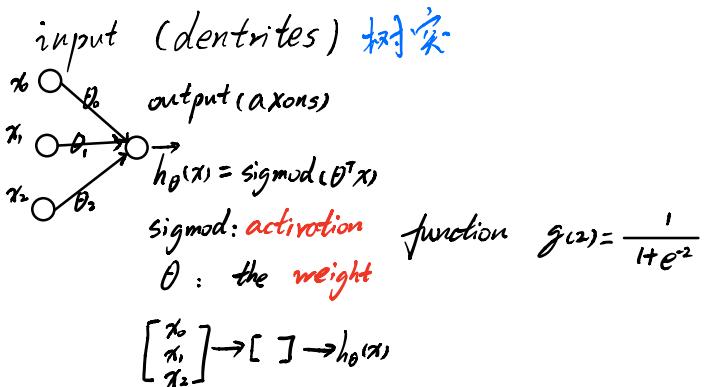
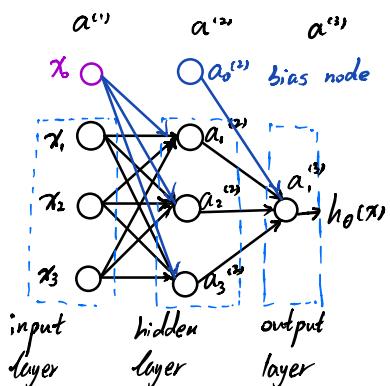


Neuron Network



$L = 3$, 3 layers



$a_0^{(0)}, a_0^{(1)}$: the "bias node", always 1

$a_0^{(2)}, a_1^{(2)}, a_2^{(2)}$: the "activation" unit

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} \rightarrow h_\theta(x)$$

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

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

each layer get its own matrix of weight $\theta^{(j)}$

If network has s_j units in layer j and s_{j+1} units in layer $j+1$, then $\theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$, the "+1" comes from "bias nodes"

e.g. layer 1 has 2 input and layer 2 has 4 activation nodes, dimension of $\theta^{(1)}$ is

$$4 \times 3 \quad (s_j = 2, s_{j+1} = 4, s_{j+1} \times (s_j + 1) = 4 \times 3)$$

$$a_1^{(2)} = g(\theta_{1,0}^{(2)}x_0 + \theta_{1,1}^{(2)}x_1 + \theta_{1,2}^{(2)}x_2 + \theta_{1,3}^{(2)}x_3) = g(z_1^{(2)})$$

$$a_2^{(2)} = g(\theta_{2,0}^{(2)}x_0 + \theta_{2,1}^{(2)}x_1 + \theta_{2,2}^{(2)}x_2 + \theta_{2,3}^{(2)}x_3) = g(z_2^{(2)})$$

$$a_3^{(2)} = g(\theta_{3,0}^{(2)}x_0 + \theta_{3,1}^{(2)}x_1 + \theta_{3,2}^{(2)}x_2 + \theta_{3,3}^{(2)}x_3) = g(z_3^{(2)})$$

$$h_\theta(x) = g(\theta_{0,0}^{(2)}a_0^{(2)} + \theta_{0,1}^{(2)}a_1^{(2)} + \theta_{0,2}^{(2)}a_2^{(2)} + \theta_{0,3}^{(2)}a_3^{(2)}) = g(z_0^{(2)})$$

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

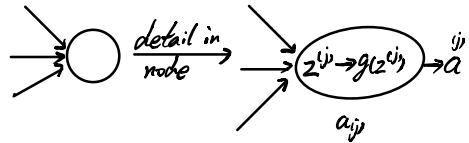
vectorize x and z^j :

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad z^{(j)} = \begin{bmatrix} z_0^{(j)} \\ z_1^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

Setting $x = a^{(0)}$

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

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



where $g(\cdot)$ applied element-wise to vector $z^{(j)}$,

when computed $a^{(j)}$, we add a bias node $a_0^{(j)} = 1$, compute another z vector

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

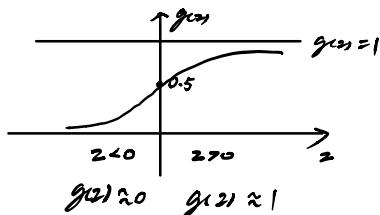
$$a^{(j+1)} = g(z^{(j+1)})$$

when comes to the last layer, $j = L-1$, $\theta^{(j)} = \theta^{(L-1)}$ has only one row,
the multiplication result is a single number, the final hypothesis:

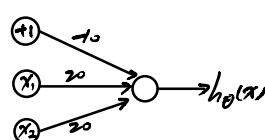
$$h_\theta(x) = a^{(L)} = g(z^{(L)})$$

Examples:

neural network used to simulate logical gates:



