

## Model and Cost function

Not's:

$m$  : no. of training ex's

$x$ 's : input vars/features

$y$ 's : output vars/features target var. that we are trying to predict

$(x, y)$  : one training ex./observ

$(x^{(i)}, y^{(i)})$ :  $i$ th training ex./observ

$(x^{(i)}, y^{(i)})$ ;  $i = 1, 2, \dots, m$  : training set

$i$  : index into the training set

$X$  : space of I/P vals

$Y$  : space of O/P vals

$n$  : no. of features

$$\text{Hypothesis} : h_{\theta}(x) = \theta_0 + \theta_1 x$$

$\theta_0, \theta_1$  : para. of model

#  $h$  maps from  $x$ 's to  $y$ 's

## Neural Networks

Neural Network is a series of algo.s that endeavours to recognize underlying relationships in a set of data through a process that mimics the way the human brain operates.

- NNs were widely used in 80s and 90s; popularity diminished in late 90s
- Recent resurgence; state-of-art techniques for many applic's
- NNs are much more expensive in terms of comput'. One of the reason for their resurgence is that the comput' needed is only recently matched by powerful computers built in last decade.

Need : Algo.s like Lin. Reg. & Log Reg. (simple or polynomial) is not a good way to learn complex non-linear hypotheses when 'n' is large bc it generates too many features.

$2500 \times 3$  if RGB

ex:  $50 \times 50$  pixel image  $\rightarrow 2500$  pixels

$$n = 2500$$

$$\text{grayscale val. b/w } 0-255 \quad X = \begin{bmatrix} \text{pixel 1 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix} \quad O(n^q) \approx \frac{n^q}{a}$$

$n \rightarrow$  features  
 $a \rightarrow$  polynomial

NNs are a better way to learn complex non-linear hypo., even when  $n$  is large

## Model Represent

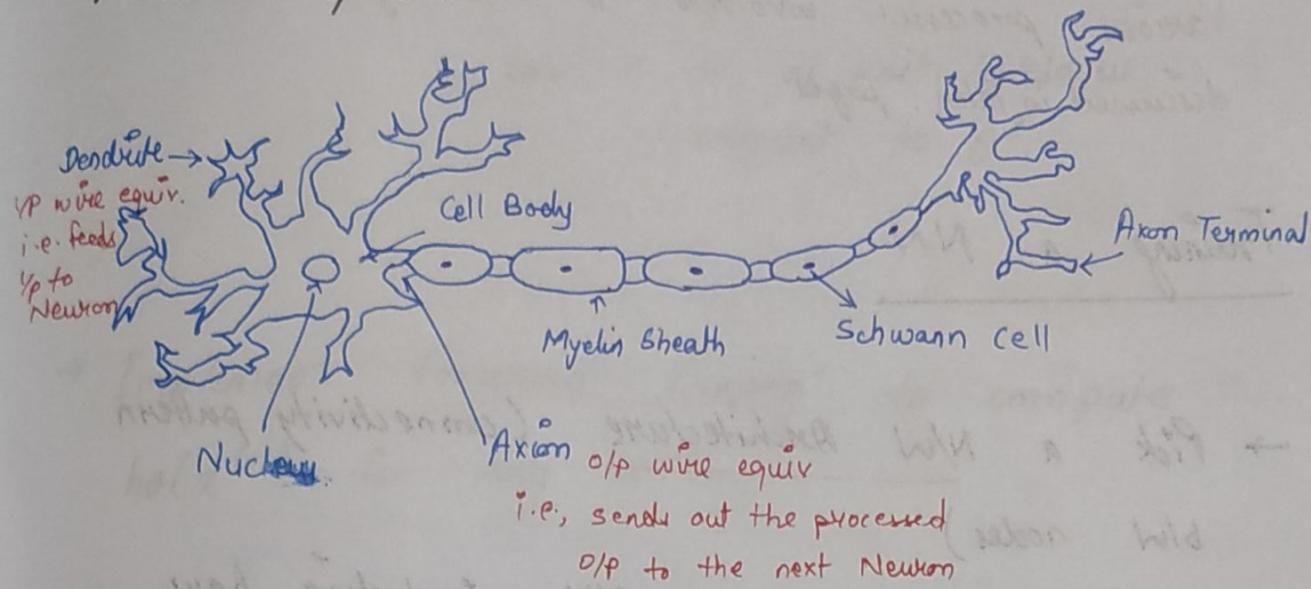


fig. Neuron

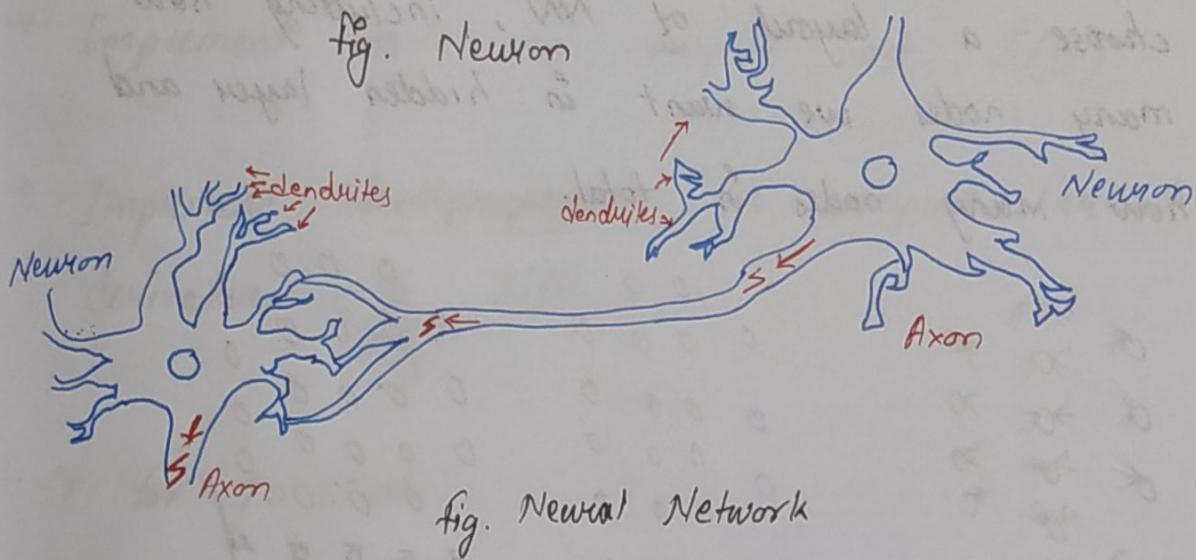


fig. Neural Network

Neurons communicate w/ each other using little pulses of electricity called spikes.

