

## Week 4 - Notes

### Neural - networks representation: Non-linear Hypothesis

Consider a supervised learning classification problem with training set as follows:

#### Non-linear classification



possibly apply logistic regression with heaps of non-linear features such as:

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

sigmoid function.

Logistic Regression is o.k for say ~ two features. But what happens if we have more?

$x_1$  = size

$x_2$  = # bedrooms

$x_3$  = # floors

$x_4$  = age.

...

100

$\left. \begin{array}{l} h=100 \\ h=100 \end{array} \right\}$

$x_1^2, x_1 x_2, x_1 x_3, \boxed{x_1 x_4}, \dots, x_1 x_{100}$   
 $x_2^2, x_2 x_3$

"called second order terms"

collecting all second order terms for the case when  $h=100$  results in

~ 5000 features.  $\left( \approx \frac{n^2}{2} \right) = O(n^2)$

- could overfit  $\mathcal{J}$  training set

- and be computationally expensive.

one thing you could do is say include:

$\rightarrow x_1^2, x_2^2, x_3^2, \dots, x_{100}^2$  so only  $\approx 100$

If you were to include 3rd order polynomial features

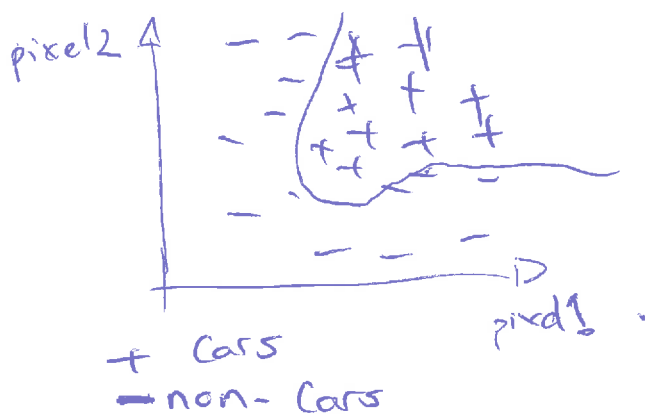
e.g.  $x_1 x_2 x_3, x_1^2 x_2, x_{10}, x_{11}, x_{12}, \dots$

$O(n^3) \approx 170,000$ .

## Week 4 - notes

### neural-networks representation: non-linear hypothesis

For ~~too~~ many machine learning problems  $n$  will be large.



let's say:

50x50 pixel images  $\rightarrow$  2500 pixels  
 $n = 2500$  (7500 in RGB)

$$X = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix} \quad 0-255$$

All quadratic features  $(x_i \times x_j)$   $\approx$  3 million features

$\therefore$  Logistic regression is not a good way to extract information when  $n$  is large.

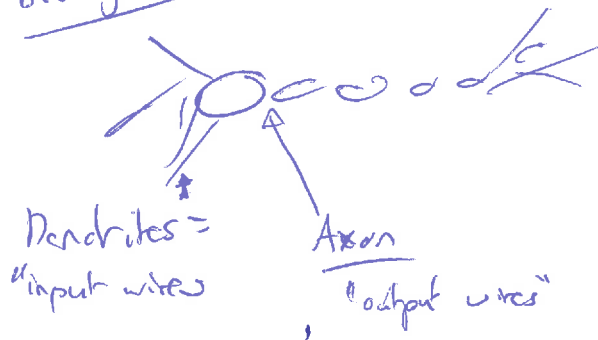
# Week 4 - notes

## Neurons and the brain

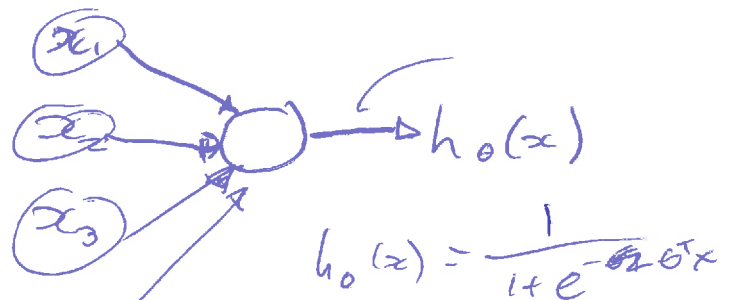
Origins: Algorithms that try to mimic the brain.

## Neural Network Representation I

biological model



Neuron model: Logistic unit.



