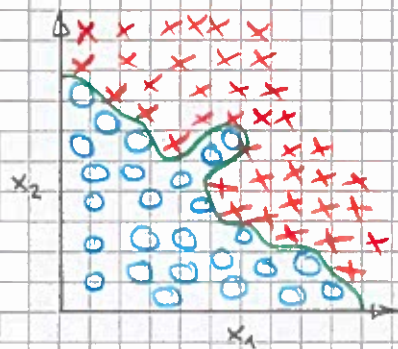


8. Neural Networks

- Non Linear Hypothesis

Neural Networks - Non Linear hypothesis Chapter 8 [8]



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

by using only second order terms

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

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 2000 features

\Rightarrow grows $O(n^2) \Rightarrow \frac{n^2}{2}$ \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

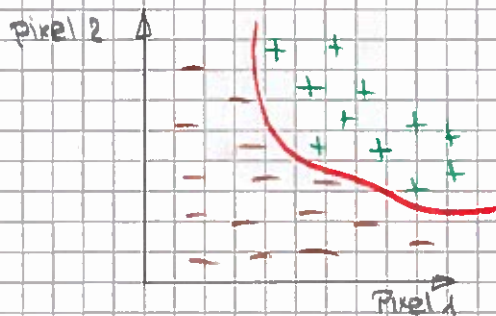
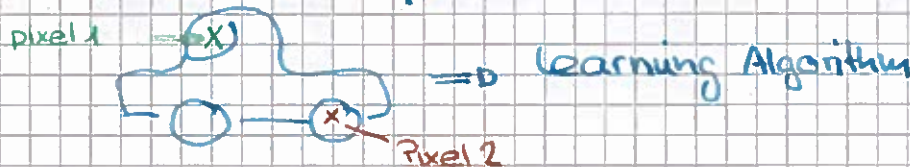
more complex functions, would not fit the complex dataset

while using cubic features it would $O(n^3)$

be about 170'000 features for

$n=100$ features

Why is non linear hypothesis relevant?



$+$ Cars
 $-$ Non Cars

\Rightarrow 50 x 50 pixels images \Rightarrow 2500 pixels

$n = 2500$ (7500 RGB)

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

$\underbrace{\hspace{10em}}_{2500}$

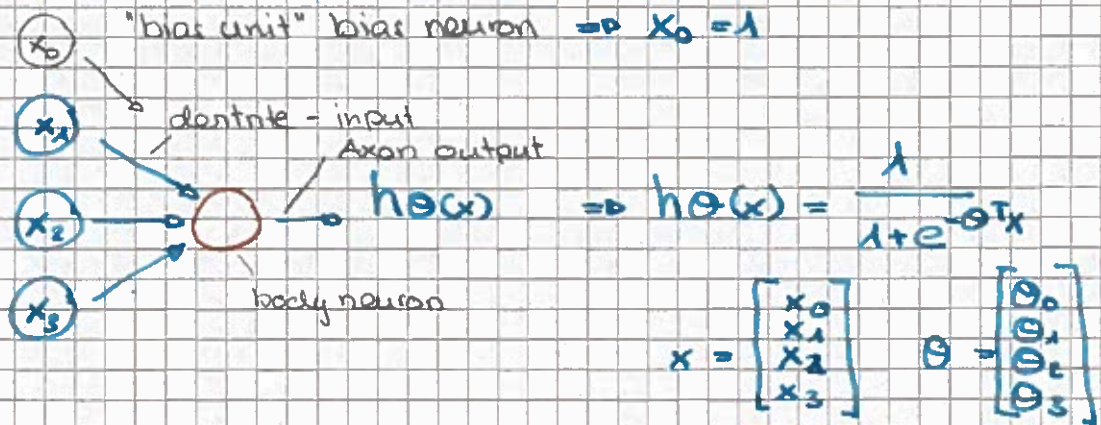
by using all quadratic features $(x_i x_j)$