So, a neuron basically sends a spike and the other neurons basically accept the info., performs some computation on it & passes to the next Neuron.

This process is repeated until the desired process is done.  
ex: Contraction of Muscle.

# Neural Network Implement<sup>n</sup> (step-by-step Algo.)

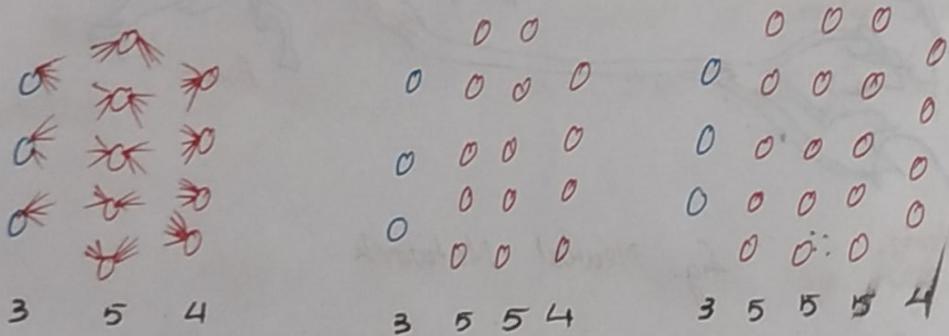
# this is an overview to generate an idea of the whole picture.

Various processes involved in this Algo. are discussed in later pages.

## Training a NN

→ Pick a N/W Architecture (connectivity pattern b/w nodes)

choose a layout of NN, including how many nodes we want in hidden layers and how many nodes in total.



• no. of  $\text{I/p}$  units : dimensions of features  $x^{(i)}$

• no. of  $\text{o/p}$  units : no. of classes  $y = \{0, 1\}$  or  $y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$

• reasonable default : 1 hidden layer, or

if  $>1$  hidden layers, have same no. of hidden units in every layer (more expensive / accurate)

no. of hidden layers =  $n * \text{no. of I/p units}$ ;  $n \in \mathbb{N}$

→ Randomly initialize weights ( $\theta$ ) in range  $[-\epsilon, \epsilon]$

where,  $\epsilon = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$

$L_{in}, L_{out}$  = no. of units in layer adjacent to  $\theta^l$

→ Implement forward Propag to compute  $h_\theta(x^{(i)})$  for any  $x^{(i)}$

→ Implement code to compute Cost  $J(\theta)$

→ Implement Backpropagation to compute Partial derivatives  $\frac{\partial}{\partial \theta_k^{(l)}} J(\theta)$

\* for  $i = 1 : m$

perform fwd propag & back prop using  $(x^{(i)}, y^{(i)})$

get activa's  $a^{(l)}$  & delta terms  $(\delta^{(l)}, \Delta^{(l)})$

$\delta^{(l)}$  for  $l = 2, \dots, L$

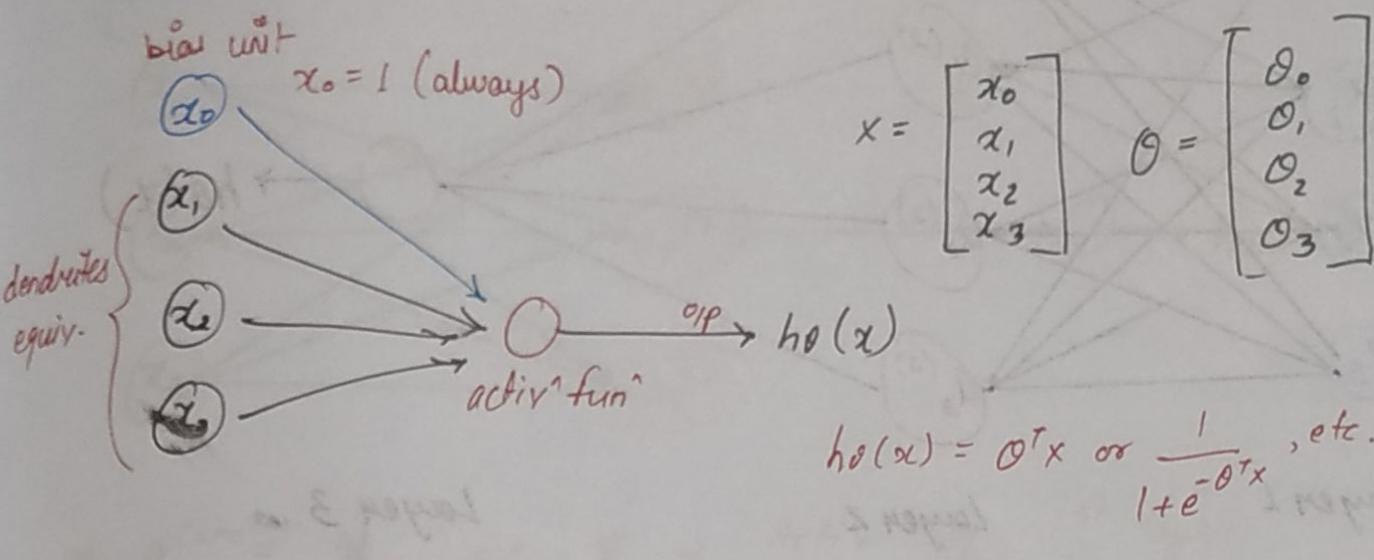
$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