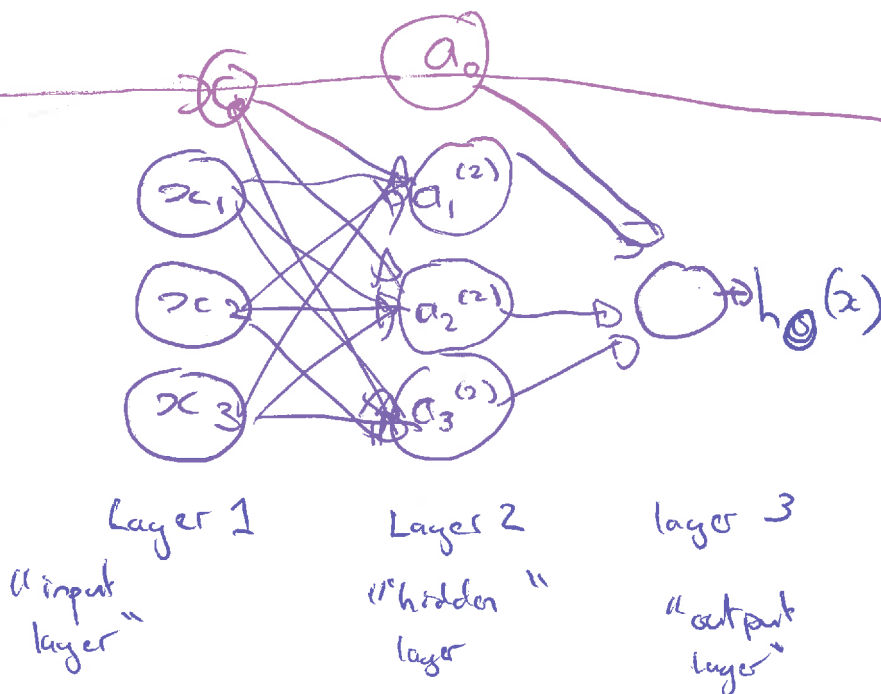
bias unit.

Always  
( $x_0 = 1$ )

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

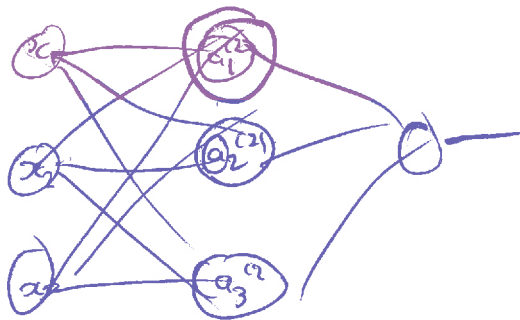
Sigmoid (logistic) activation function

Also known as "weights" or parameters.



# Week 4 - notes

## Neural Network Representation I



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

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

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

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

$a_3^{(2)} = g(\text{~~~~~})$

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

→ If a network has  $s_j$  units in layer  $j$ ,  ~~$s_{j+1}$~~   $s_{j+1}$  units in

layer  $j+1$ , then  $\Theta^{(j)}$  is

will be of dimension  $s_{j+1} \times (s_j + 1)$

eg

subscript...

$2 \times$

input layer 1  
layer 2  
 $s_1 = 2$   $s_2 = 4$

$\therefore \Theta = 4 \times 3$   
↑  
units of

## Week 4 - Notes

### Neural Networks representation - Model representation $\bar{A}$

$$a_1^{(2)} = g(\text{~~~~~}) \rightarrow z_1^{(2)}$$

$$a_2^{(2)} = g(\text{~~~~~}) \rightarrow z_2^{(2)}$$

$$a_3^{(2)} = g(\text{~~~~~}) \rightarrow z_3^{(2)}$$

$$\therefore a_2^{(1)} = g(z_{a2}^{(2)})$$

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

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \theta^{(1)} x \quad \text{step 1} \quad (*)$$

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

$\mathbb{R}^3$

$\mathbb{R}^3$  (3 dimensional vector)

previous output from layer 2.

add bias step 3

then to Add the bias neuron

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

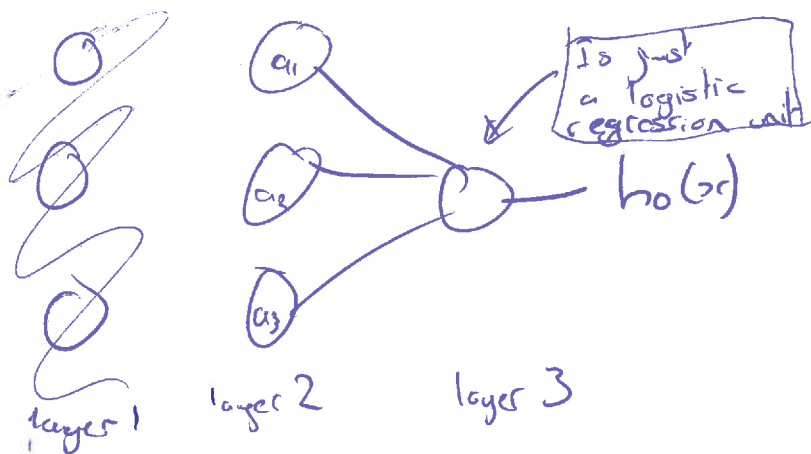
repeat  
of step  
1 & 2  
(from previous) (\*)

$$z^{(3)} = \theta^{(2)} a^{(2)}$$
$$h_0(x) = a^{(3)} = g(z^{(3)})$$

## week 4 notes

### Neural Network representation II

- Neural network learning its own features.