\Rightarrow 3 million features per training example

is computational expensive

Not good for algorithm with lot of features

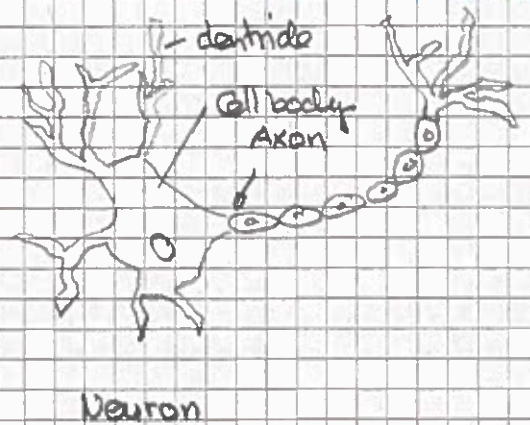
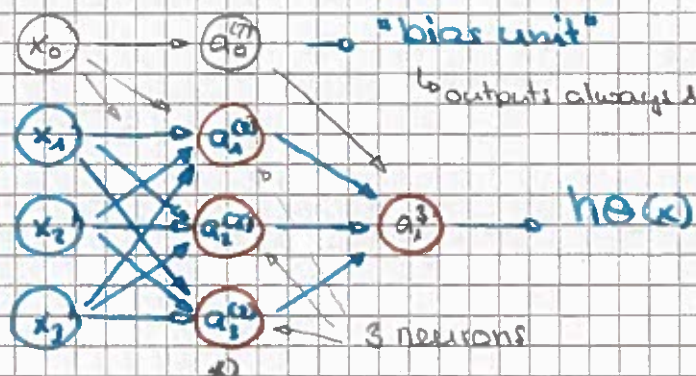
Neuron model: Logistic unit



Sigmoid (logistic) activation function.

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

Neural Network

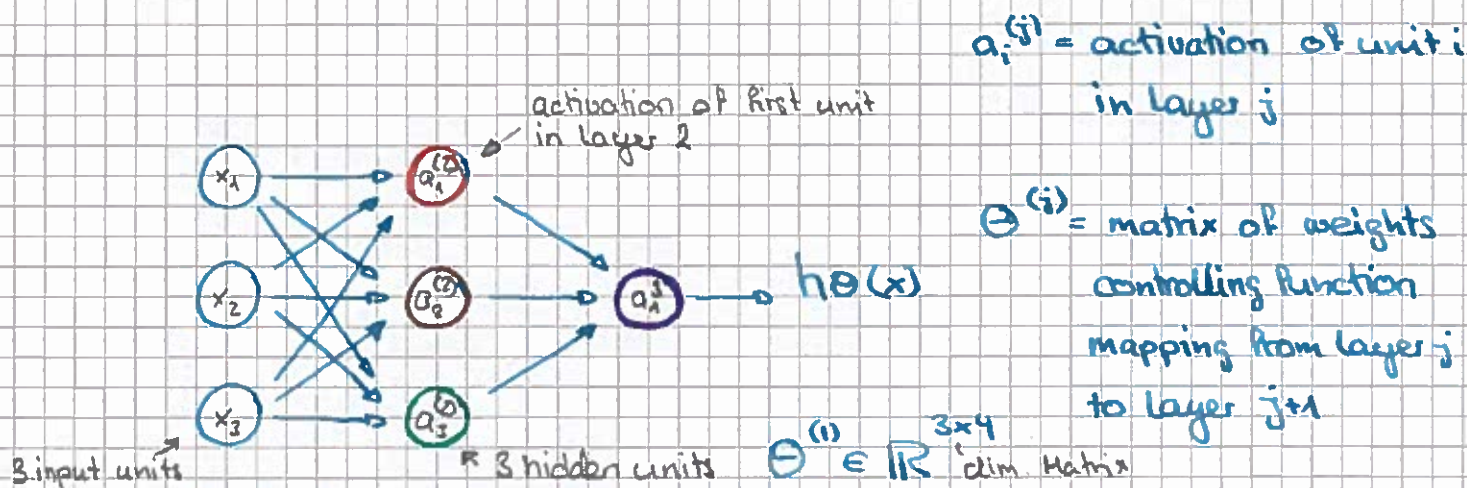


Layer 1 Layer 2 Layer 3

x z y

input lay. hidden lay. output lay.

