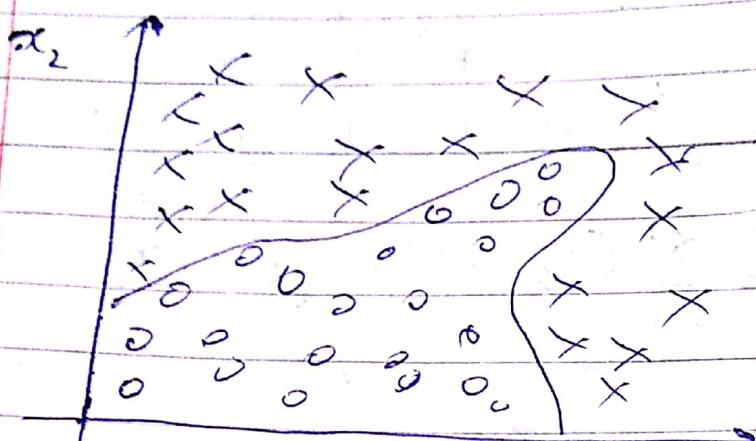


My  
Wing  
know  
our score.

## Week 4

Suppose we need to classify this:



Many non-linear features will be needed to apply logit's regression:

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2)$$

In order to properly fit the decision boundary. The features only increase dramatically in real world situations due to high # of features.

Suppose real world feature set of houses:

$$x_1 = \text{size}$$

$$x_2 = \# \text{ of bedrooms}$$

$$x_3 = \# \text{ of floors}$$

$$x_4 = \text{age}$$

...

$$x_{100}$$

} 100 features  
(n = 100),

Counting the second order features from the data set fields:

$$\left. \begin{array}{c} x_1, x_2, x_1 x_2, x_1 x_3, \dots, x_1 x_{100} \\ x_2, x_1 x_2, x_2 x_3, x_2 x_4, \dots, x_2 x_{100} \\ \vdots \\ x_1 x_2 x_3, x_1 x_2 x_4, \dots, x_1 x_2 x_{100} \end{array} \right\} \text{for } n=100$$

$\approx 5000$  features  
 $[O(n^2)]$

Creates overfitting → asymp. growth rate of quadratic features as  $\#$  of features increases.

finding out parameter values begins to become computationally expensive

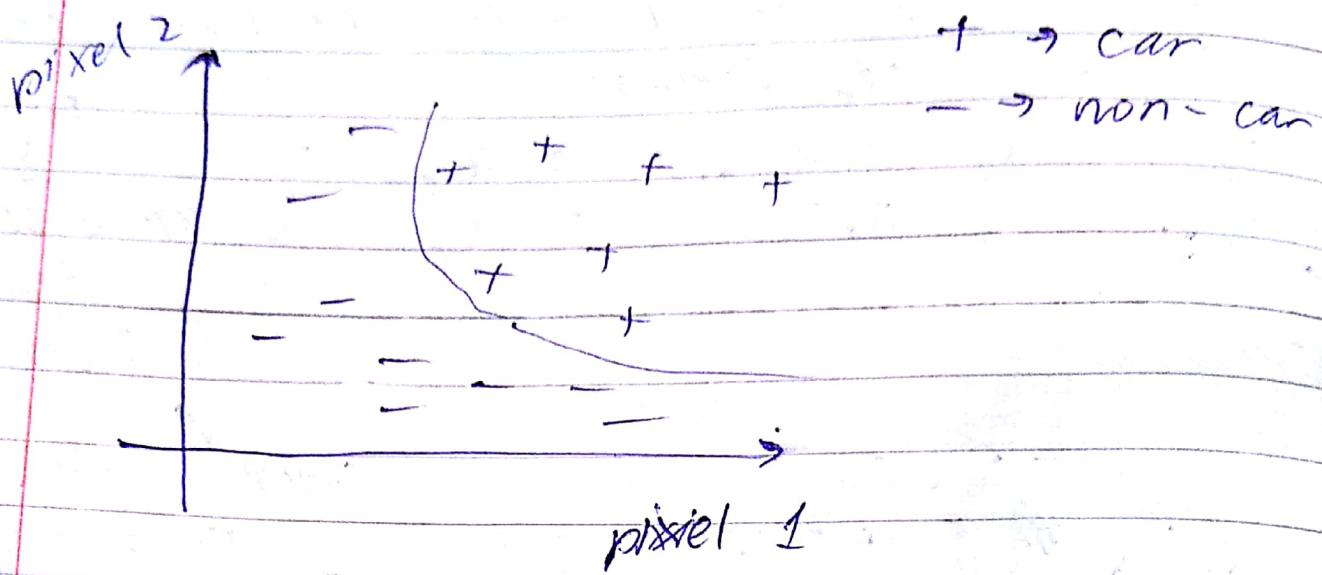
Counting the third order features from the data set fields:

$$\left. \begin{array}{c} x_1 x_2 x_3, x_1^2 x_2, \dots \end{array} \right\} 170\,000 \text{ features}$$

$[O(n^3)]$  for  $n=100$

~~Too many~~ Feature space is blown up when  $n$  increases as data is complex (as the complexity increases, the ~~more~~ higher the order needed becomes)

Suppose computer vision problem where we are trying to understand whether a picture is a car or not



~~Each~~ Each pixel constitutes a feature  
 So if  $50 \times 50$  pixel image then 2500 pixels ~~total~~/features in total  
 Consider that each pixel is a value of intensity - not RGB)