Compute  $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

\* /

- Use Gradient Checking to compare  
 $\frac{\partial}{\partial \theta_{jk}} J(\theta)$  computed using backprop v/s  
using Numerical estim' of grad. of  $J(\theta)$ .
- # If these vals. are similar, our implement' is correct & we can proceed.  
If not, debug the prev. implement'  
Then disable gradient checking code.
- Use Gradient Descent or any other Advanced optimiz' method w/ backprop to try to minimize  $J(\theta)$  as a fun' of para.s ( $\theta$ )
- # for NN,  $J(\theta)$  is a Non-Convex fun' and so algo.s like grad. desc. can theoretically be stuck in the local minima.  
But in practice, it is found that it is usually not a great problem and usually algo.s like grad. desc. do a great job at minimizing  $J(\theta)$

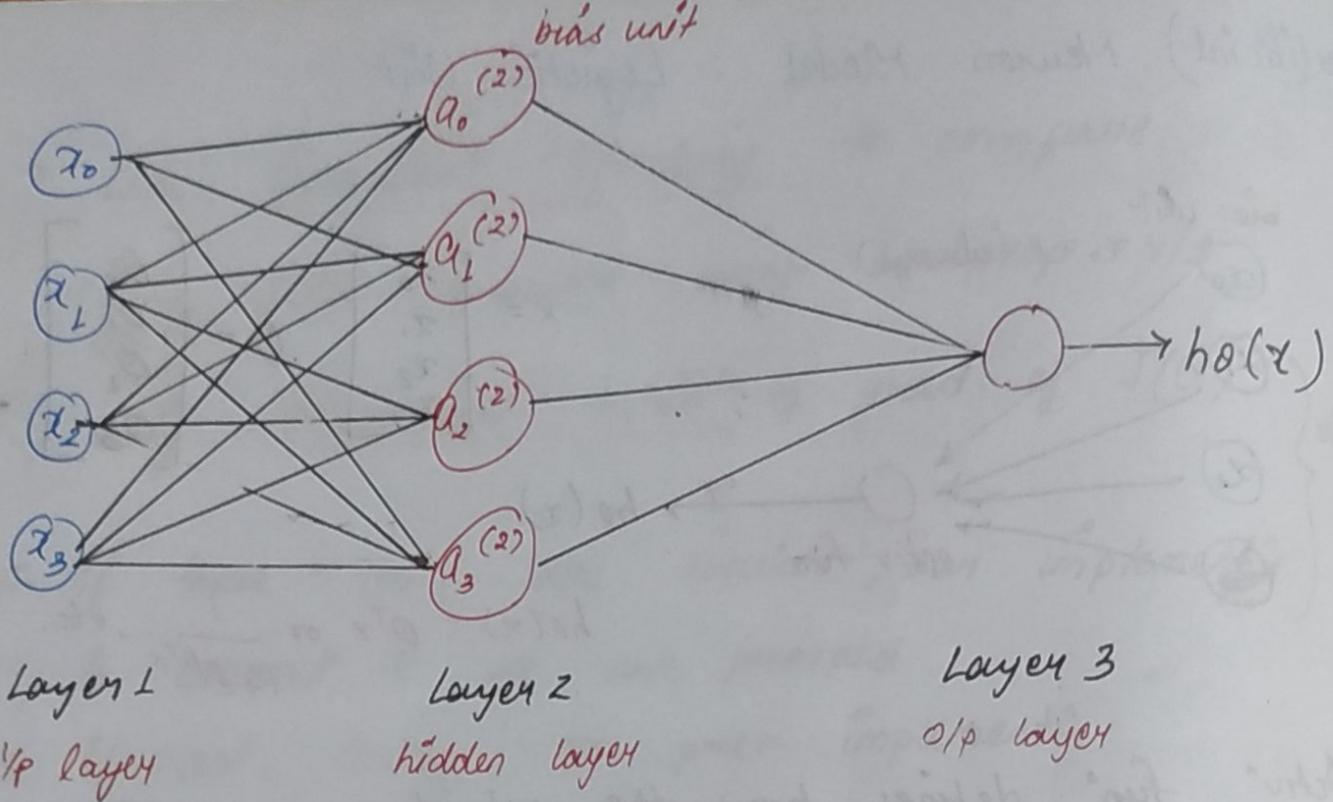
# (Artificial) Neuron Model : Logistic Unit



"Activ' fun" defines how the wt-ed sum of I/P is transformed into an O/P from a node i.e., it delivers an O/P based on YPs.

The purpose of "Activ' fun" is to add non-linearity to the N/W.

"Activ' fun" decides whether or not a neuron should be activated i.e., it decides whether Neuron's Yp to the N/W is imp. to the process of predict.



$a_j^{(j)}$  → activ<sup>n</sup> of unit "j" in layer "j"

$\theta^{(j)}$  → matrix of wt.s controlling fun<sup>n</sup> mapping from layer "j" to layer "j+1"

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

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

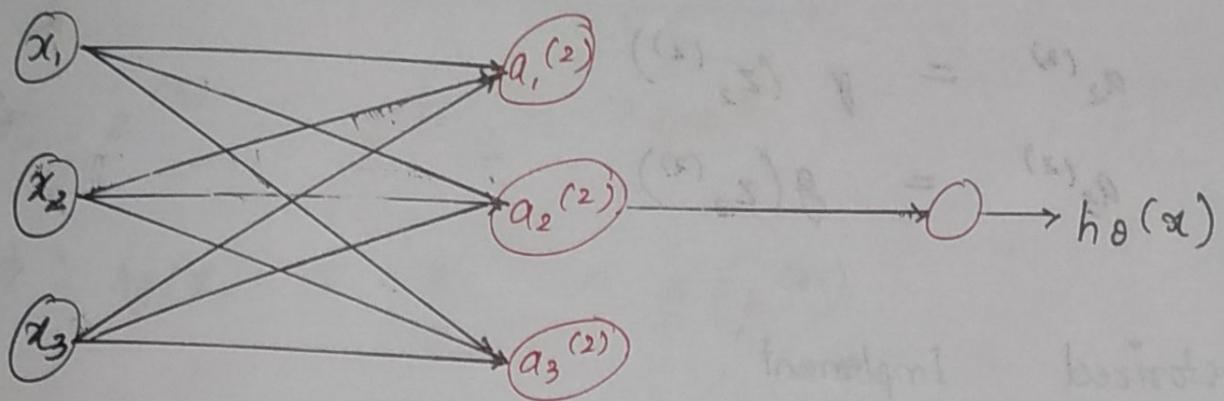
$$a_3^{(2)} = \sigma(\theta_{30}^{(1)}x_0 + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3)$$

