural Networks: Learning (Backpropagation Algorithm)

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

$$\min_{\theta} J(\theta) \quad \theta_j^{(l)} \in \mathbb{R}$$

We need code to compute:  $-\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$

Gradient Computation

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

Forward Propagation:

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

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

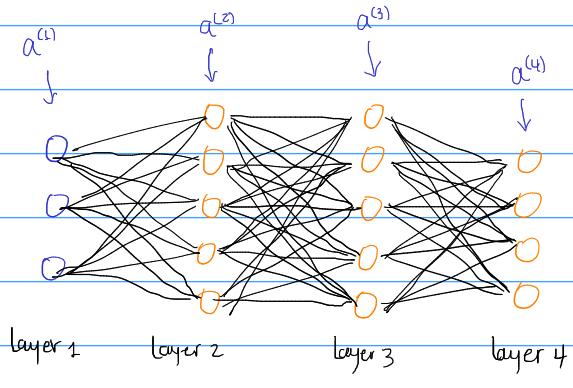
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

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

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

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

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



## Gradient Computation: Backpropagation algorithm

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

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

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \delta^{(4)} = a^{(4)} - y$$

vectorized form

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \rightarrow a^{(3)} \cdot (1-a^{(3)})$$

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

(No  $\delta^{(1)}$ )

$$\frac{\partial J(\Theta)}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} \cdot \delta_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda=0)$$

calculating  $g'(z^{(3)})$   $g(z^{(3)}) = \frac{1}{1+e^{-z}} \wedge g(z^{(2)}) = a^{(2)}$   
 let's  $z^{(3)} = z$

$$g'(z) = -1(1+e^{-z})^{-2} \cdot (-e^{-z})$$

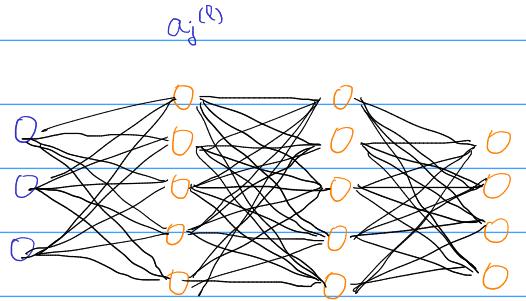
$$g'(z) = \frac{e^{-z} + 1 - 1}{(1+e^{-z})^2} \quad g(z)$$

$$g(z) = \frac{(1+e^{-z})}{(1+e^{-z})^2} - \frac{1}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} - \frac{1}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right)$$

$$g'(z) = g(z)(1-g(z)) \Rightarrow g'(z^{(3)}) = g(z^{(3)})(1-g(z^{(3)}))$$

$\star g(z^{(2)}) = a^{(2)}(1-a^{(2)})$

$\star g(z^{(2)}) = a^{(2)}(1-a^{(2)})$  \*



## Backpropagation algorithm

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

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

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

set  $a^{(1)} = \mathbf{x}^{(i)}$

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

using  $y^{(i)}$ , compute  $\delta^{(l)} = a^{(l)} - y^{(i)}$

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

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \quad \xrightarrow{\text{matrix}} \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

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

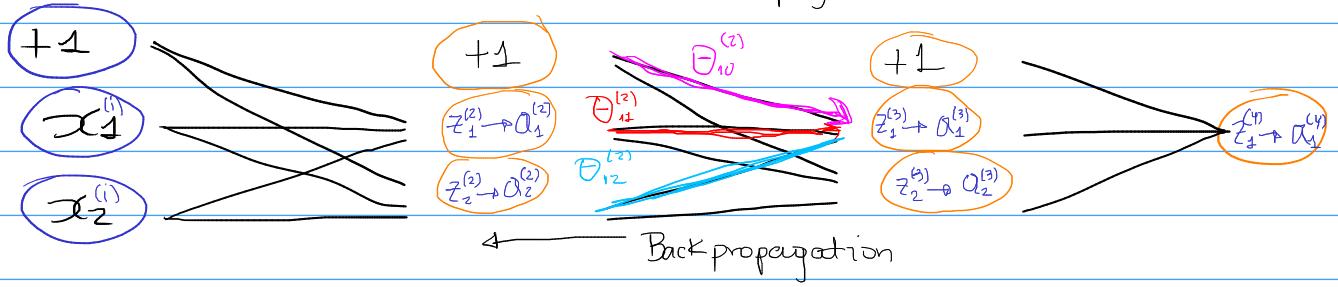
$$\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\Theta) = \Delta_{ij}^{(l)}$$

## — Neural Networks: Learning (Backpropagation Intuition)

Forward Propagation

→ Forward Propagation



$$(x^{(i)}, y^{(i)}) \quad z_1^{(2)} = \theta_{00}^{(2)} \times 1 + \theta_{11}^{(2)} \times x_1 + \theta_{12}^{(2)} \times x_2$$

What is backpropagation doing?