Considering the Quadratic features of a 2500 pixel image  $(x_i \times x_j)$  a 3 million features

→ Thus neural networks

complex non-linear hypothesis even when  $n$  is large

## Neural networks

Main goal: have machines that can mimic the working of a brain

Origins: Algorithms that try to mimic the brain

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

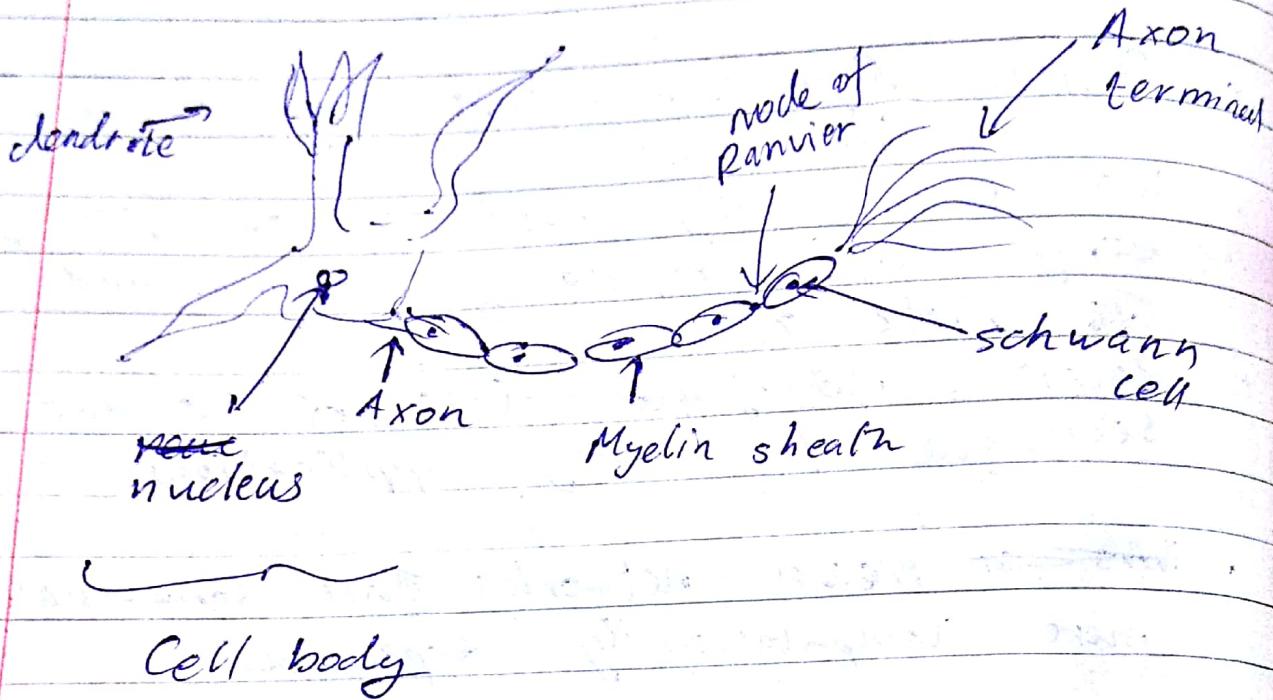
Recent resurgence: state-of-the-art technique for many applications

~~However~~ Neural networks are somewhat more computationally expensive

The "one learning algorithm" hypothesis - brain does all of its tasks (hearing, process sense of touch, consciousness, reasoning) through a single learning algorithm as opposed to many different algorithms running in unison

## Model representation

Neuron in the brain:



Dendrite - "input wires" - receive inputs  
from other locations

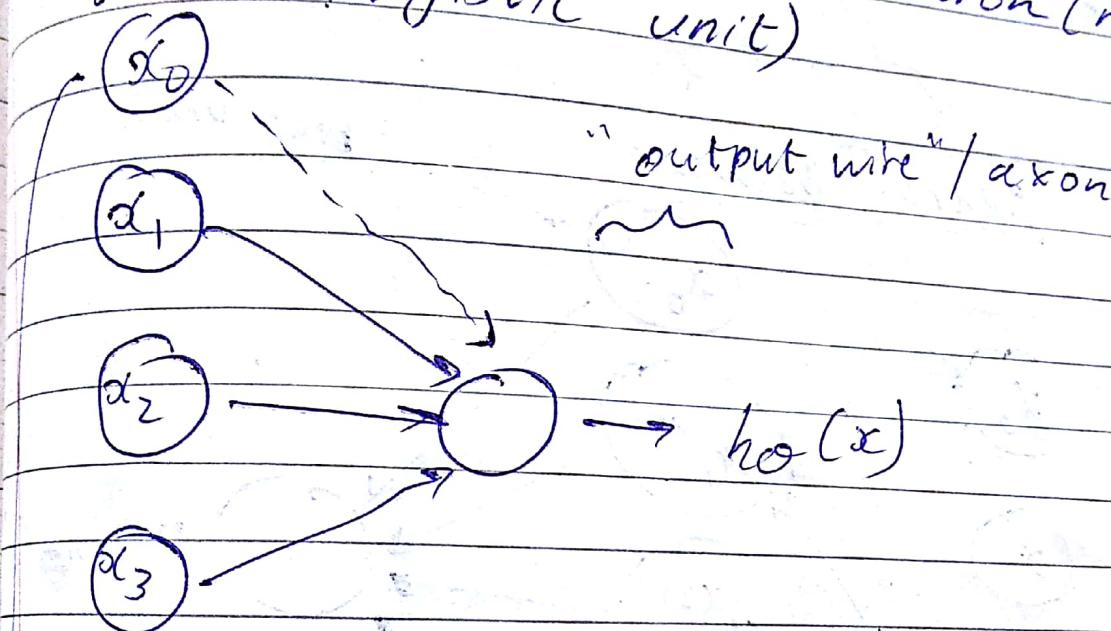
Axon - "output wire" - send signals to  
other neurons