activ<sup>n</sup> fun<sup>n</sup>

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

If n/w has  $s_j$  units in layer "j",  $s_{j+1}$  units in layer "j+1" then  $\theta^{(j)}$  will be of dimensions  $s_{j+1} \times (s_j + 1)$

# Forward Propagation : Vectorized Implement^



#  $a^{(l)} = \text{X}$  i.e.,  $\text{X}$  is Activ^ of  $1^{st}$  layer

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

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

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

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

let us define a new var.  $z_k^{(j)}$  that encompasses the para.s  $(\theta^j / a_i^{(j)})$  inside activ^ fun^  $g()$ .

$$z_k^{(j)} = \theta_{k,0}^{(j)}x_0 + \theta_{k,1}^{(j)}x_1 + \theta_{k,2}^{(j)}x_2 + \dots + \theta_{k,n}^{(j)}x_n$$

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

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)})$$

Vectorized Implement^

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

setting  $x = a^{(1)}$ , we can rewrite the eq<sup>1</sup> as

$$\begin{aligned} z^{(j)} &= \theta^{(j-1)} a^{(j-1)} \\ a^{(j)} &= g(z^{(j)}) \end{aligned}$$

so,  $z^{(1)} = \theta^{(0)} a^{(1)}$  fun<sup>n</sup>  $\theta$  is applied el.-wise on  ~~$z$~~   $z^{(j)}$

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

$\hookrightarrow \in \mathbb{R}^3$        $\hookleftarrow \in \mathbb{R}^3$

Now, adding a bias unit ( $=1$ ) to layer  $j$  after we have computed  $a^{(j)}$  ( $a_0^{(j)} = 1$ )

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

$\therefore z^{(3)} = \theta^{(2)} a^{(2)} \approx EIR^4$

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

We get this final  $z$  vec. by multiplying the next  $\theta$  maty. after  $\theta^{(j-1)}$  w/ the vals. of all the activ' nodes we get

This last theta maty.  $\theta^{(j)}$  will have only 1 row wh. is multiplied w/ 1 col.  $a^{(j)}$  so that our result is a single val. (predicted val.)

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

\* At each neuron at each layer after the l/p layer,

$$\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \quad \text{or} \quad \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \quad \text{or} \quad \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \quad \text{or} \quad \begin{array}{c} 0 \\ 0 \\ 0 \end{array}$$

we are doing exactly what we do in other algos (ax,  $\frac{1}{1+e^{-ax}}$ , etc.) here:  $g(z)$

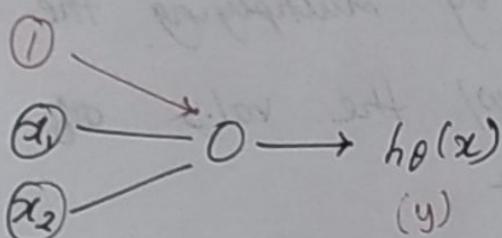
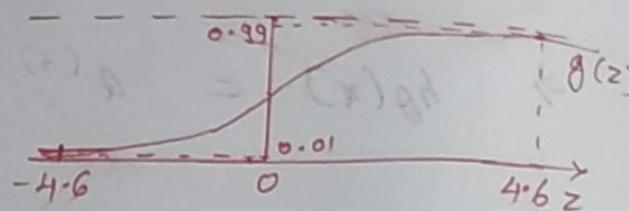
Adding all these intermediate layers in NN's allows us to more elegantly produce more complex and interesting hypo. \*

# Neural Network example

→ (simple) Logical AND

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

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



$$\text{let, } \theta = [-30 \ 20 \ 20]$$

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

$$h_\theta(x) = g(\theta^T x) = g\left(\frac{-30}{-30+20x_1+20x_2}\right)$$

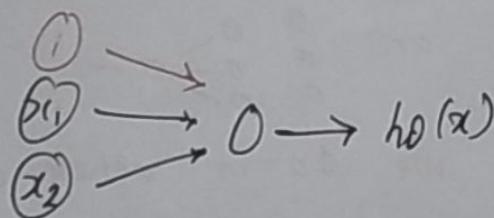
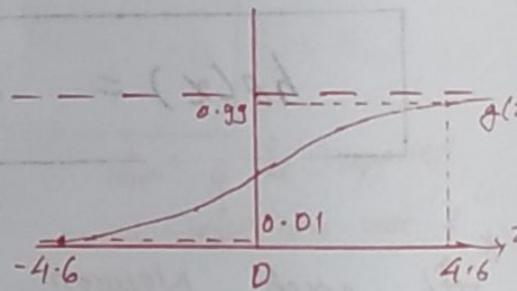
Logical AND

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

(simple) Logical OR

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

$$\rightarrow y = x_1 \text{ OR } x_2$$



$$\text{let, } \theta = [-10 \ 20 \ 20]$$

$x_1$	$x_2$	$h_\theta(x)$
0	0	$g(-10) \approx 0$
0	1	$g(20) \approx 1$
1	0	1
1	1	1

$$\Rightarrow h_\theta(x) = g(\theta^T x) = g(-10 + 20x_1 + 20x_2) \text{ logical OR}$$