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

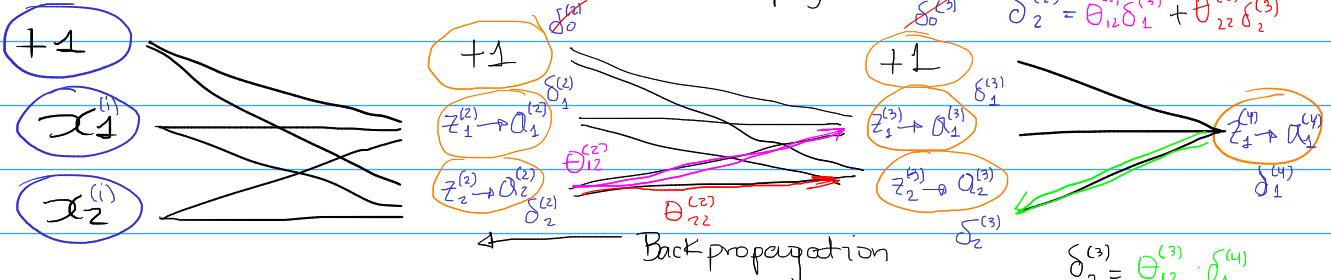
Focusing on a single example  $x^{(i)}, y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda=0$ ),

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

(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

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

## — Neural Networks: Learning (Implementation note: Unrolling parameters)

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.

$$\text{thetaVec} = [\Theta_{11}(:); \Theta_{12}(:); \Theta_{13}(:)]$$

$$DVec = [D_{11}(:); D_{12}(:); D_{13}(:)]$$

To get back to matrix representation

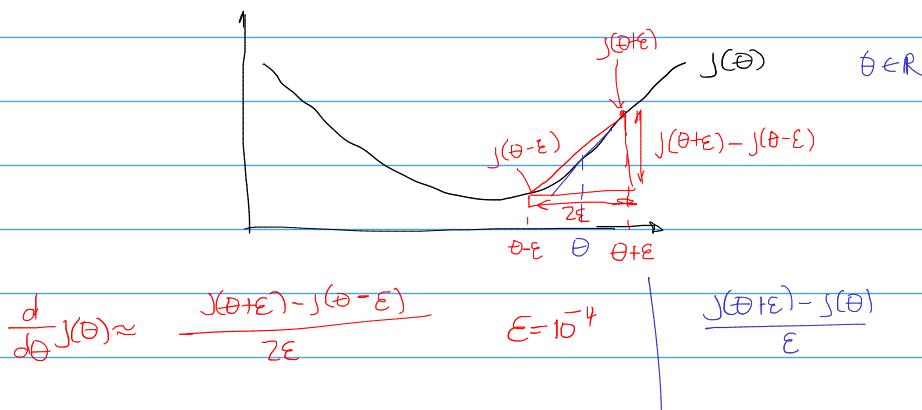
$$\Theta_{11} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$$

$$\Theta_{12} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$$

$$\Theta_{13} = \text{reshape}(\text{thetaVec}(221:231), 1, 11);$$

## — Neural Networks: Learning (Gradient Checking)

Numerical estimation of gradients



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

$$y = mx + c$$

the slope is  $y' = m$   
 or  $m = \frac{f(b) - f(a)}{b - a}$   
 where  $f(b) = mb + c$   
 $f(a) = ma + c$   
 $m = \frac{mb + \epsilon - ma - \epsilon}{b - a}$   
 $m = \frac{m(b - a)}{b - a} = m$

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

$\Theta = \Theta_1, \Theta_2, \Theta_3, \dots, \Theta_n$

$$\frac{\partial}{\partial \Theta_1} J(\Theta) \approx \frac{J(\Theta_1 + \epsilon, \Theta_2, \Theta_3, \dots, \Theta_n) - J(\Theta_1 - \epsilon, \Theta_2, \Theta_3, \dots, \Theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \Theta_2} J(\Theta) = \frac{J(\Theta_1, \Theta_2 + \epsilon, \Theta_3, \dots, \Theta_n) - J(\Theta_1, \Theta_2 - \epsilon, \Theta_3, \dots, \Theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \Theta_n} J(\Theta) = \frac{J(\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_n + \epsilon) - J(\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_n - \epsilon)}{2\epsilon}$$

From backprop.

```
for i = 1:n,
    thetaplus = theta;
    thetaplus(i) = thetaplus(i) + epsilon;
    theminus = theta;
    theminus(i) = theminus(i) - epsilon;
    gradapprox(i) = (J(thetaplus) - J(theminus)) / (2 * epsilon);
end.
```

