

function [jVal, gradient] = costFunction(theta)
jVal = [code to compute $J(\theta)$];

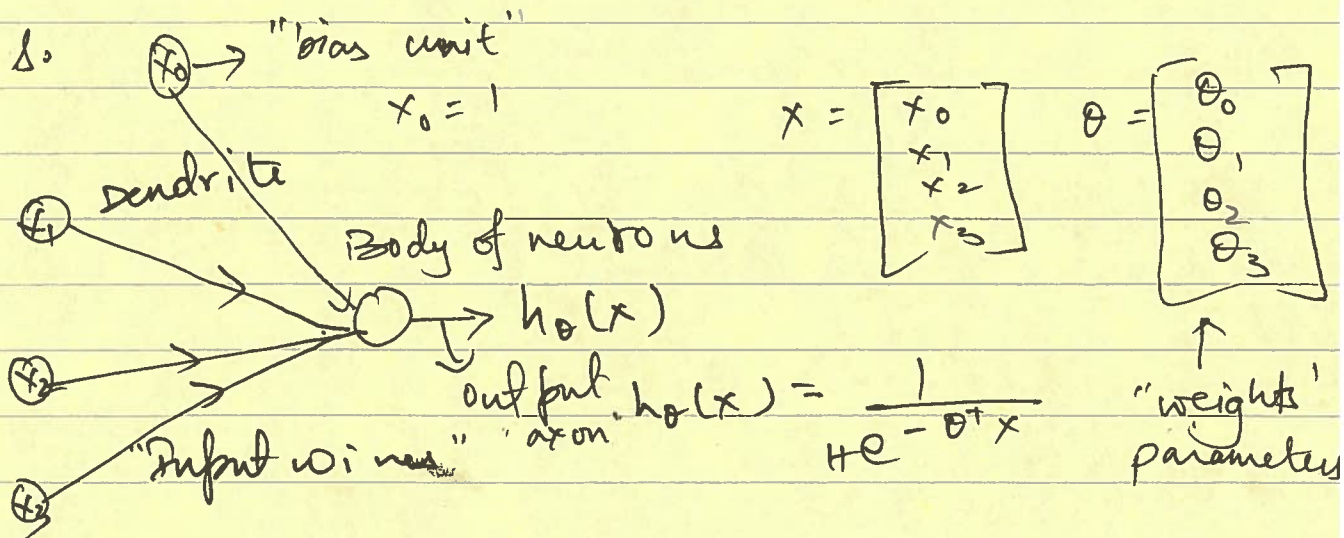
gradient(1) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

gradient(2) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$];

