# Chapter 8

## Neural Networks - Non Linear hypothesis

$$g(\Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 + \Theta_3 x_1 x_2 + \Theta_4 x_1^2 x_2$$
$$+ \Theta_5 x_1^3 x_2 + \Theta_6 x_1 x_2^2 + \ldots)$$

by using only second order terms

$$x_1^2, x_1 x_2 \ldots \qquad x_{100}$$

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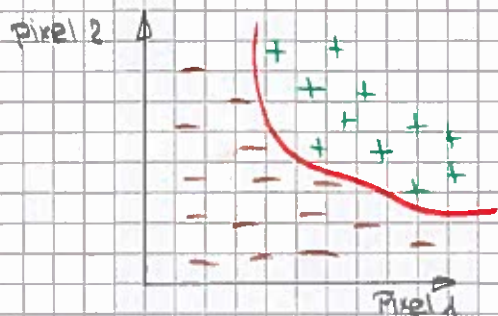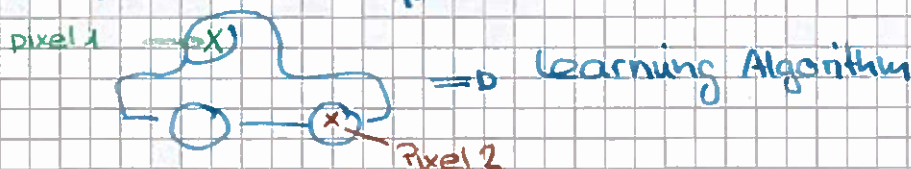$\vdots$
$x_{100}$

$n = 100$

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

just using the $x_1^2, x_2^2, x^3 \ldots x_{100}^2$
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

$\Rightarrow$ Learning Algorithm

Pixel 2



pixel 2

Pixel 1

+ Cars
− Non Cars

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

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

by using all quadratic features $(x_i \times x_j)$
= 3 million features per training example
is computational expensive
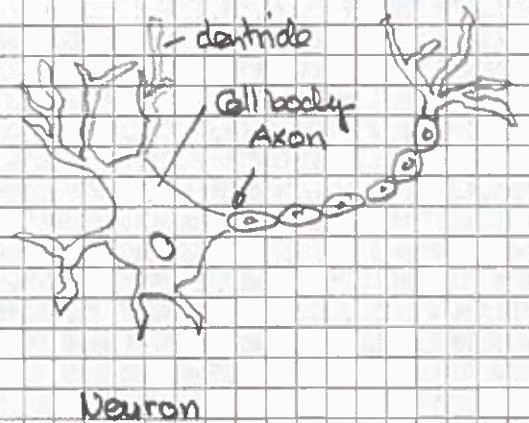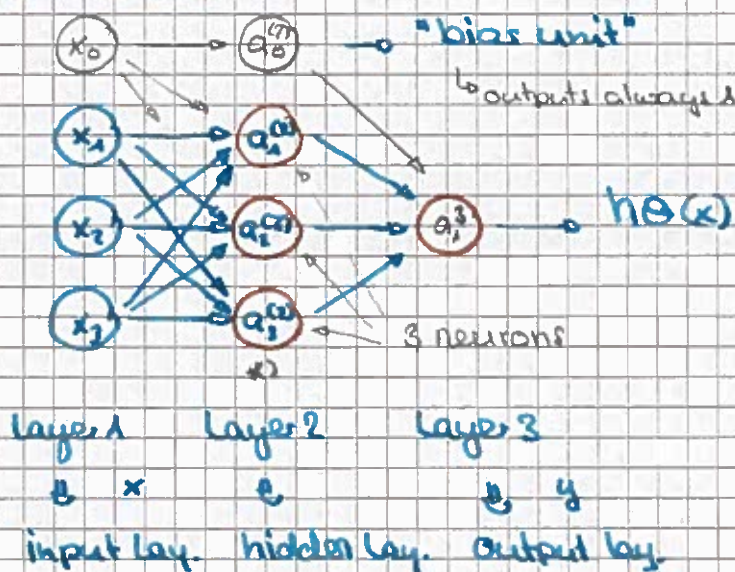Not good for algorith. with lot of features

## Neuron model: Logistic unit



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

dentrite - input

Axon output

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

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.