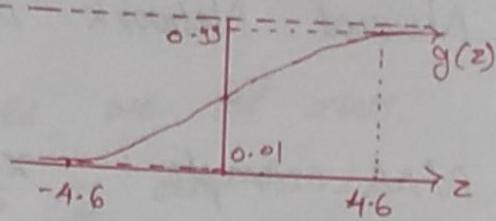
→ (simple) Logical NOT

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

$$\rightarrow y = h_\theta(x) = \text{NOT } x_1$$

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

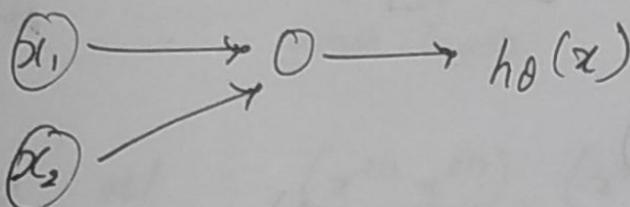
$$\Rightarrow h_\theta(x) = g(\theta^T x) = g(10 - 20x_1)$$



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

Logical NOT

→ (complex) ( $\text{NOT } x_1$ ) AND ( $\text{NOT } x_2$ )



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

$$\Rightarrow h_\theta(x) = g(\theta^T x)$$

$$= g(10 - 20x_1 - 20x_2)$$

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

# Multiclass Classification using Neural Networks

To classify data into multiple classes, we let our hypo. fun<sup>n</sup> return a vec. of vals. of dim ( $n \times 1$ ) where  $n$  is the no. of classes.

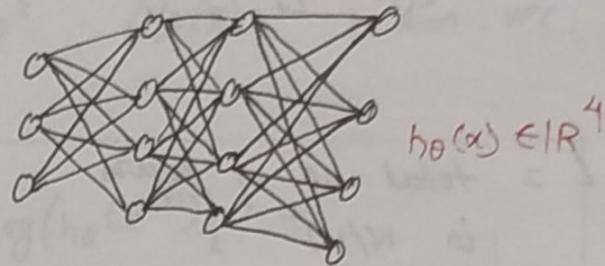
say,  $y \in \{\text{Pedestrian, Car, Bike, Truck}\}$

want,

$$h_{\theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when Pedestrian

when Car



$$h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when Bike

when Truck

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

where  $y \in \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right\}$

so,  $y, h \in \mathbb{R}^n$

where,  $n \rightarrow \text{no. of discrete classes}$

∴ Generalizing

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

where  
 $\leftarrow h_{\theta}(x_i) \in \mathbb{R}^n$

# Neural Network Learning

## NN (Classification)

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

$L$  = total no. of layers  
in N/W

here,  $L = 4$

$s_L$  = no. of units in layer  $L$   
(not including bias unit)

here,

$$s_1 = 3, s_2 = 5, s_3 = 5$$

$$s_4 = s_L = 4$$

$k$  = no. of units in o/p layer

## Binary Classification

$y = 0 \text{ or } 1 \quad \{y \in \{0, 1\}\} \text{ i.e., 1 o/p unit}$

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

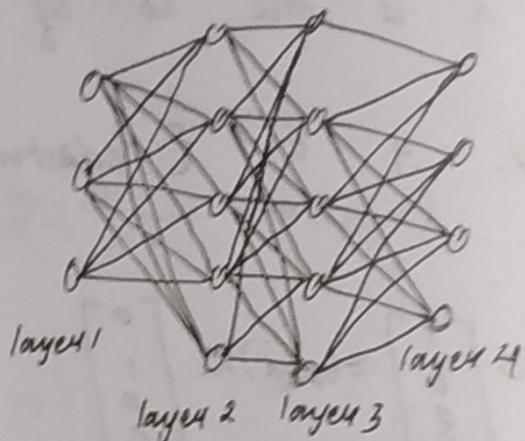
$$s_L = 1 \quad k = 1$$

## Multi-class Classification ( $K$ classes)

$y \in \mathbb{R}^k$  ( $k$  dimensional vec.)

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

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



$$\text{ex: } y \in \left[ \begin{array}{c} 1 \\ 0 \\ 0 \end{array} \right], \left[ \begin{array}{c} 0 \\ 1 \\ 0 \end{array} \right], \left[ \begin{array}{c} 0 \\ 0 \\ 1 \end{array} \right], \left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ 1 \end{array} \right]$$

## Cost Function

log. Reg.

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

Neural Network:  $h_\theta(x) \in \mathbb{R}^k$ ;  $(h_\theta(x))_i = i^{\text{th}} \text{ o/p in } k\text{-dim. vec.}$

$$J(\theta) = \underbrace{-\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K}_{\text{1. add'l nested summ' that loops through the no. of o/p nodes}} \left[ y_k^{(i)} \cdot \log(h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{j;i}^{(l)})^2$$

1.  $\rightarrow$  add'l nested summ' that loops through the no. of o/p nodes

2. no. of cols in  $\theta$  maty = no. of cols in curr layer (including bias unit)

3. no. of rows in  $\theta$  maty = no. of nodes in next layer (excluding bias unit)

NOTE:

$\rightarrow$  double sum simply adds the log-reg. cost calc'd for each cell in the o/p layer

$\rightarrow$  triple sum simply adds up the sq's of all individual  $\theta$ 's in the entire N/W

$\rightarrow$  the ' $i$ ' in triple sum doesn't refer to training ex.  $i$

## Backpropagation

"Backpropag" is a technique to minimize Cost  
in Neural Networks

Goal : min  $J(\theta)$

we need to compute: ①  $J(\theta)$

$$\textcircled{2} \quad \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \quad \left| \begin{array}{l} \theta_{ij}^{(l)} \in \mathbb{R} \\ \text{i.e., a real no.} \end{array} \right.$$

we calc. the Partial Derivative of Cost fun' for optimiz' algo.s like Grad. Desc. using the following Algo :