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

code

## Implementation Note:

- Implement backprop to compute D<sub>θ</sub> (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 backpropagation code for learning

## 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.

## — Neural Networks : learning (Random initialization)

Initial value of  $\Theta$

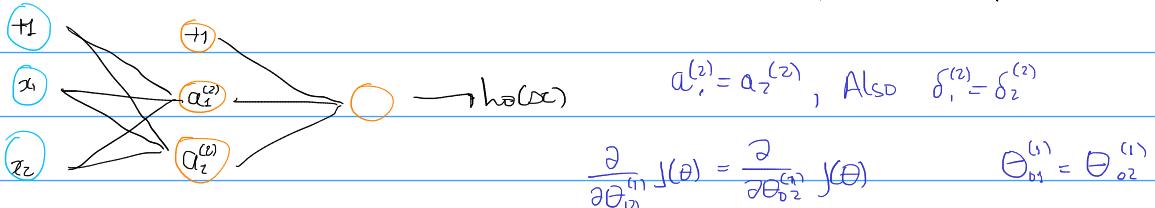
For gradient descent and advanced optimization method, need initial value of  $\Theta$   
 $\text{optTheta} = \text{fminunc}(@\text{costfunction}, \text{initialTheta}, \text{options})$

Consider gradient descent

Set  $\text{initialTheta} = \text{zeros}(n, 1)$  ?

Zero initialization

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



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

## 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.

$$\text{theta1} = \text{rand}(10, 1) \times (2 \times \text{INIT\_EPSILON}) - \text{INIT\_EPSILON}$$

$$\text{theta2} = \text{rand}(1, n) \times (2 \times \text{INIT\_EPSILON}) - \text{INIT\_EPSILON}$$

$$\begin{aligned} &x \xrightarrow{\text{EPR in } [0, 1]} \\ &x \cdot (2\epsilon) - \epsilon = \epsilon (2x - 1) \\ &\xrightarrow{\text{EPR in } [-1, 1]} \\ &\xrightarrow{\text{EPR in } [-\epsilon, \epsilon]} \end{aligned}$$

## — Neural Networks: Learning (Putting it together)

Training a neural network.

Pick a network architecture (connectivity pattern between neurons)

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

No. of output units: Number of classes.

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

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 derivative  $\frac{\partial}{\partial \theta_j^{(k)}} J(\theta)$

for  $i=1:m$  {

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

(Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l=1, \dots, L$ )

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

3. compute  $\frac{\partial}{\partial \theta_j^{(k)}} J(\theta)$

5. Use gradient checking to compare  $\frac{\partial}{\partial \theta_j^{(k)}} J(\theta)$  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  $\theta$ .  $J(\theta)$  - non-convex

## Quiz Neural Networks: Learning

$$\textcircled{1} \quad \Delta_{ij}^{(2)} := \Delta_{ij}^{(2)} + \delta_i^{(3)} \times (\alpha_j^{(2)})^T \quad \Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \alpha_j^{(l)} \delta_i^{(l+1)^T} \rightarrow \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (\alpha^{(l)})^T$$

$$\Delta^{(2)} = \Delta^{(2)} + \delta^{(3)} \times (\alpha^{(2)})^T$$

where  
i = training example  
j = j unit in layer.

$$\textcircled{2} \quad \text{Theta}_1 \in \mathbb{R}^{5 \times 3}, \text{Theta}_2 \in \mathbb{R}^{4 \times 6}, \text{theta vec} = [\text{theta}_1(:); \text{theta}_2(:)]$$

$$\Rightarrow \text{Theta}_1 = \text{reshape}(\text{theta vec}(1:15), 5, 3)$$

$$\Rightarrow \text{Theta}_2 = \text{reshape}(\text{theta vec}(16:39), 4, 6)$$

$$\textcircled{3} \quad J(\theta) = 2\theta^4 + 2 \quad \text{and} \quad D=1 \quad \text{and} \quad \epsilon = 0.01$$

$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} = \frac{J(1.01) - J(0.99)}{2 \cdot 0.01} = \frac{0.160016 - 0.0008}{2 \cdot 0.01}$$

$\textcircled{4}$  - using gradient checking can help verify if one's implementation of back propagation is bug-free

- For computational efficiency, after we have performed gradient checking to verify that our back propagation code is correct, we usually disable gradient checking before using backpropagation to train the network.

$\textcircled{5}$  - Suppose you are training a neural network using gradient descent. Depending on your random initialization, your algorithm may converge to different local optima...

- If we are training a neural network using gradient descent, one reasonable 'debugging' step to make sure it is working is to plot  $J(\theta)$  as a function of the number of iteration, and make sure it is decreasing ...