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

## Neural Network



"bias unit"

outputs always 1

$h_\Theta(x)$

3 neurons

*)

dentrite

Cell body

Axon

Neuron

Layer 1     Layer 2     Layer 3

   x                 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 $j$

$\Theta^{(j)}$ = 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

sigmoid

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3\right)$$
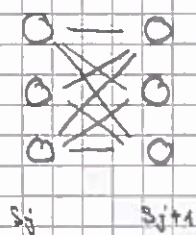
$$a_2^{(2)} = g\left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{13}^{(1)} x_3\right)$$

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

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

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

Example =>

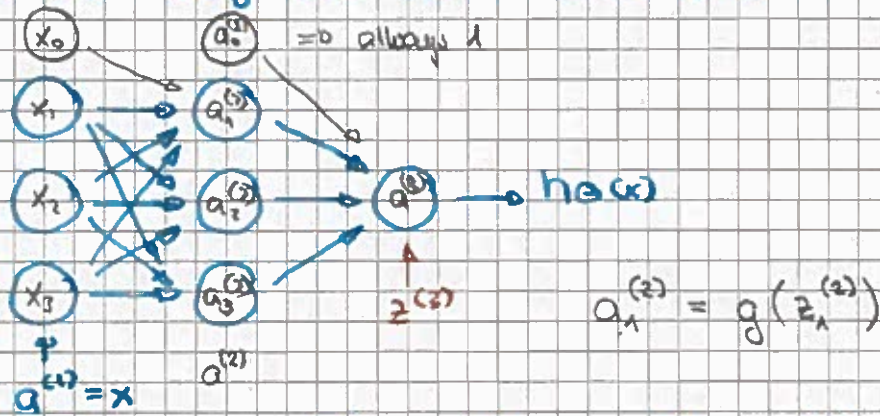$=> \mathbb{R}\ 3 \times 4$     $=> 4 \times 3$

$s_j$    $s_{j+1}$        $s_j$    $s_{j+1}$

Forward propagation: vectorized implementation



$a^{(1)} = x$

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

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

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

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

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

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

Vectorized implementation

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

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

$\underset{\mathbb{R}^3}{\qquad} \qquad \underset{\mathbb{R}^3}{\qquad} \implies$ 3 dim. vector

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

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

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

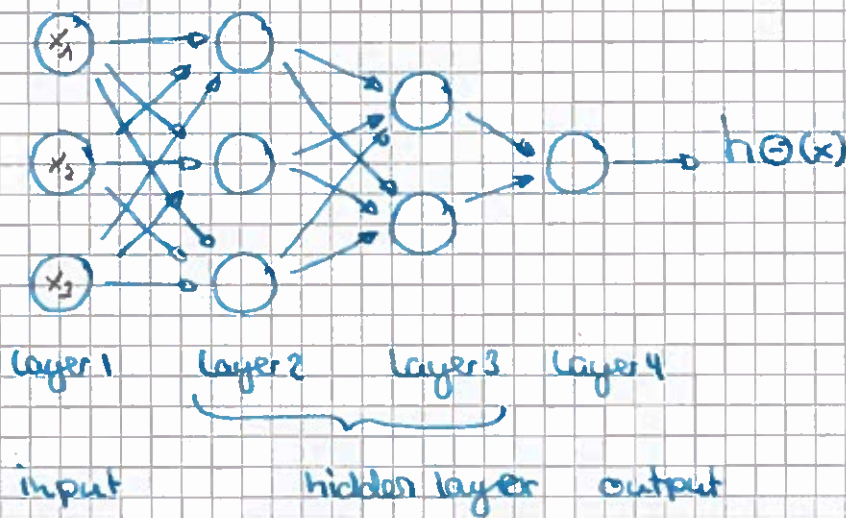Neural Network learning its own features.