## Backpropagation Algorithm

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

→ set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ )  $\left| \begin{array}{l} \text{used to compute} \\ \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \end{array} \right.$

→ for  $\ell = 1 \text{ to } m$  i.e., for  $i^{\text{th}}$  item, we'll work w/  $(x^{(i)}, y^{(i)})$

- set  $a^{(0)} = x^{(i)}$  <sup>activ<sup>n</sup></sup>

- perform "fwd propag" to compute for  $\ell = 2, 3, \dots, L$

- using  $y^{(t)}$ , compute  $\delta^{(\ell)} = a^{(\ell)} - y^{(t)}$

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

$$\delta^{(\ell)} = \underbrace{\text{np.dot}\left((\theta^{(\ell+1)})^T \cdot \delta^{(\ell+1)}\right)}_{\substack{\text{element-wise} \\ \text{multiplic}^n}} \cdot \underbrace{g'(z^{(\ell)})}_{\substack{\hookrightarrow \text{derivative} \\ \text{of activ<sup>n</sup> fun<sup>n</sup>}}}$$

Intu<sup>n</sup>:  $\delta_j^{(\ell)}$  = error of node  $j$  in layer  $\ell$

$\delta^{(L)}, \delta^{(L-1)}, \dots, \delta^{(2)}$  → only until layer 2 bc we don't calc. error for l/p layer

for sigmoid fun<sup>n</sup>, we know that,

$$g'(z^{(\ell)}) = a^{(\ell)} \cdot * (1 - a^{(\ell)})$$

$$\therefore \delta^{(\ell)} = ((\theta^{(\ell)})^T \delta^{(\ell+1)}) \cdot * a^{(\ell)} \cdot * (1 - a^{(\ell)})$$

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

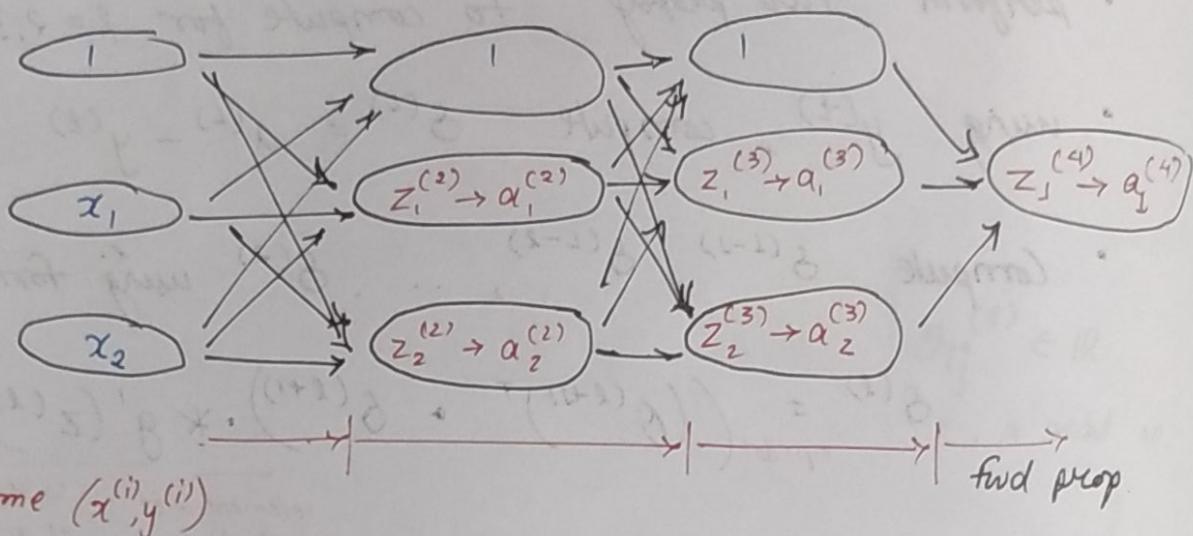
$$D_{ij}^{(\ell)} := \frac{1}{m} \Delta_{ij}^{(\ell)} \quad \begin{matrix} \text{if } j = 0 \\ \text{+ corresponds} \\ \text{to bias term} \end{matrix}$$

The matrx 'D' is used as an "accumulator" to add up our vals as we go along & eventually compute our Partial Derivative. Thus, we get (upon mathematical calc.)

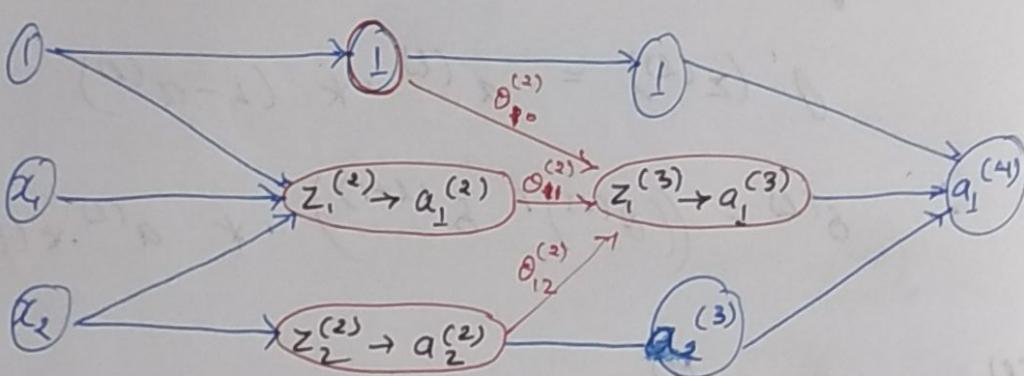
$$\frac{\partial}{\partial \theta_{ij}^{(\ell)}} = D_{ij}^{(\ell)}$$

