# 8.   NEURAL NETWORKS – NON LINEAR HYPOTHESIS

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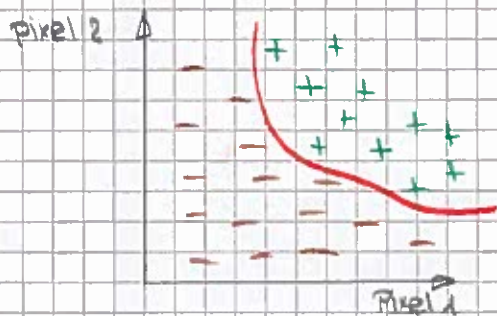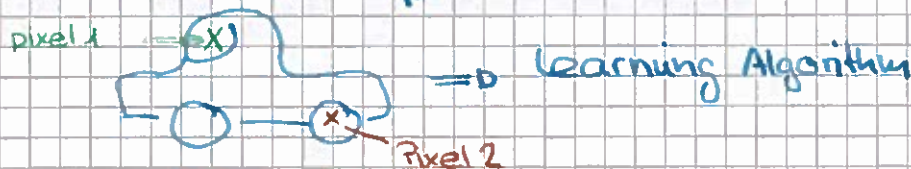$$+ \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots )$$

by using only second order terms
$x_1^2, x_1 x_2 \dots \quad x_{100}$

$x_1$ = size
$x_2$ = # bedrooms
$x_3$ = # floors
$x_4$ = age
$\vdots$
$x_{100}$

$n=100$

Including all the quadratic features
would mean to have a 5000 features
$\Rightarrow$ grows $O(n^2) \Rightarrow \dfrac{n^2}{2}$  $\Rightarrow$ could lead to overfitting

Just using the $x_1^2, x_2^2, x_3^2 \dots x_{100}^2$
does not allow to create a function as
more complex functions, would not fit the complex dataset

while using cubic features it would $O(n^3)$
be about 170'000 features for
$n = 100$ features

Why is non linear hypothesis relevant?



pixel 1 $\longrightarrow$ (x)
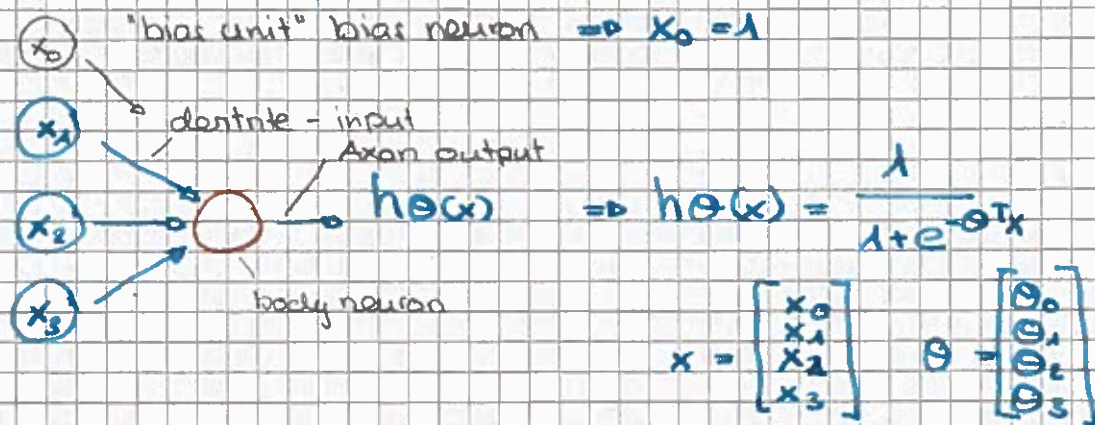
$\Rightarrow$ Learning Algorithm