$$h_0(x) = g(\theta_0^{(2)} a_0^{(2)} + \theta_1^{(2)} a_1^{(2)} + \theta_2^{(2)} a_2^{(2)} + \theta_3^{(2)} a_3^{(2)})$$

\* just the same as logistical regression but instead of using the features we are ~~are~~ now using

features  $x_1, x_2, x_3$   
 $a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$

the new instead of feeding the features  $x_1, x_2, x_3$  into logistical regression, the network gets to learn its own parameters  $a_1, a_2, a_3$  to do it.

other network architecture

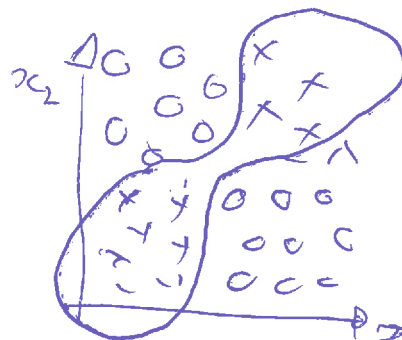
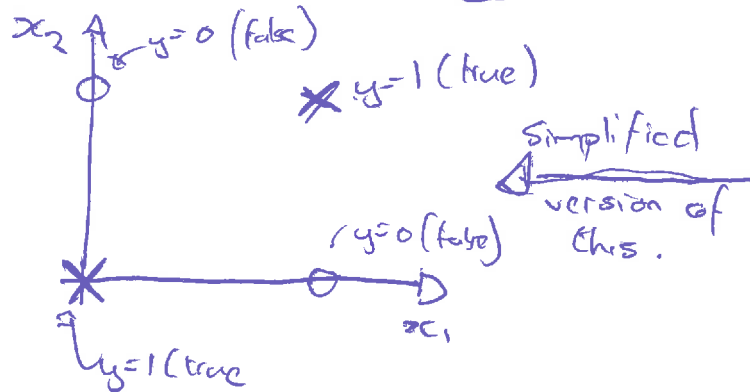


# Week 4 - Notes.

## Neural Networks representation - examples & intuitions I

Non-linear classification example: XOR/XNOR.

$\rightarrow x_1, x_2$  are binary (0 or 1)

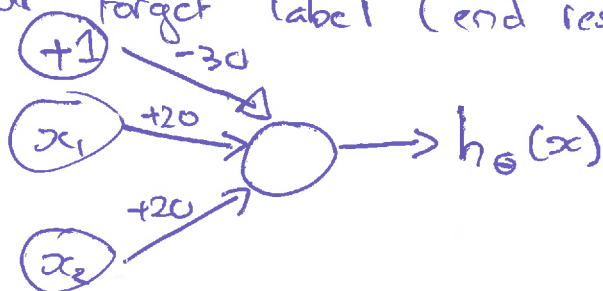


what we would like to do is learn a non-linear boundary, that separates the positive and the negative examples.

## Simple example "AND"

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

our target label (end result) is:  $y = x_1 \text{ AND } x_2$



$$-10 + 20x_1 + 20x_2$$

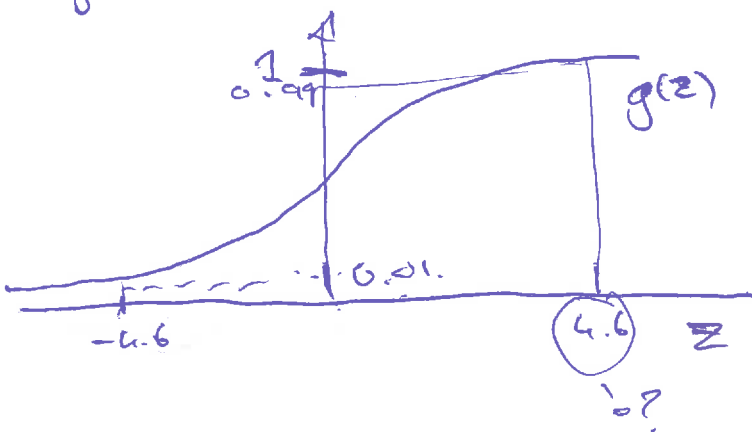
FOR  $\uparrow$

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

$\uparrow$   $\uparrow$   $\uparrow$   
 $\theta_{10}$   $\theta_{11}$   $\theta_{12}$