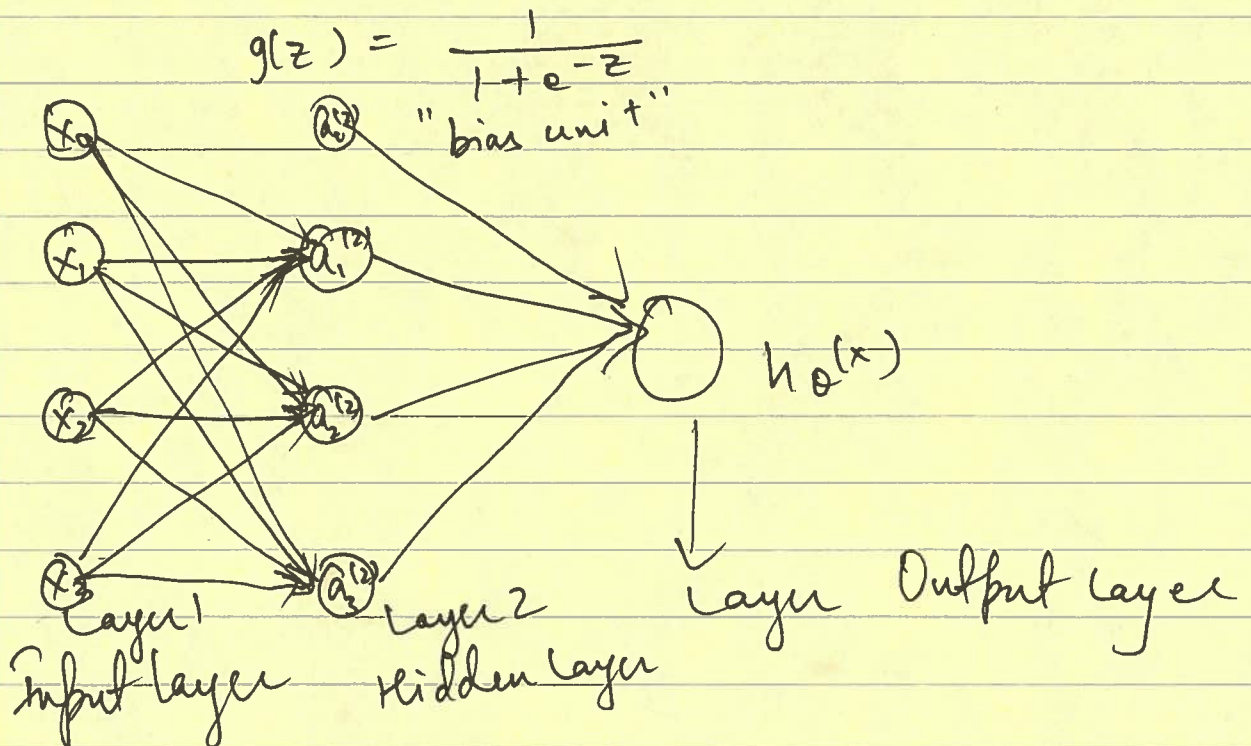
⋮

gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

Neurons → Dendrite (Input Wires)
 → Axon (Output Wire)



Sigmoid (logistic) activation function



Neural Network Rep Model Rep-I

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

θ^j = matrix of weights controlling
function mapping layer j to $j+1$

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

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

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

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

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

$$\begin{array}{cc} (s_{j+1}) \times (s_j + 1) \\ s_2 \times (s_1 + 1) \\ \underbrace{\text{Layer 1}} & \underbrace{\text{Layer 2}} \\ s_j = 2 & s_{j+1} = 4 \end{array}$$

$$\underbrace{s_{j+1}}_4 \times s_{j+1} \\ 4 \times (2+1)$$

Neural Network Ref Model Ref-II

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

$$a_1^{(2)} = g(z_1^{(2)})$$

$$z_k^{(2)} = \theta_{k,0}^{(1)} x_0 + \theta_{k,1}^{(1)} x_1 + \dots + \theta_{k,n}^{(1)} x_n$$

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

$$z_j = \theta^{(j-1)} a^{(j-1)}$$

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

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

$$\downarrow \mathbb{R}^3 \quad \uparrow \mathbb{R}^3$$

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

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

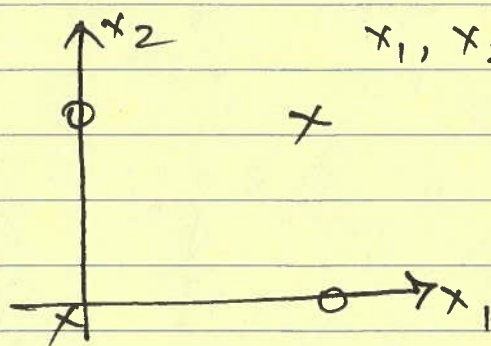
$$z^{(3)} = \theta^{(2)} a^{(2)}$$

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

$$h_\theta(x) = a^{(j+1)} = g(z^{(j+1)})$$

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

Examples & Intuition.