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



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

Sigmoid

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

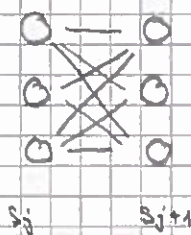
$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

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

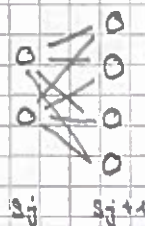
$$h(\theta(x)) = a_4^{(3)} = g(\Theta_{40}^{(2)} a_0^{(2)} + \Theta_{41}^{(2)} a_1^{(2)} + \Theta_{42}^{(2)} a_2^{(2)} + \Theta_{43}^{(2)} a_3^{(2)})$$

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

Example \Rightarrow

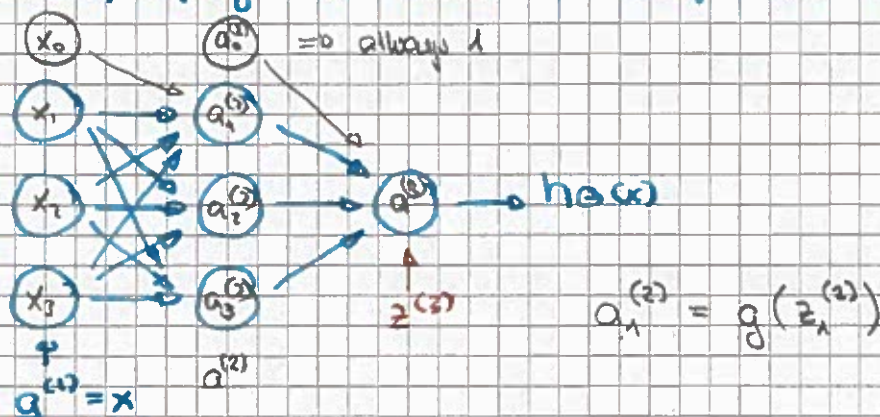


$$\Rightarrow \mathbb{R}^{3 \times 4}$$



$$\Rightarrow 4 \times 3$$

Forward propagation: vectorized implementation



$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3) \rightarrow z_1^{(2)}$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) \rightarrow z_2^{(2)}$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) \rightarrow z_3^{(2)}$$

$$h(x) = a^{(2)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) \rightarrow z^{(3)}$$

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

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

Vectorized
implementation

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

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

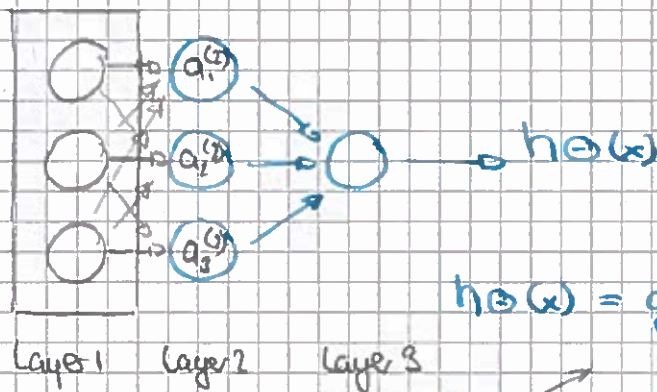
\mathbb{R}^3 \mathbb{R}^3 \Rightarrow 3 dim Vector

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

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

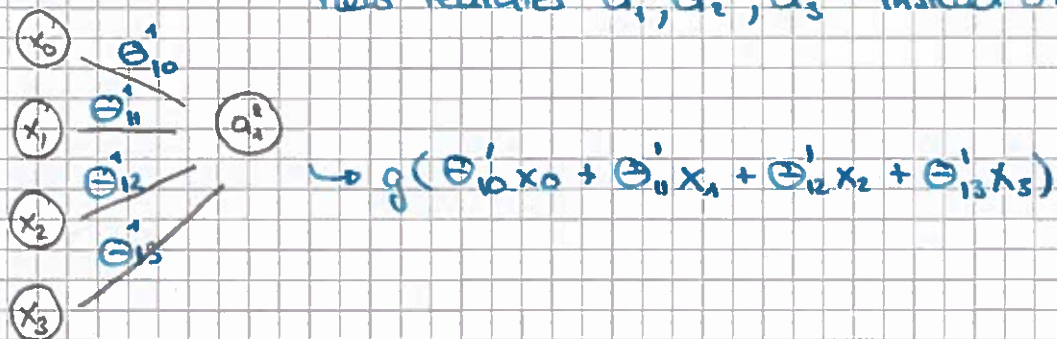
$$h(x) = a^{(3)} = g(z^{(3)})$$

Neural Network learning its own features.