Pixel 2



+ Cars
− Non Cars

$\Rightarrow$ 50 x 50 pixels images $\Rightarrow$ 2500 pixels
$n = 2500$  (7500 RGB)

$$x = \begin{bmatrix} \text{pixel 1 intensity} \\ " \quad 2 \quad " \\ " \quad 3 \quad " \\ \vdots \\ \text{pixel n intensity} \end{bmatrix} \begin{matrix} \text{Value } 0-255 \\ \\ \\ \\ \\ 2500 \end{matrix}$$

by using all quadratic features $(x_i \times x_j)$
$\approx$ 3 million features per training example
is computational expensive
Not good for algorith. with lot of features
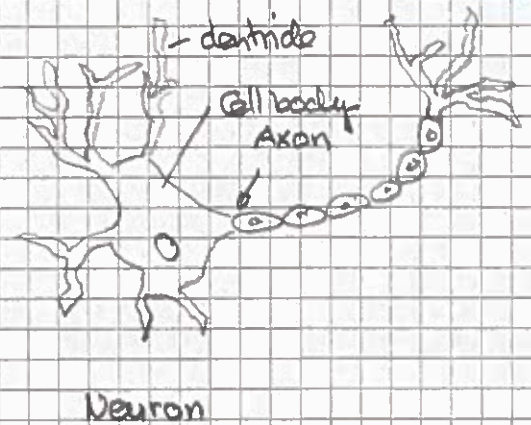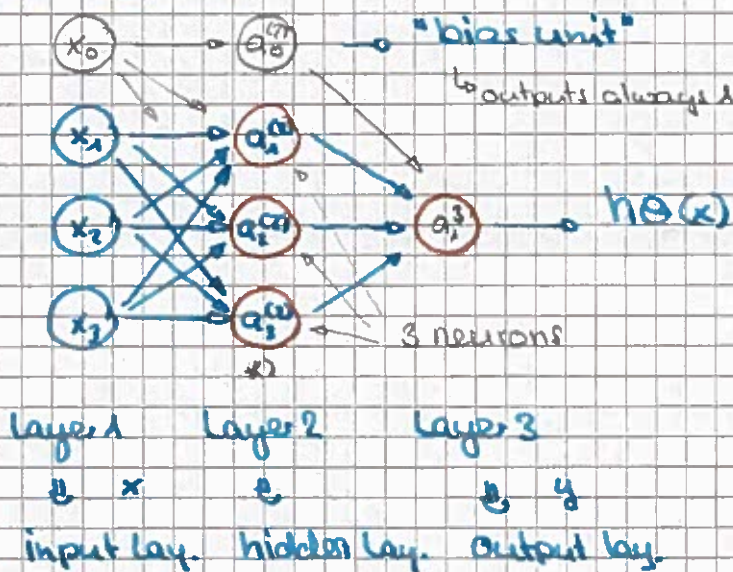
## Neuron model: Logistic unit



"bias unit" bias neuron $\Rightarrow x_0 = 1$

dentrite - input

Axon output

$$h_\Theta(x) \quad \Rightarrow \quad h_\Theta(x) = \frac{1}{1+e^{-\Theta^T x}}$$

body neuron

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad \Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \Theta_3 \end{bmatrix}$$

Sometimes called "weights" - parameter of model

## Sigmoid (logistic) activation function.

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

## Neural Network



"bias unit"

↳ outputs always 1

$h_\Theta(x)$

3 neurons

*)

- dentride
- Cell body
- Axon

Neuron

Layer 1     Layer 2     Layer 3

   x        e        y

input lay.   hidden lay.   output lay.

\*) you don't observe the values processed in the hidden layer

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

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

activation of first unit in layer 2



3 input units     3 hidden units     $\Theta^{(1)} \in \mathbb{R}^{3\times4}$ dim. Matrix

a "activation" is the value which is computed and output by the node

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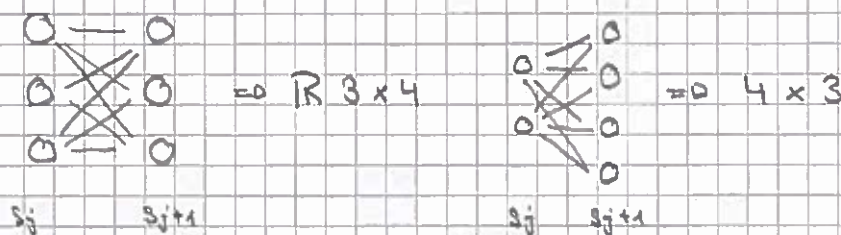
sigmoid