Neuron  
does  
actual  
computation

- Neuron - a computational unit that gets some inputs & sends output through axon
- Neurons communicate with one another via small pulses of electricity called spikes.

through these spikes, muscles contract, thought occurs  
eyes blinks; all human thought occurs

~~Artificial Neural Network (ANN)~~  
is a simple model of a neuron (modeled as a logistic unit)



$$h_0(x) = \frac{1}{1+e^{-\theta^T x}}$$

where:

dendrites body of  
or input neuron

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}, \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

At times this is drawn  
as well (called the bias unit / bias  
neuron)

$$\theta_0 = 1$$

"weights"  
of model

Sigmoid (logistic) activation function:

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

$L_1$  used in hypothesis  
of neural networks

In neural networks the parameters are called "weights" as well.

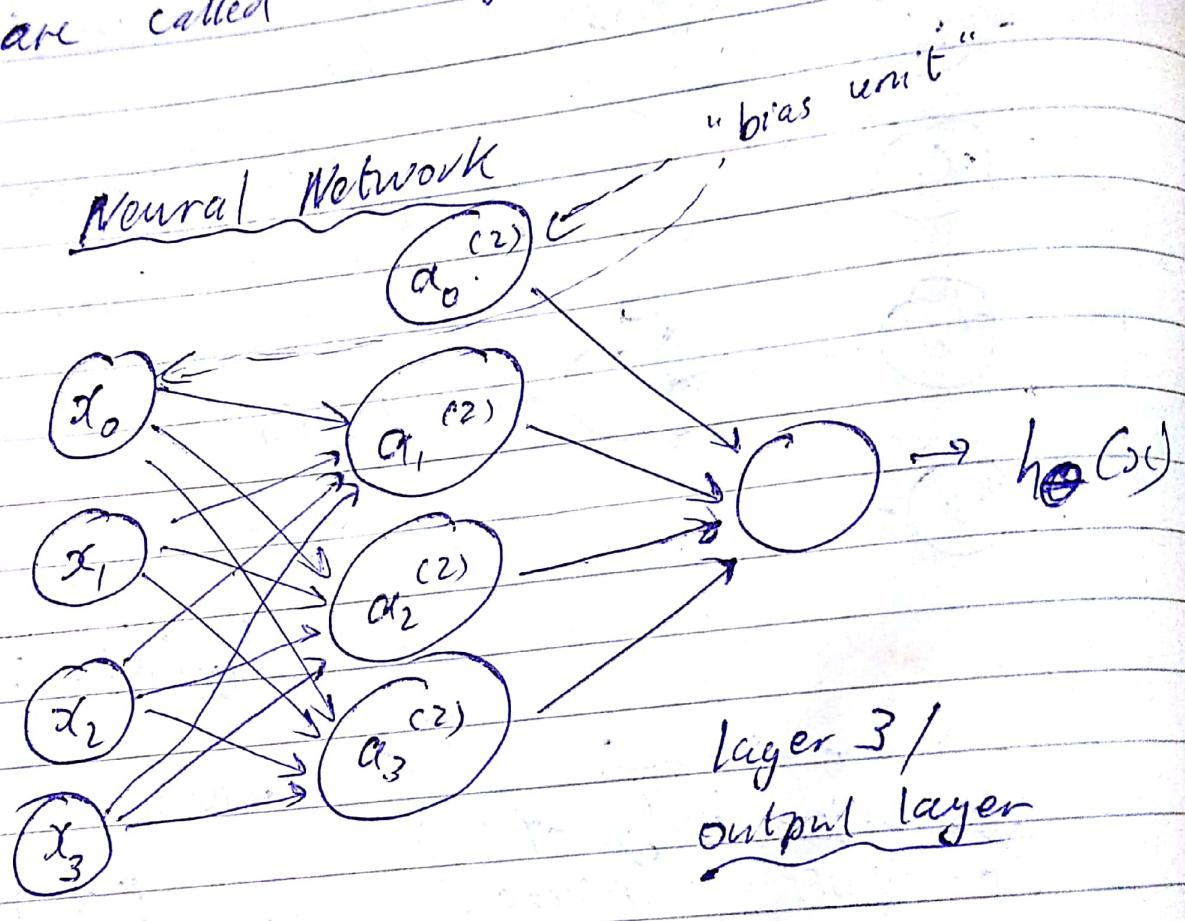
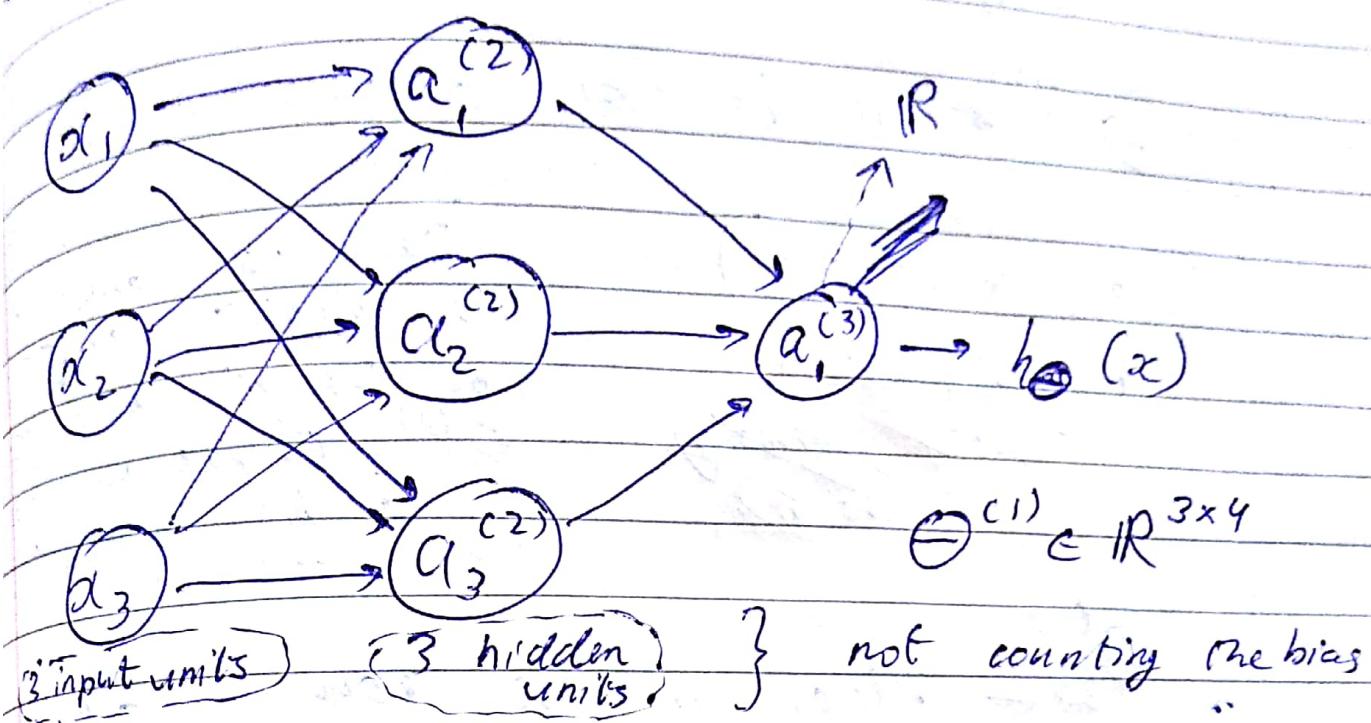