Layer 1        Layer 2        Layer 3

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

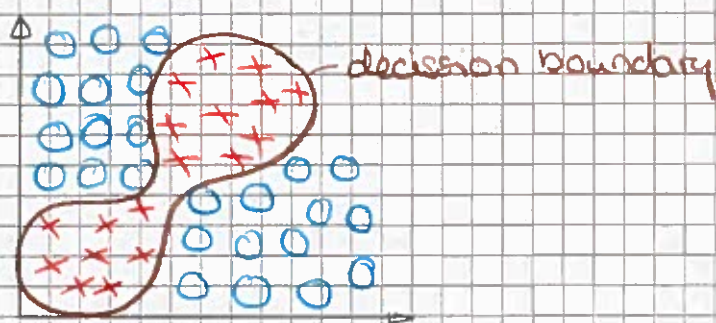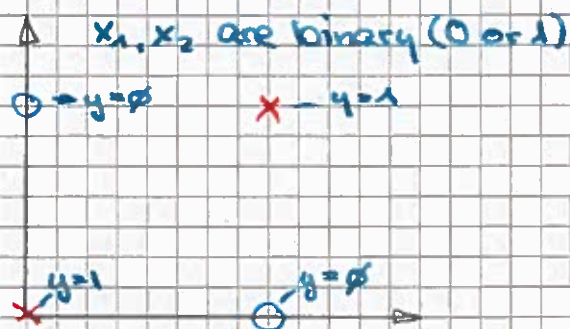it's just logistic regression except using the new features $a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$ instead of $x_1, x_2, x_3$



$$\rightarrow g\left(\Theta_{10}' x_0 + \Theta_{11}' x_1 + \Theta_{12}' x_2 + \Theta_{13}' x_3\right)$$

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)

$\bigoplus = y = \emptyset$  $\quad \times - y = 1$

$\times$ $y=1$  $\quad \bigoplus y=\emptyset$
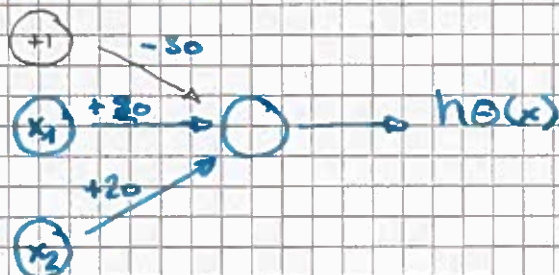
decision boundary

$y = x_1 \text{ XOR } x_2$

$x_1 \text{ XNOR } x_2$

$\text{NOT } (x_1 \text{ XOR } x_2)$

## Simple example AND:

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

$y = x_1 \text{ AND } x_2$

OR:

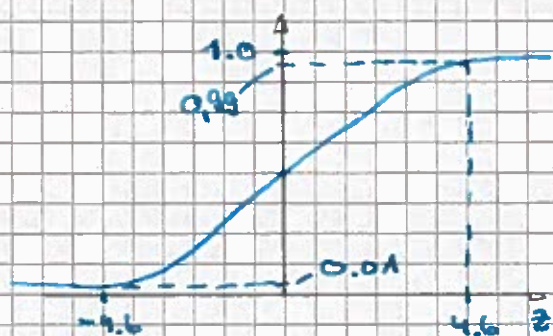| $x_1$ | $x_2$ | $h_\Theta(x)$ |
|---|---|---|
| 0 | 0 | $g(-10) = 0$ |
| 0 | 1 | $g(+10) = 1$ |
| 1 | 0 | $g(+10) = 1$ |
| 1 | 1 | $g(30) = 1$ |

$x_0$  $-10$
$x_1$  $+20$
$x_2$  $+20$

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

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

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

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

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

$x_1 \text{ AND } x_2$

Negation:      NOT $x_1$

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

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

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

Putting it together:    $x_1$ XNOR $x_2$

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

$$\overline{x_1} \text{ AND } \overline{x_2}$$

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

truck

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

image

$y^{(i)}$ 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

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

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