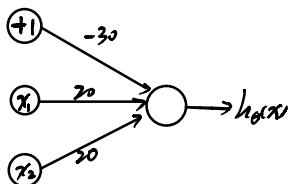
$$h_\theta(x) = g(z)$$



"OR" function

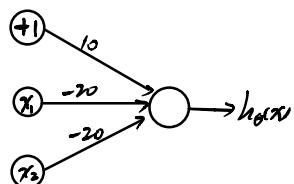
$$\theta = [-10 \ 20 \ 20]$$

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



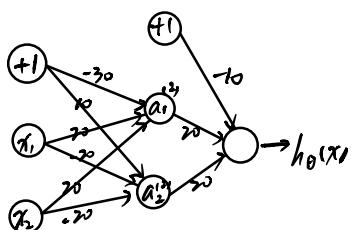
"AND" function

$$\theta^{(1)} = [-3.0 \ 2.0 \ 2.0]$$



"NOR" function

$$\theta^{(1)} = [1.0 \ -2.0 \ -2.0]$$



"XNOR" function

(combine AND and OR)

$$\theta^{(1)} = [-3.0 \ 2.0 \ 2.0]$$

$$\theta^{(2)} = [-1.0 \ 2.0 \ 2.0]$$

Multiclass Classification:

to classify class into multiple classes, we let hypothesis return a vector.

$$y^{(i)} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad (\text{4 classes})$$

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_n^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ \vdots \\ a_n^{(2)} \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} h_\theta(x_0) \\ h_\theta(x_1) \\ h_\theta(x_2) \\ \vdots \\ h_\theta(x_n) \end{bmatrix}$$

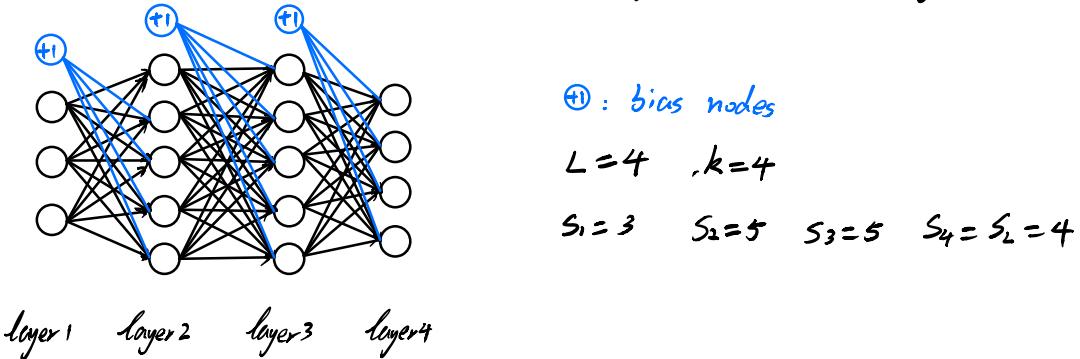
Backpropagation :

Cost function :

e.g. $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$, $k = \text{total number of output}$,

$L = \text{total number of layers in the network}$

$s_l = \text{number of units (not counting bias unit) in layer } l$



Binary classification

$$y = 0, 1$$

$$s_i = 1$$

Multi-class classification (k classes)

$$y \in \mathbb{R}^k, k \text{ output units}$$

$$s_i = k$$

$$k \geq 3$$

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 Networks: $J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^k y_k^{(i)} \log h_\theta(x^{(i)})_k + (1-y_k^{(i)}) \log (1-h_\theta(x^{(i)})_k) \right]$

$$+ \frac{\lambda}{2m} \sum_{i=1}^L \sum_{j=1}^{s_i} \sum_{l=1}^{s_{i+1}} (\theta_{ji}^{(l)})^2$$

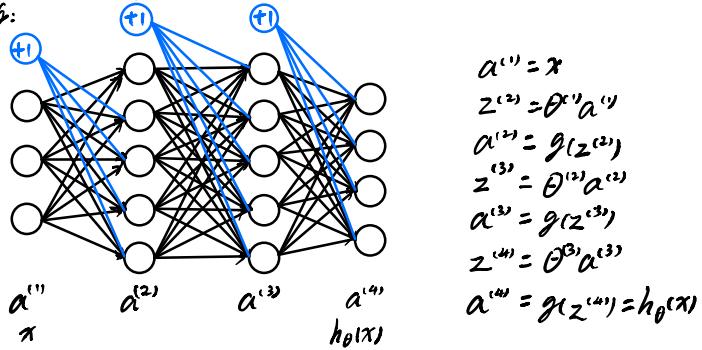
$h_\theta(x) \in \mathbb{R}^k$ ($h_\theta(x)_i = i^{\text{th}}$ output)

i, j start from 1: ignore the bias value

to compute $\min J(\theta)$, we need to compute $J(\theta)$ and $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$.

Gradient computation: Back propagation algorithm

e.g.:



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

for the neural network above: $L=4$

vectorize

$$\begin{cases} \delta_j^{(4)} = a_j^{(4)} - y_j = h_{\theta}(x)_j - y_j \\ \delta^{(4)} = a^{(4)} - y \end{cases}$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)}) \quad g'(z^{(3)}) = a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)}) \quad \delta^{(2)} = ((\theta^{(2)})^T \delta^{(3)}) * a^{(2)} * (1 - a^{(2)})$$

no $\delta^{(1)}$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda, \lambda=0)$$

$$\left[\begin{array}{c} \vdots \\ \vdots \end{array} \right]$$

Give the training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$,

set $\Delta_{ij} = 0$ (for all i, j, l) accumulator for computing $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

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

set $a^{(0)} = x^{(i)}$

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

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

compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ no $\delta^{(1)}$

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

?

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

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

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$$

backpropagation intuition:

recall the cost function for a neural network is:

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

if non-multiclass classification ($k=1$) and disregard regularization ($\lambda=0$)

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

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(t)$$