Diagram illustrating a neural network structure:

- Layer 1 / input layer**
- Layer 2 / hidden layer**
- values while training** (indicated by a bracket under Layer 2)

Hidden layer called like this may because you can see the input values going into the model & output coming out but you can't see the values in one ~~the~~ hidden layer

ANNs can have multiple hidden layers



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

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

sigmoid activation function

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

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

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

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

If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then

$\Theta^{(j)}$  will be of dimension

$$s_{j+1} \times (s_j + 1)$$

~~not counting  
the bias~~

This +1 comes from  
the addition in  $\Theta^{(j)}$   
of the bias nodes  $\theta_0$  &  
 $\theta_0^{(j)}$ . Output will not  
include bias nodes while  
inputs will.

$$\Theta^{(j)}$$

$\theta_{30}$  ← indices indicate row & column

$$\Theta^{(1)} = \begin{bmatrix} \theta_{10} & \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{20} & \theta_{21} & \theta_{22} & \theta_{23} \\ \theta_{30} & \theta_{31} & \theta_{32} & \theta_{33} \end{bmatrix}$$

$$\Theta^{(2)} = [\theta_{10} \quad \theta_{11} \quad \theta_{12} \quad \theta_{13}]$$

## Forward propagation: Vectorised implementation

$$z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3$$

$$z_2^{(2)} = \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3$$

$$z_3^{(2)} = \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3$$

Hence:

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

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

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

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} \quad z^{(2)} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix}$$

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

$$Z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} = \begin{bmatrix} \Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 \\ \Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3 \\ \Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3 \end{bmatrix}$$

$$= \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} & \Theta_{13} \\ \Theta_{20} & \Theta_{21} & \Theta_{22} & \Theta_{23} \\ \Theta_{30} & \Theta_{31} & \Theta_{32} & \Theta_{33} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \Theta^{(1)} x$$

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

$$a^{(2)} = g(\alpha Z^{(2)})$$

↙  $\mathbb{R}^3$

$\alpha \mathbb{R}^3$

The inputs to the neural network can be expressed as  $x$  or  $a^{(1)}$

Each layer has ~~one~~ a bias node:  $a_0^{(j)}$

Each layer has a bias node:  $a_0^{(j)} = 1$

Hence when evaluating  $Z^{(3)}$  we must add  $a_0^{(2)} = 1 \rightarrow a^{(2)} \in \mathbb{R}^4$