# Backpropagation Intuition

Forward propagation



Understanding  $\perp$  of the node in  $l$  of the layers (selected at random)



$$z_1^{(3)} = \theta_{10}^{(2)} \times \perp + \theta_{11}^{(2)} \times a_1^{(2)}$$

$$+ \theta_{12}^{(2)} \times a_2^{(2)}$$

for Intu<sup>1</sup>, let us consider a simple non-multiclass classifier ( $k=1$ ) ~~to~~

ignoring regulariz<sup>n</sup> ( $\lambda=0$ ), we get

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

$$\text{cost}(t) \approx (h_{\theta}(x^{(t)}) - y^{(t)})^2$$

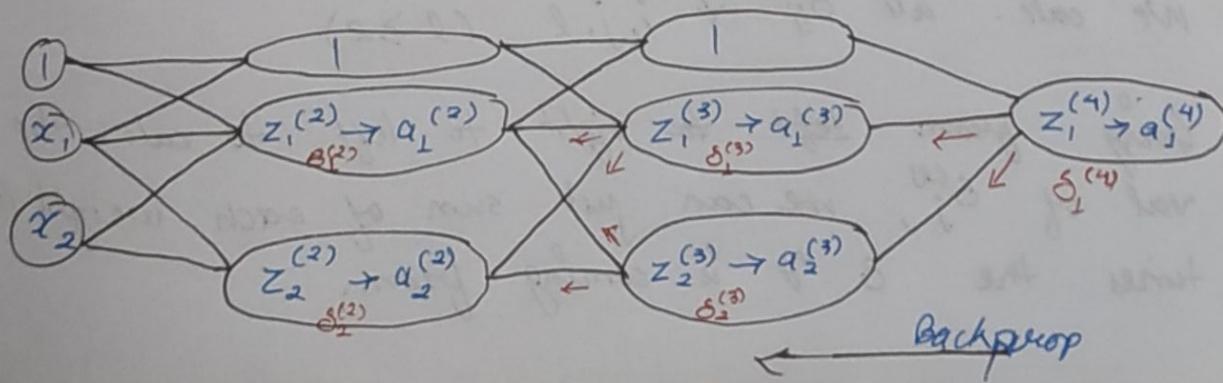
\* Intuitively, Backpropagation computes  $\delta_j^{(l)}$  i.e., the "error" of cost of  $a_j^{(l)}$

More formally, delta vals are actually the derivatives of the Cost Fun<sup>1</sup>.

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(t) \quad (\text{for } j \geq 0)$$

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

\* Error,  $\delta_j^{(l)} = a_j^{(l)} - y^{(i)}$  (calculated val - actual val.)

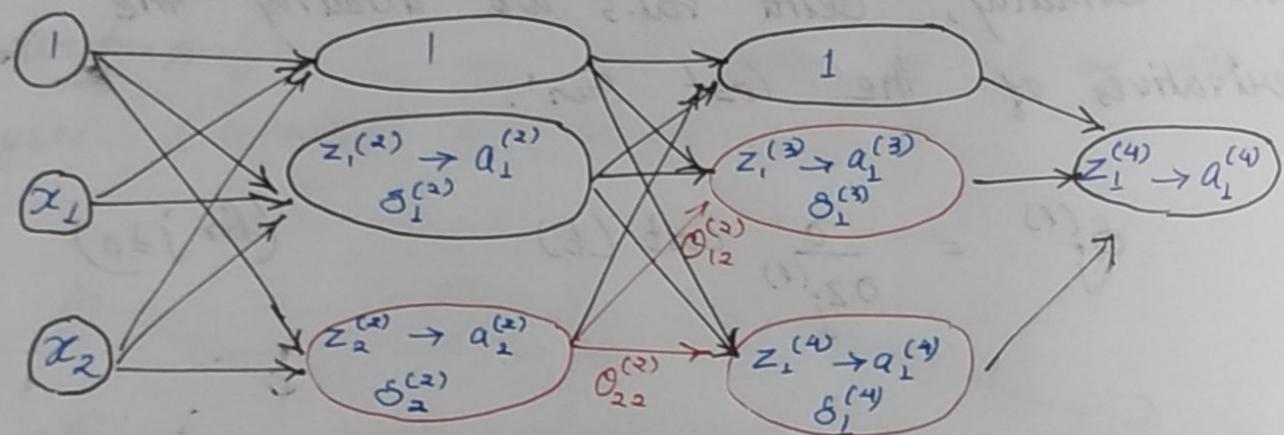


first, we calc. the error for the o/p layer

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

then, we propagate backwards to calc. the error for nodes of prev. layers (until layer 2)

Understanding  $\delta_1^{(2)}$  of the node of  $l$  of the column, (selected at random)



$$\delta_2^{(2)} = \theta_{12}^{(2)} \times \delta_1^{(3)} + \theta_{22}^{(2)} \times \delta_2^{(4)}$$

$$\text{Similarly, } \delta_3^{(3)} = \delta_1^{(4)} \times \theta_{12}^{(3)}$$

We calc. all  $\delta_{ij}^{(l)}$  &  $i, j, l$  ( $l \geq 2$ )

Going from left to right to left, to calc. the val of  $\delta_j^{(l)}$ , we can just sum of each weight ( $\theta$ ), times the  $\delta$  it is coming from.

## Implement<sup>^</sup> Note : Chaining Params

w/ NNs, we work w/ set of matrices :

$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots$  (params / wt.s) and  
 $D^{(1)}, D^{(2)}, D^{(3)}, \dots$  (grad.s), etc.

We can unroll these matrices into vectors.

$$\# [D]_{m \times 3} = [D]_{m \times 1} + [D]_{m \times 1} + [D]_{m \times 1}$$

The advantage of matr. represent<sup>^</sup> of paras is that it allows us to take advantage of Vectorized Implement<sup>^</sup> and is more convenient while performing forward and backpropagation.

The advantage of vec. repr. of paras is that some advanced optimiz<sup>^</sup> algo.s / methods (in Octave) assume that the paras passed to them are "unrolled" into a vec. and thus, vec. represent<sup>^</sup> of paras allows us to use these Advanced Optimiz<sup>^</sup> algo.s / methods.

## Gradient Checking

## Numerical Gradient Checking

Backpropag' is a tricky algo.