x_1, x_2 are binary (0 or 1)

$$y = x_1 \text{ XOR } x_2$$

$$x_1 \text{ XNOR } x_2$$

$$\text{NOT}(x_1 \text{ XOR } x_2)$$

XOR

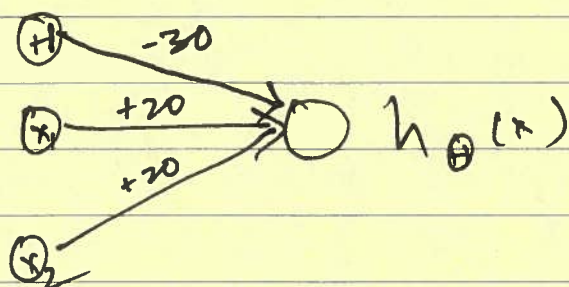
A	B	A XOR B
0	0	0
0	1	1
1	0	1
1	1	0

XNOR

A	B	A XNOR B
0	0	1
0	1	0
1	0	0
1	1	1

$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2$$



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

\downarrow \downarrow \uparrow
 $\theta_{10}^{(1)}$ $\theta_{11}^{(1)}$ $\theta_{12}^{(1)}$

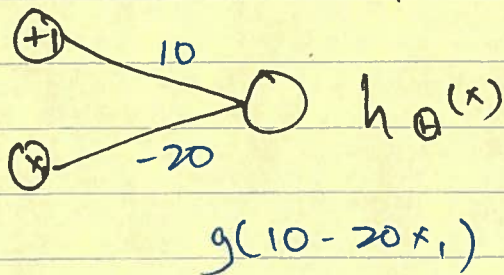
x_1	x_2	$h_{\theta} x$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(0) \approx 1$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

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

x_1	x_2	
0	0	$g(-10) \approx 0$
0	1	$g(20) \approx 1$
1	0	$g(20) \approx 1$
1	1	$g(40) \approx 1$

x_1 AND x_2
Not x_1

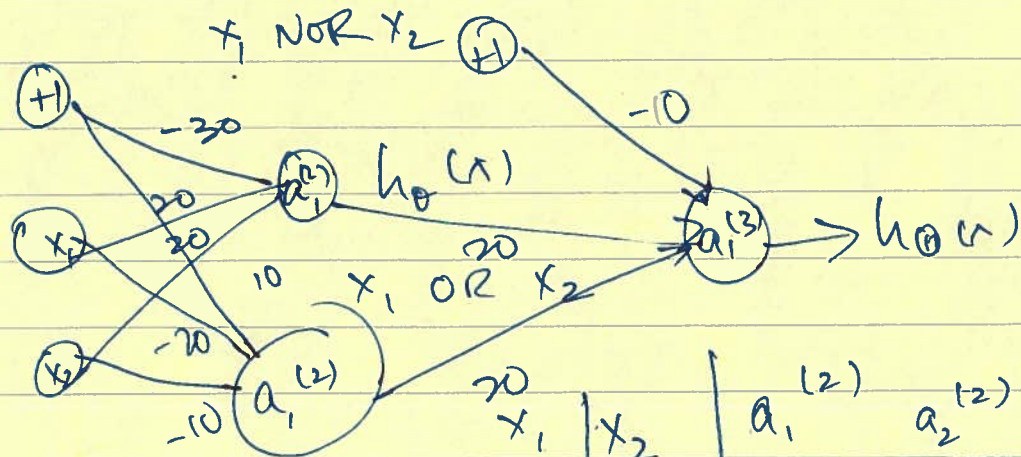
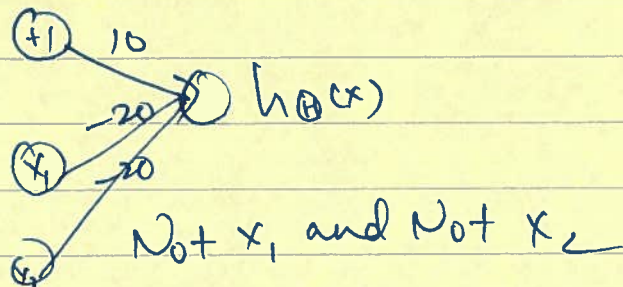
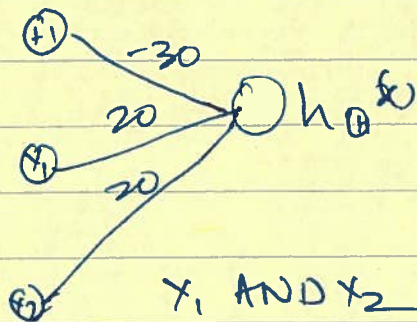