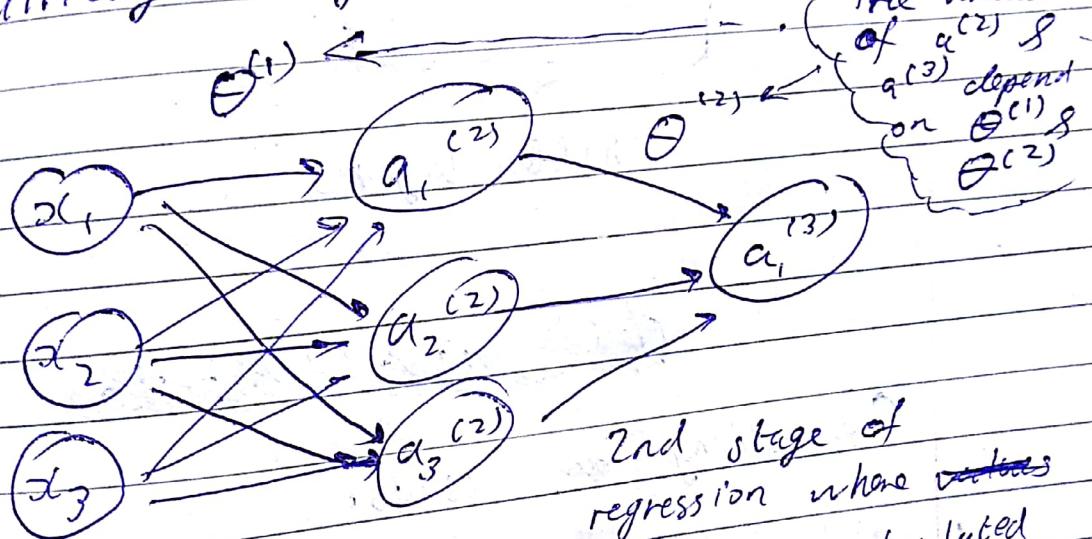
~~sk~~

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

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

Process of computing  $h_{\theta}(x)$  is called forward propagation because the inputs are used propagated forward till the output & activations are evaluated at each stage

Neural networks can be seen as a multistage logistic regression

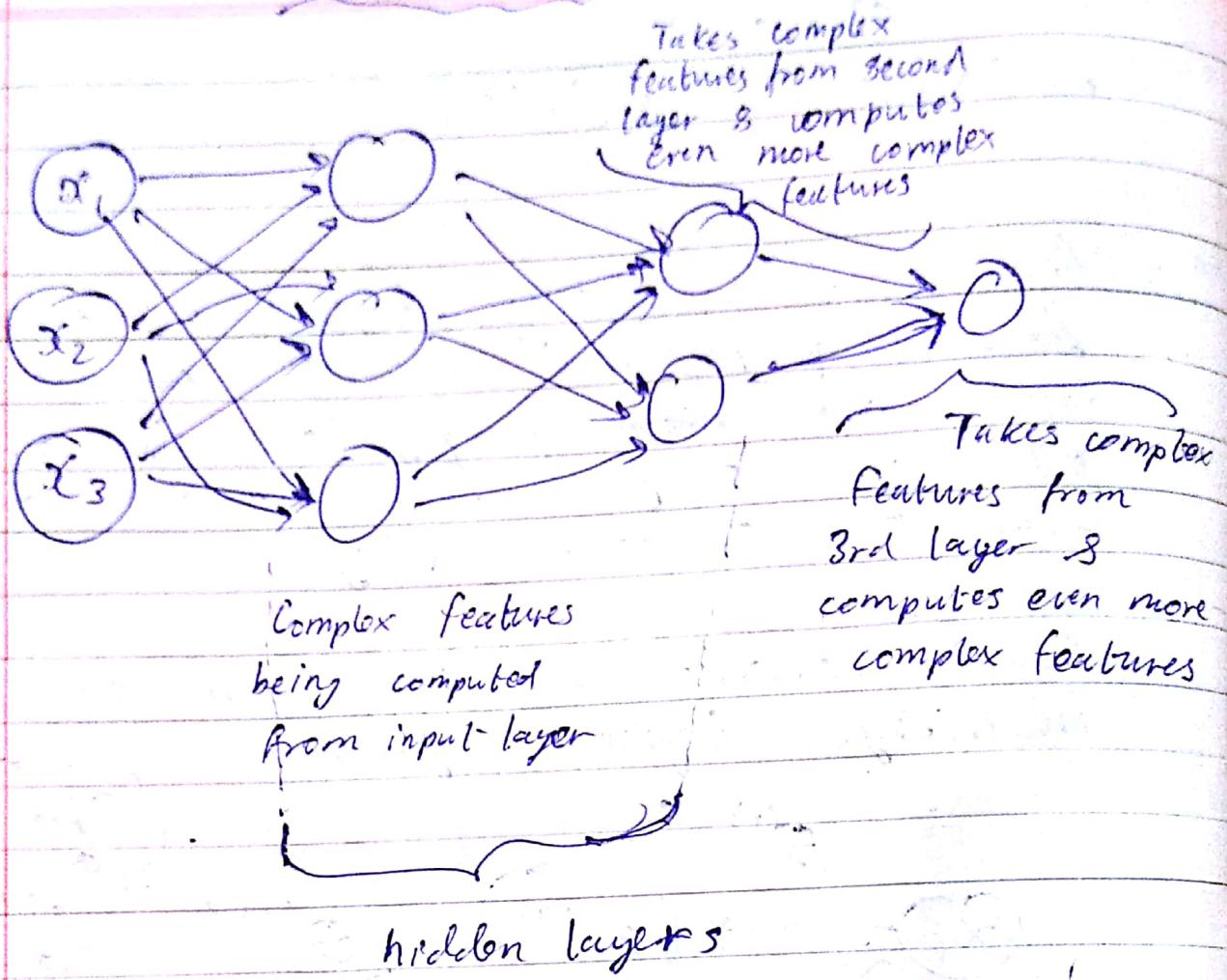


1st stage  
of regression  
where values  
 $a^{(2)}$  are calculated

2nd stage of  
regression where values  
 $a^{(3)}$  is calculated  
from  $a^{(2)}$  as inputs