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{13}^{(1)} x_3\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3\right)$$
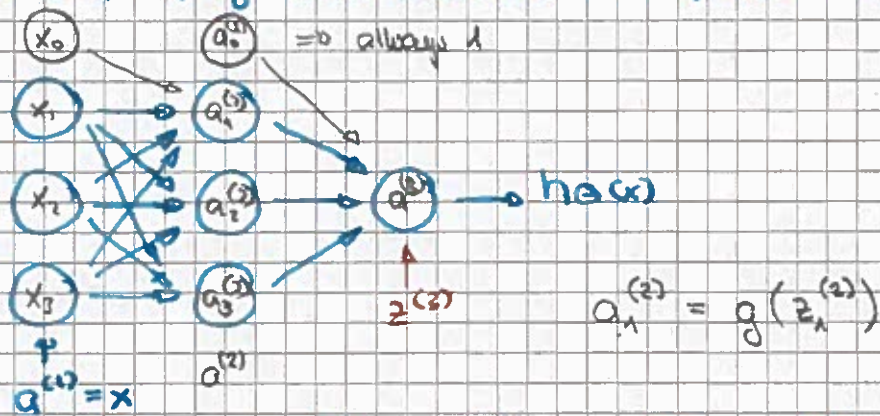
$$h_\Theta(x) = a_1^{(2)} = g\left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}\right)$$

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

Example =o



=o $\mathbb{R}\, 3\times4$

$s_j$    $s_{j+1}$



=o $4\times3$

$s_j$    $s_{j+1}$

### Forward propagation: vectorized implementation



$a^{(1)} = x$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3\right) \to z_1^{(2)}$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3\right) \to z_2^{(2)}$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3\right) \to z_3^{(2)}$$

$$h_\Theta(x) = g\left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}\right) \quad z^{(3)}$$

$a_1^{(2)} = g\left(z_1^{(2)}\right)$

$a_3^{(2)} = g\left(z_2^{(2)}\right)$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \\ \cancel{z_4^{(2)}} \end{bmatrix}$$

**Vectorized implementation**

$$z^{(2)} = \Theta^{(1)} x \quad \Rightarrow \quad \Theta^{(1)} \cdot a^{(1)}$$

$$a^{(2)} = g\left(z^{(2)}\right)$$
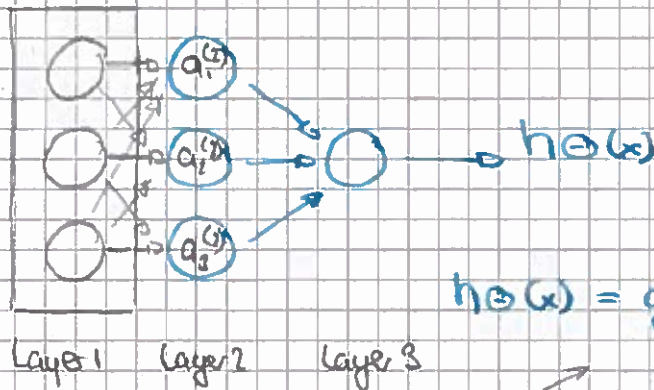
$\underbrace{\quad}_{\mathbb{R}^3} \qquad \underbrace{\quad}_{\mathbb{R}^3} \Rightarrow 3 \text{ dim. vector}$