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

Sigmoid function:



AND function

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

# Week 4 - Notes

## Neural networks representation - examples and intuitions II

negation | NOT  $x_1$



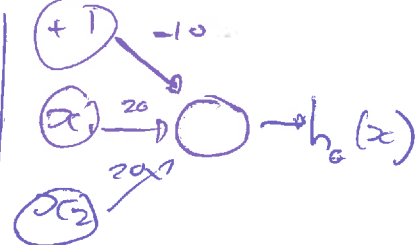
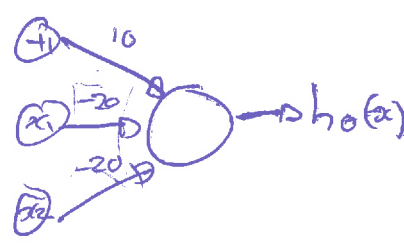
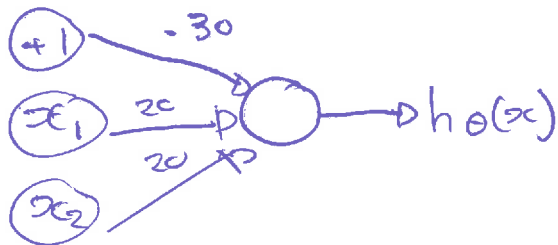
$x_1$	$h_\theta(x)$
0	$g(10) = 1$
1	$g(-10) = 0$

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

(NOT  $x_1$ ) AND (NOT  $x_2$ )

(= 1 if and only if  $x_1 = x_2 = 0$ )

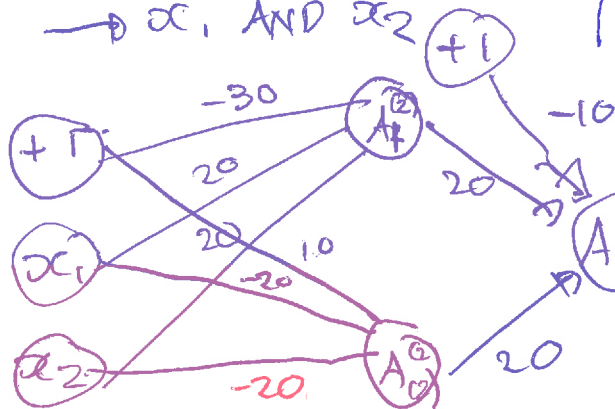
Putting it together:  $x_1$  XOR  $x_2 \rightarrow$  looks like



$\rightarrow x_1$  AND  $x_2$

(NOT  $x_1$ ) AND (NOT  $x_2$ )

$\rightarrow x_1$  OR  $x_2$



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



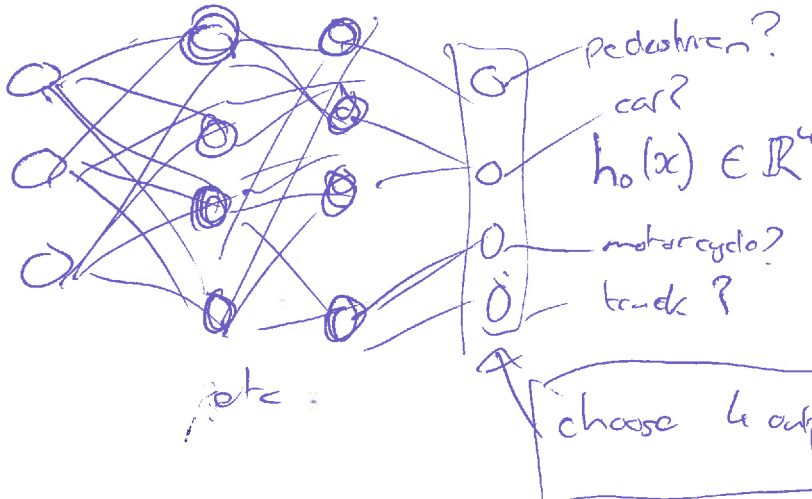
# Week 6 - Notes

## Multiclass classification

### Neural networks representation - ~~softmax~~ ~~neural networks~~

example: images of

- pedestrian
- car
- motorcycle
- truck



eg. we want

$$h_0(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{pedestrian}$$

$$h_0(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{car}$$

$$h_0(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{motorcycle}$$

etc

training set:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(n)}, y^{(n)})$

$y^{(i)}$  one of

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

pedestrian car motorcycle truck

previously

$y = \{1, 2, 3, 4\}$   
now representing  
binarily...