~~Net~~ Neural Networks got to learn their  
own features:  $a^{(2)}$  in this example

## Other network architectures



## Neural Network Derivation

Generalising our results:

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

$$\alpha = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_l^{(j)} \end{bmatrix}$$

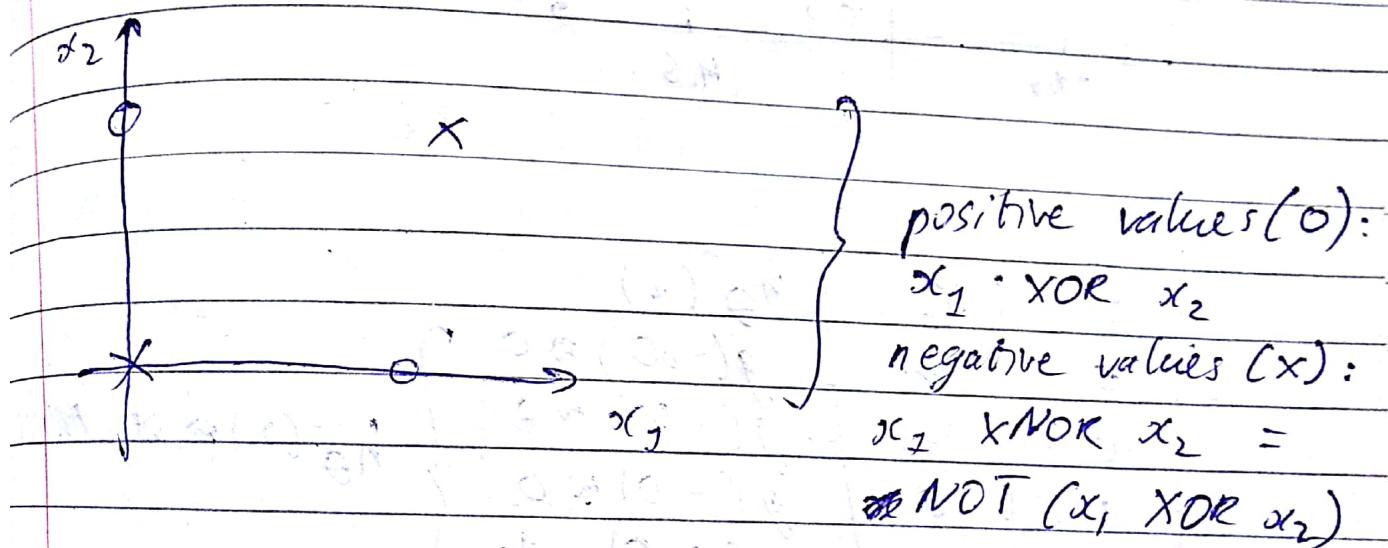
$l$ : # of nodes in layer  $j$

Given that  $z^{(j)} = \Theta^{(j-1)} a^{(j-1)}$

Given that  $x = a^{(1)}$  then:  
 $z^{(j)} = \Theta^{(j-1)} a^{(j-1)}$

## Neural Network Intuition

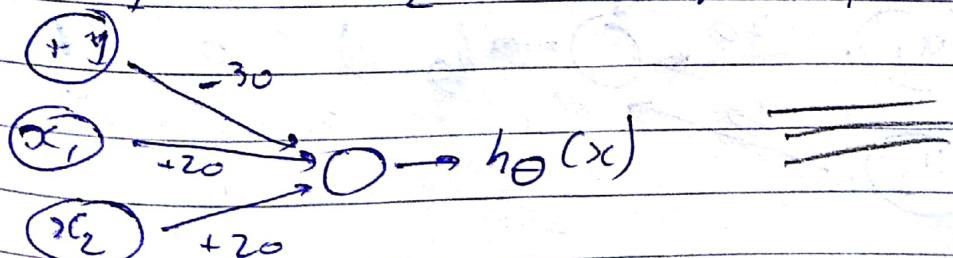
Non-linear classification example : XOR / XNOR  
 $x_1, x_2$  ~~are~~ are binary (0 or 1)



## Simpler example (AND)

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

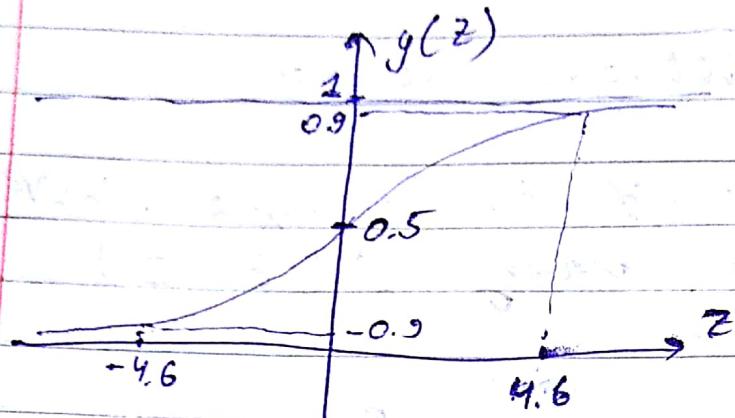
$y = x_1 \text{ AND } x_2 \leftarrow \text{target function}$