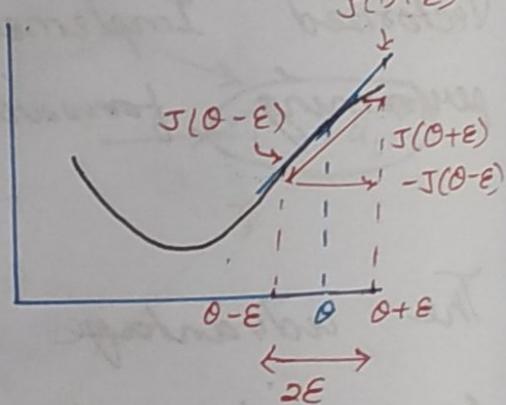
As a result, upon implement' and use of an optimiz' algo like grad. desc., it may appear that the para.s are optimizing but the bugs may not eventually minimize  $J(\theta)$

So, (Numerical) Gradient Checking is a technique (an expensive technique) to check if implement' of backpropag' works as intended.

we can approximate the grad. of our Cost Fun' w/ :

$$\frac{\partial J(\theta)}{\partial \theta} \approx \frac{J(\theta+E) - J(\theta-E)}{2E}$$

↖ 2 sided  
difference



↖ 1 sided  
difference

$$\frac{\partial J(\theta)}{\partial \theta} \approx \frac{J(\theta+E) - J(\theta)}{2E}$$

w/ multiple theta matrices,  $\theta \in \mathbb{R}^n$  (Unrolled/Vectorized), we can approximate the derivative as

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

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon}$$

i.e.,

$$\frac{\partial}{\partial \theta_j} J(\theta) \approx \frac{J(\theta_1, \dots, \theta_j + \epsilon, \dots, \theta_n) - J(\theta_1, \dots, \theta_j - \epsilon, \dots, \theta_n)}{2\epsilon}$$

# a small val. for  $\epsilon$  such as  $\epsilon = 10^{-4}$  guarantees that the math works out properly

# if  $\epsilon$  is too small, we can end up w/ numerical problems

Implement :

# epsilon = 1e-4

for i from 1 to n :

thetaPlus = theta

thetaPlus(i) = thetaPlus(i) + epsilon

thetaMinus = theta

thetaMinus(i) = thetaMinus(i) - epsilon

gradApprox(i) =  $\frac{J(\text{thetaPlus}) - J(\text{thetaMinus})}{2 \times \text{epsilon}}$

# Check that gradApprox ≈ DVec

Implement Note :

→ Implement backprop to compute DVec

→ Implement grad. checking to compute gradApprox

→ Check if gradApprox and DVec gives similar vals.

If yes, our implement<sup>n</sup> of backprop is correct

If not, debug backprop implement

→ Turn off grad. checking bc it is expensive in comput<sup>n</sup> power & time

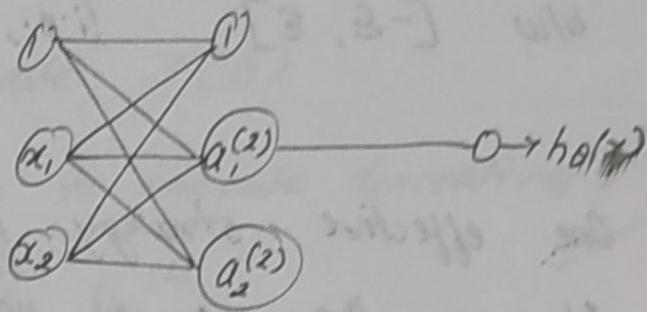
→ Use backprop code for learning i.e., for further optimiz<sup>n</sup>.

## Random Initialization

Need : Zero Initialization

Initializing all the wts i.e.,  $\theta$  to 0 does not work.  
as all nodes will update the same val. repeatedly.

After each update, para.s  
corresponding to '1's going  
into each of 2 hidden  
units are identical, i.e.,



$$a_1^{(2)} = a_2^{(2)} \Rightarrow \delta_1^{(2)} = \delta_2^{(2)}$$

$$\frac{\partial}{\partial \theta_{01}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{02}^{(1)}} J(\theta)$$

$$\theta_{01}^{(1)} = \theta_{02}^{(1)}$$

so no matter how many units in the hidden  
layer, each of them would compute the same  
val as they are fed the same no.s ([0])

## Random Initialization : Symmetry Breaking

to avoid the aforementioned problem, we initialize each  $\theta_{ij}^{(l)}$  to a random val.

b/w  $[-\epsilon, \epsilon]$  (i.e.,  $-\epsilon \leq \theta_{ij}^{(l)} \leq \epsilon$ )

One effective strategy for selecting  $\epsilon$  is to base it on the no. of units in the NW.

A good choice is

$$\epsilon = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$$

where,  $L_{in} = s_e$  and  $L_{out} = s_{e+1}$  are the no. of units in the layer adjacent to  $\theta_e$

Implement :

# init-epsilon = some val.

Theta1 = np.random.rand(row\_size, col\_size) \* (2 \* init-epsilon) - init-epsilon

Theta2 = "

Theta3 = "

# Neural Network learning (Summary)

- pick a N/W Architecture
- randomly initialize wts ( $\theta$ ) in range  $[-\epsilon, \epsilon]$
- Implement forward propagation to compute  $h_{\theta}(x)$
- implement code to compute  $J(\theta)$
- implement backpropagation to compute derivatives  
$$\frac{\partial}{\partial \theta_k^{(l)}} J(\theta) \quad | \quad \delta^{(l)}, \Delta^{(l)}, D_{ij}^{(l)}$$
- use gradient checking to compare gradient computed w/ backpropagation w/ gradient computed w/ numerical estimation
- if not, debug backpropagation  
If yes, disable gradient checking code
- use Gradient Descent or any other Advanced Optimiz" Technique w/ Backpropagation to ~~to~~  
~~min~~  $J(\theta)$
- ~~train~~ ~~#~~