Add $a_0^{(2)} = 1 \quad \Rightarrow a^{(2)} \in \mathbb{R}^4$
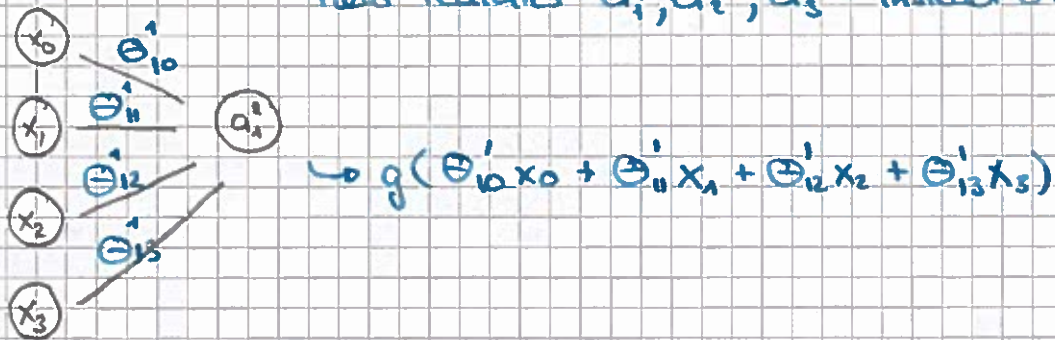
$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$h_\Theta(x) = a^{(3)} = g\left(z^{(3)}\right)$$

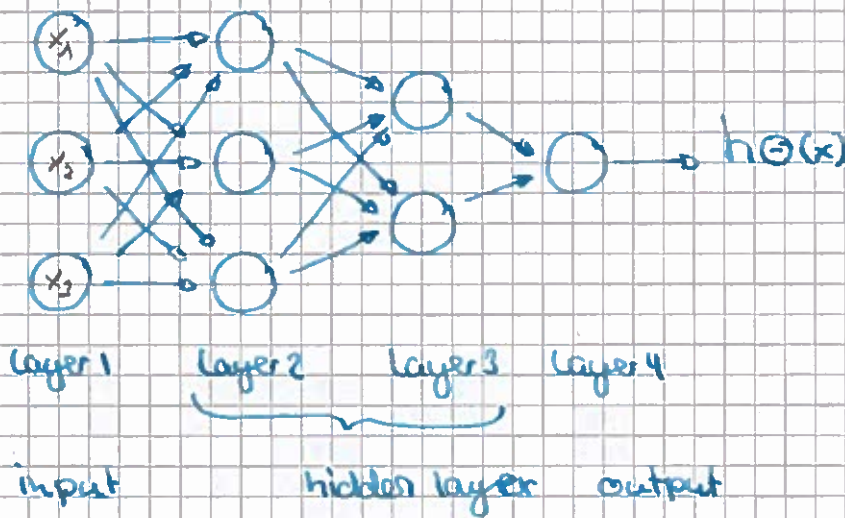Neural Network learning its own features.



Layer 1    Layer 2    Layer 3

$$h_\Theta(x) = g\left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}\right)$$
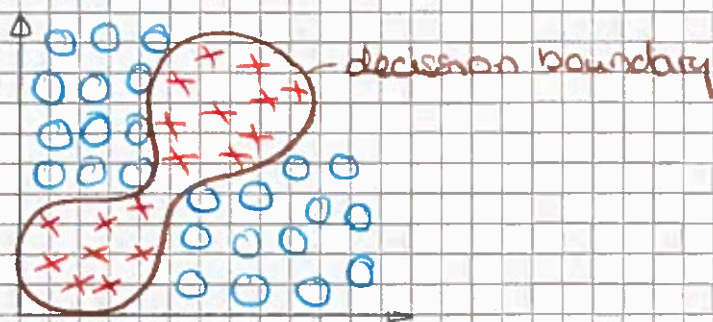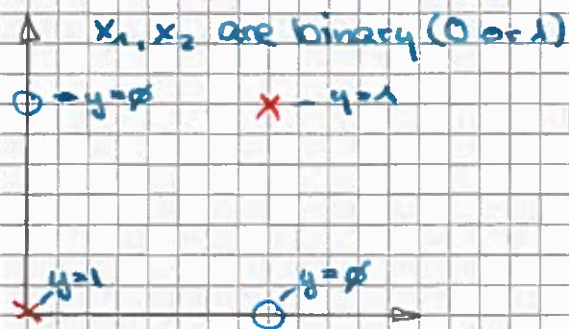
it's just logistic regression except using the new features $a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$ instead of $x_1, x_2, x_3$



$$\hookrightarrow g\left(\Theta_{10}^{1} x_0 + \Theta_{11}^{1} x_1 + \Theta_{12}^{1} x_2 + \Theta_{13}^{1} x_3\right)$$