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

$\uparrow \quad \uparrow \quad \uparrow$

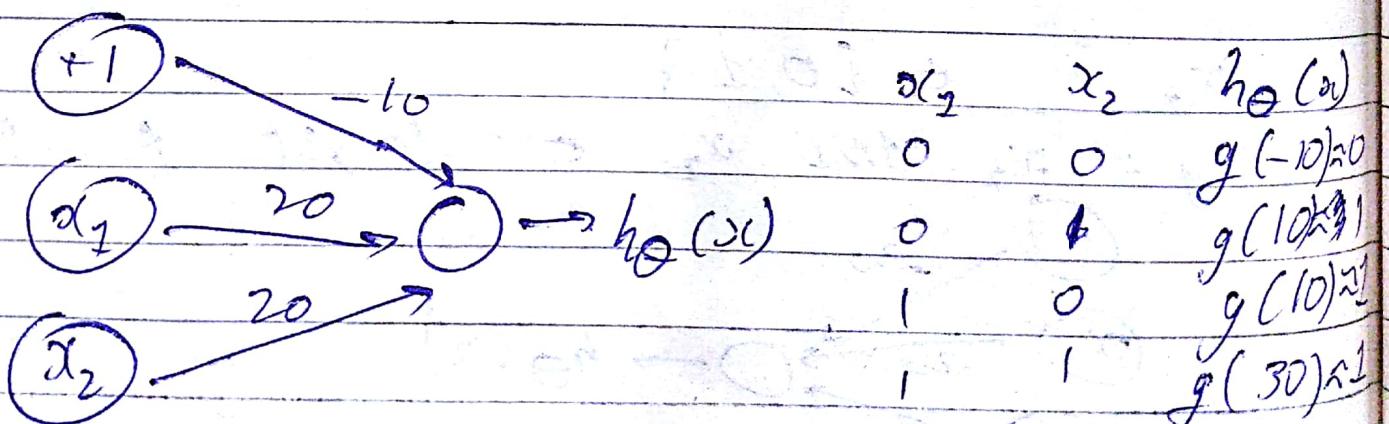
$\theta_{10}^{(1)}$      $\theta_{11}^{(1)}$      $\theta_{12}^{(1)}$



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

$\left. \begin{array}{l} h_{\theta}(x) \approx x_1 \text{ AND } x_2 \\ h_{\theta}(x) \approx x_1 \text{ OR } x_2 \end{array} \right\}$

Simple example: OR



Simple example : AND  
in mathematical terms

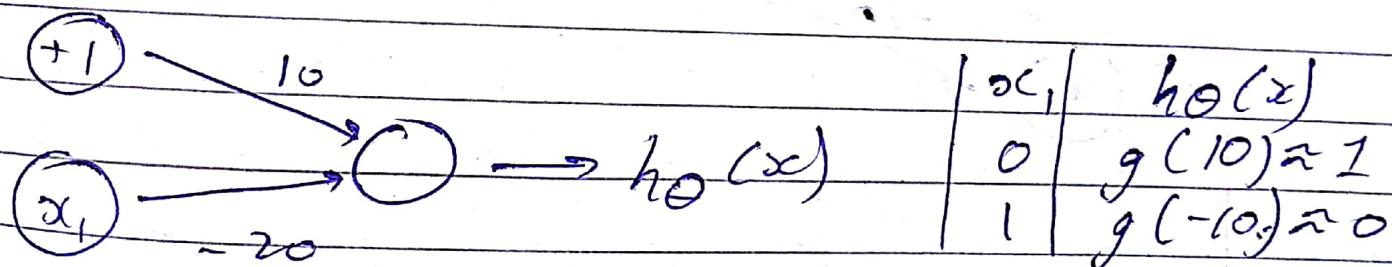
$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow [g(z^{(2)})] \rightarrow h_{\Theta}(x)$$