x_1 OR x_2 . $\{0, 1\}$

x_1	$h_\theta(x)$
0	$g(10) \approx 1$
1	$g(-20) \approx 0$

Not x_1 AND Not x_2
= 1 if any only if
 $x_1 = x_2 = 0$

Large



$a_1^{(3)} \rightarrow h_\theta(x)$
↑
output

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

$$g(-30 - 20(x_1) - 20(x_2))$$

$$g(-20 + 30(x_1) + 30(x_2))$$

x_1	x_2	
0	0	0
0	1	1
1	0	1
1	1	1

x_1	x_2	$h_0(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(-10) \approx 0$

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

L = total # of layers in network

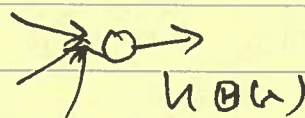
S_l = no. of units (not counting bias unit) in layer l

$$S_1 = 3, S_2 = 5, S_4 = S_L = 4$$

Binary classification

$y = 0$ or 1

1 output unit



$$h(\theta) \in \mathbb{R}$$

$$S_L = 1 \quad K = 1$$

Multiclass

$$y \in \mathbb{R}^K$$

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

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

Ex $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$
red car

Logistic regression

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

$K \rightarrow \text{dim}$

$i \rightarrow i \text{ element}$

Gradient computation.

$(x, y) \rightarrow \delta^{(2)} \rightarrow \delta^{(3)} \rightarrow \delta^{(4)}$

$a^{(1)}$	$a^{(2)}$	$a^{(3)}$	$a^{(4)}$
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

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

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

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

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

$$a^{(3)} = g(z^{(3)})$$

$$z^{(4)} = \textcircled{1}^{(3)} a^{(3)}$$

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

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

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

$$\delta_j^{(3)} = \left(\textcircled{H}^{(3)} \right)^T \delta^{(4)} \cdot g'(z^{(3)})$$

$$\delta^2 = \left(\textcircled{H}^{(2)} \right)^T \delta^{(3)} \cdot g'(z^{(2)})$$

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

$$\begin{aligned} & a^{(3)} \cdot (1 - a^{(3)}) \\ & a^{(2)} \cdot (1 - a^{(2)}) \end{aligned}$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

ignoring λ if $\lambda = 0$.

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

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

For $i = 1$ to m
set $a^{(1)} = x^{(i)}$

Perform forward prop to compute a^l for $l = 2, 3, \dots, L$

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

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

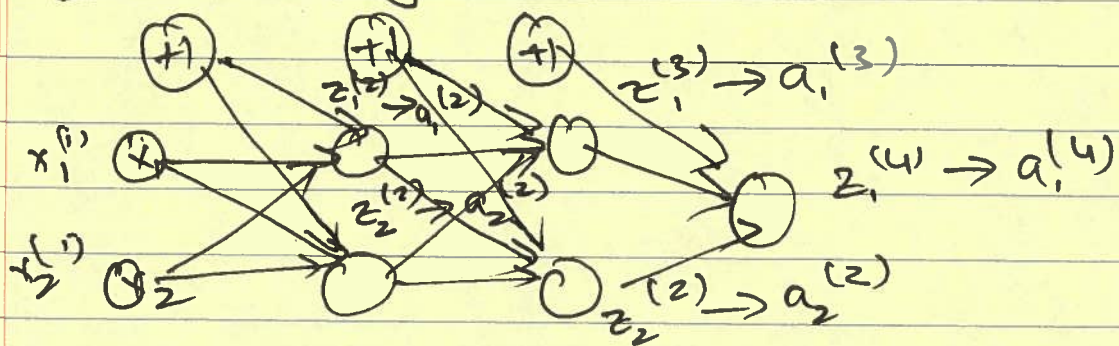
$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

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

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

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

Forward Propagation



$(x^{(1)}, y^{(1)})$

$$z_{11}^{(3)} = \textcircled{H}_{10}^{(2)} x_1 + \textcircled{H}_{11}^{(2)} x a_1^{(2)} + \textcircled{H}_{12}^{(2)} a_1^{(2)}$$