Other network architectures



Layer 1    Layer 2    Layer 3    Layer 4

input    hidden layer    output

## Non linear classification example: XOR\XNOR

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

$\bigcirc = y = \emptyset$      $\times - y = 1$

$y = 1$ (×)      $\bigcirc y = \emptyset$

decision boundary

XOR

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$y = x_1$ XOR $x_2$

$x_1$ XNOR $x_2$

NOT $(x_1$ XOR $x_2)$

### Simple example AND:

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

$y = x_1$ AND $x_2$

$(+1) \xrightarrow{-30}$

$(x_1) \xrightarrow{+20} \bigcirc \longrightarrow h_\Theta(x)$

$(x_2) \xrightarrow{+20}$

OR:

| $x_1$ | $x_2$ | $h_\Theta(x)$ |
|---|---|---|
| 0 | 0 | $g(-10) \approx 0$ |
| 0 | 1 | $g(+10) \approx 1$ |
| 1 | 0 | $g(+10) \approx 1$ |
| 1 | 1 | $g(30) \approx 1$ |

$(x_0) \xrightarrow{-10}$
$(x_1) \xrightarrow{+20} \bigcirc$
$(x_2) \xrightarrow{+10}$

$g(-10 \cdot 1 + 20 \cdot \emptyset + 20 \cdot 0) = g(-10)$

$\Rightarrow h_\Theta(x) = g(-30 + 20 x_1 + 20 x_2)$

$\underbrace{\quad}_{\Theta_{10}^{(1)}} \quad \underbrace{\quad}_{\Theta_{11}^{(1)}} \quad \underbrace{\quad}_{\Theta_{12}^{(1)}}$

$(-30 + 20 \cdot 1 + 20 \cdot 0)$

1.0
0.99

-0.01

-4.6      4.6   $z$

| | $x_1$ | $x_2$ | $h_\Theta(x)$ |
|---|---|---|---|
| | 0 | 0 | $g(-30) \approx 0$ |
| AND $\Rightarrow$ | 0 | 1 | $g(-10) \approx 0$ |
| | 1 | 0 | $g(-10) \approx 0$ |
| | 1 | 1 | $g(10) \approx 1$ |

$x_1$ AND $x_2$

## Negation:   NOT $x_1$



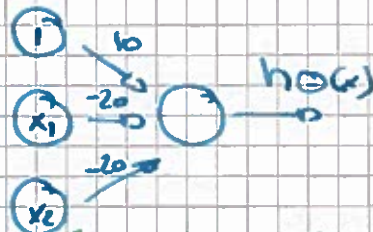| $x_1$ | $h_\Theta(x)$ |
|---|---|
| 0 | $g(10) \Rightarrow 1$ |
| 1 | $g(-10) \Rightarrow 0$ |

$$h_\Theta(x) = g(\Theta_{10} + \Theta_{11} x_1)$$

$$g(+10 - 20 x_1)$$

## Putting it together:  $x_1$ XNOR $x_2$



$x_1$ AND $x_2$       (Not $x_1$) AND (Not $x_2$)       $x_1$ OR $x_2$



|  |  |  | $\overline{X_1}$ AND $\overline{X_2}$ |  |
|  |  | AND |  | OR |
| $x_1$ | $x_2$ | $a_1^{(2)}$ | $a_2^{(2)}$ | $h_\Theta(x)$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 | 1 |

Multiple output units: One vs. all

| image | image | image | image |
|-------|-------|-------|-------|
| pedestrian | car | motorcycle | truck |



pedestrian — output would be a vector

car

motorcycle

truck

$$h_\Theta(x) \in \mathbb{R}^4$$

want $h_\Theta(x)$ $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  $h_\Theta(x)$ $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  $h_\Theta(x)$ $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$  $h_\Theta(x)$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

pedestrian      car      motorcycle      truck

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

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

pedestrian   car   motor-cycle   truck

image

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

$\mathbb{R}^4$  4 dim vectors