Graph of function

$$\Theta^{(1)} = [-30 \ 20 \ 20]$$

Simple example: NOT

Negation:  $\text{NOT}$



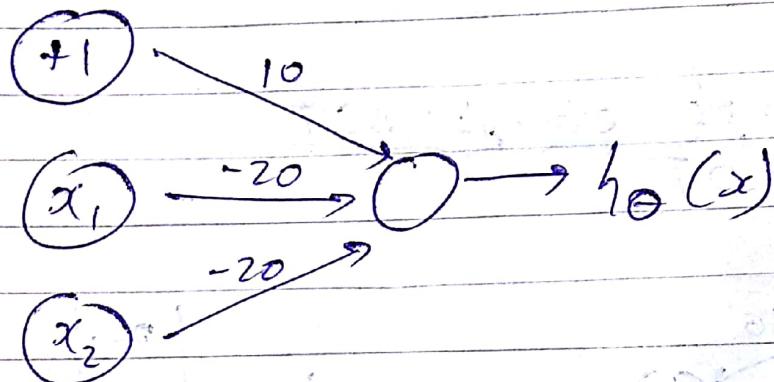
$$h_{\Theta}(x) \approx g(10 + 20x_1)$$

Simple example: try it yourself

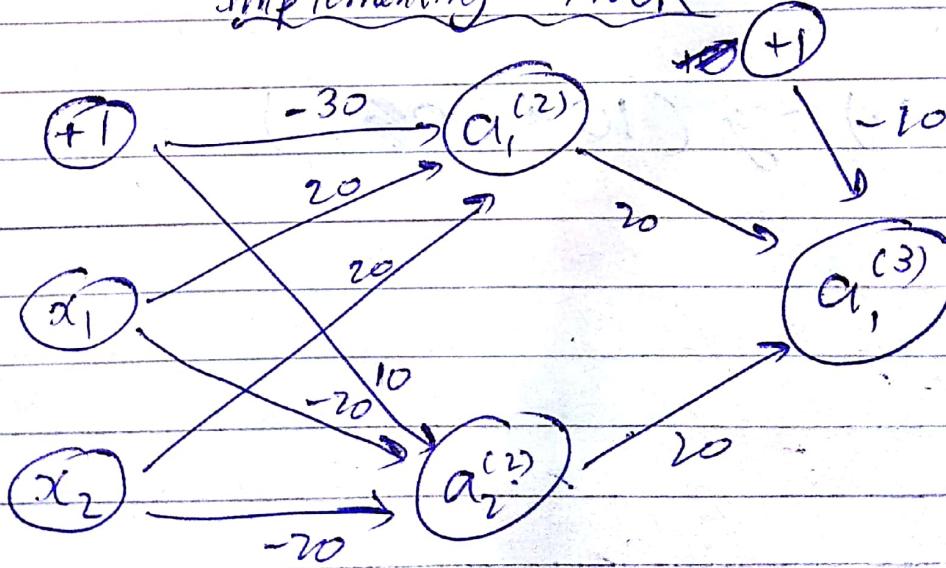
Implement this function:  $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

$x_1$	$x_2$	$f$
0	0	1
0	1	0
1	0	0
1	1	0

$$h_{\Theta}(x) = 10 - 20x_1 - 20x_2$$



Implementing XNOR



$x_1$  AND  $x_2$

~~AND form~~

$a_1^{(2)}$  OR  $a_2^{(2)}$



$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_0(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1



$(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$  or  $x_1 \text{ NOR } x_2$

~~form~~

In mathematical terms:

AND :

$$\Theta^{(1)} = [-30 \quad 20 \quad 20]$$

NOR :

$$\Theta^{(1)} = [10 \quad -20 \quad -20]$$

OR :

$$\Theta^{(1)} = [-10 \quad 20 \quad 20]$$

XNOR :

$$\begin{bmatrix} \Sigma_0 \\ \Sigma_1 \\ \Sigma_2 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \end{bmatrix} \rightarrow [a^{(3)}] \rightarrow h_0(x)$$

where :  $\Theta^{(1)} = [-30 \quad 20 \quad 20]$  &  $\Theta^{(2)} = [-10, 20, 20]$

$$\begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix}$$

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

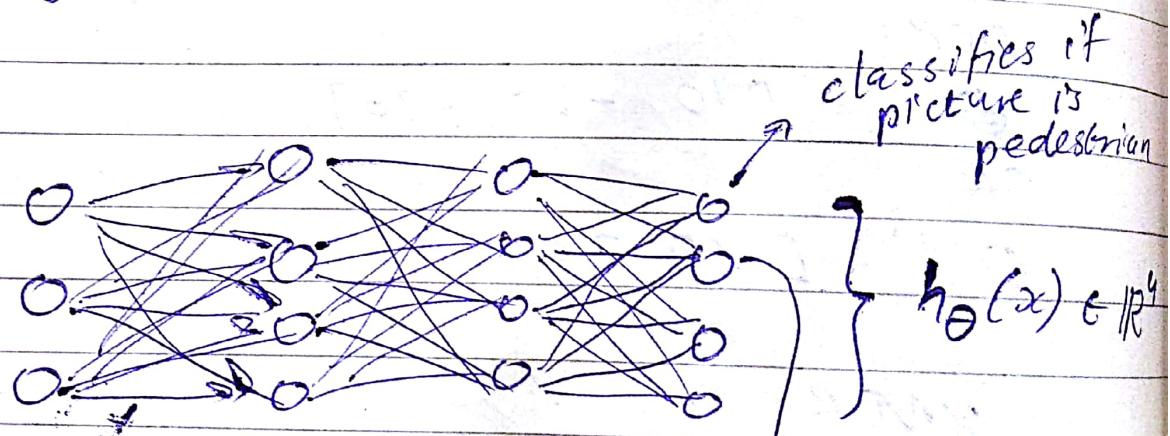
$$a^{(3)} = g(\Theta^{(2)}a^{(2)})$$

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

## Multi-class classification with Neural Networks

The approach used is an extension of one-vs-all method

Example problem: classify pedestrian, car, motorcycle & truck



$$h_{\theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ when pedestrian}$$

classifies if picture is car and etc

$$h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ when picture is car}$$

$$h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \text{ when picture is motorcycle}$$

$$h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \text{ when picture is truck}$$

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

$y^{(i)}$  is one of :  $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

pedestrian      car      motorcycle      truck

$$h_{\theta}(x^{(i)}) \approx y^{(i)}$$

$\hookrightarrow \mathbb{R}^4$

Mathematically  
 the above  
 problem can  
 be expressed as  
 such:

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_n^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ \vdots \\ a_n^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(3)} \\ a_1^{(3)} \\ a_2^{(3)} \\ \vdots \\ a_n^{(3)} \end{bmatrix} \rightarrow \begin{bmatrix} h_{\theta}(x)_1 \\ h_{\theta}(x)_2 \\ h_{\theta}(x)_3 \\ h_{\theta}(x)_4 \end{bmatrix}$$