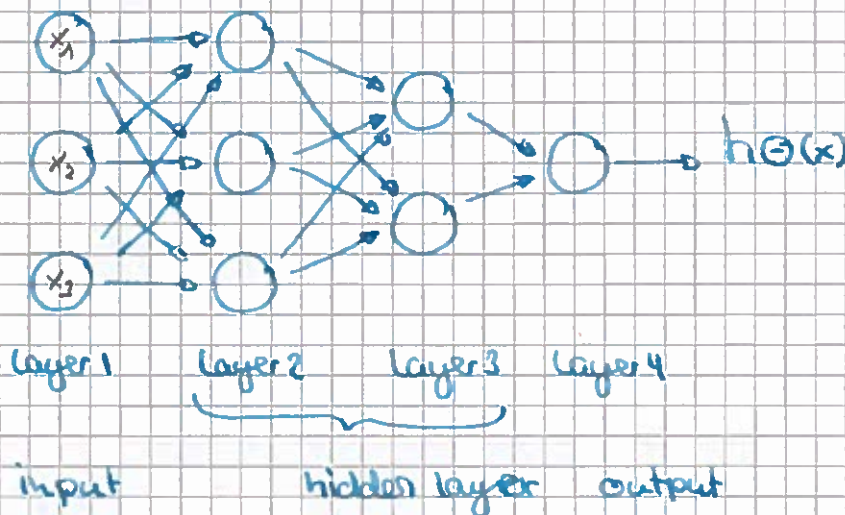


$$h(\theta(x)) = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

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

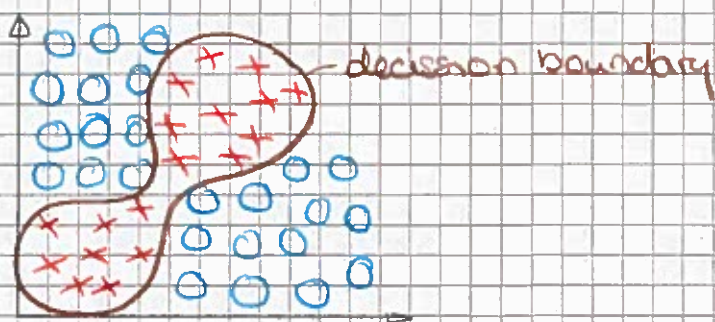
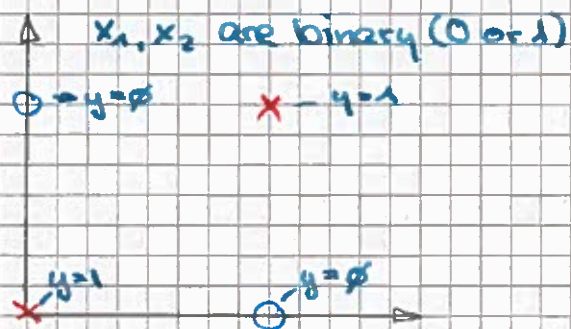


Other network architectures



input hidden layer output

Non linear classification example: XOR/XNOR



$$y = x_1 \text{ XOR } x_2$$

$$x_1 \text{ XNOR } x_2$$

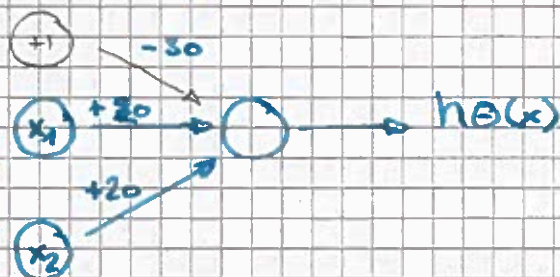
$$\text{NOT } (x_1 \text{ XOR } x_2)$$

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Simple example AND:

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

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



OR:

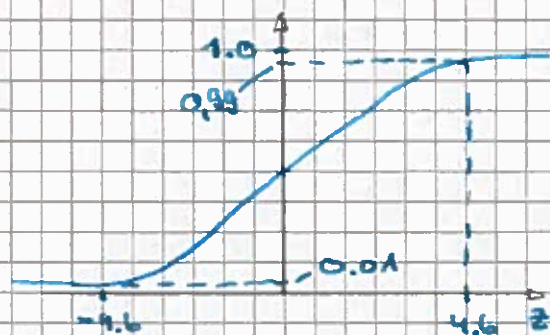
x_1	x_2	$h\theta(x)$
0	0	$g(-10) = 0$ (x_0)
0	1	$g(+10) = 1$
1	0	$g(+10) = 1$ (x_1)
1	1	$g(30) = 1$ (x_2)

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

$$\Rightarrow h\theta(x) = g(-30 + 20x_1 + 20x_2)$$

$\ominus_{x_0}^{(-)}$ $\ominus_{x_1}^{(+)}$ $\ominus_{x_2}^{(+)}$

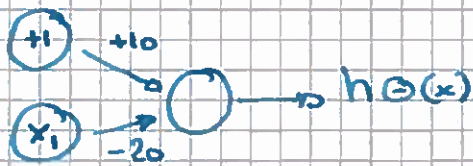
$$(-30 + 20 \cdot 1 + 20 \cdot 0)$$



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

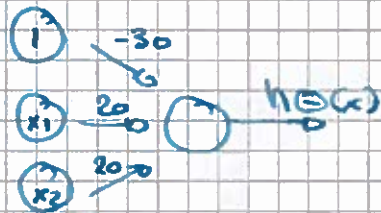


x_1	$h\Theta(x)$
0	$g(10) = 1$
1	$g(-10) = 0$

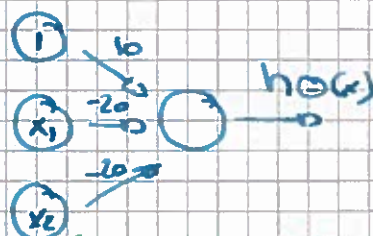
$$h\Theta(x) = g(\Theta_0 + \Theta_1 x_1)$$

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

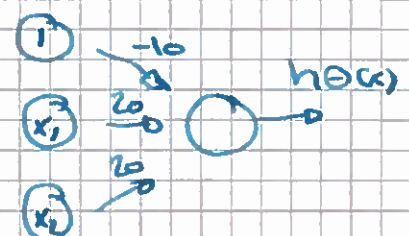
Putting it together: x_1 XOR x_2



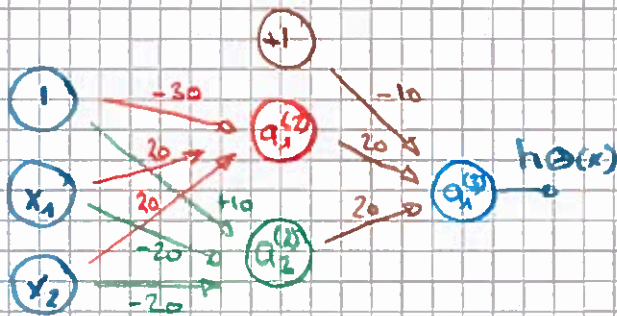
x_1 AND x_2



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



x_1 OR x_2

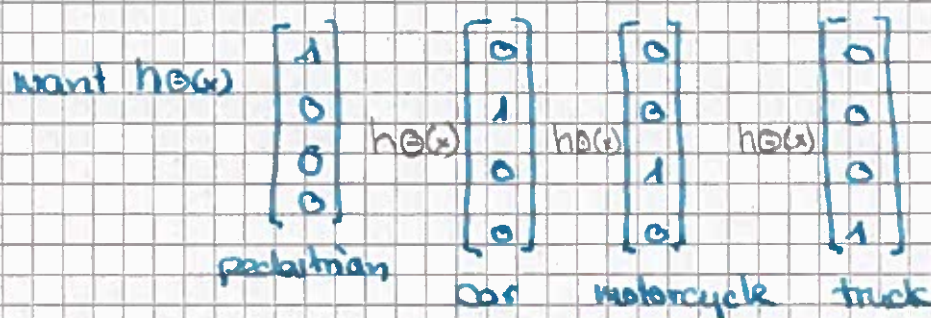
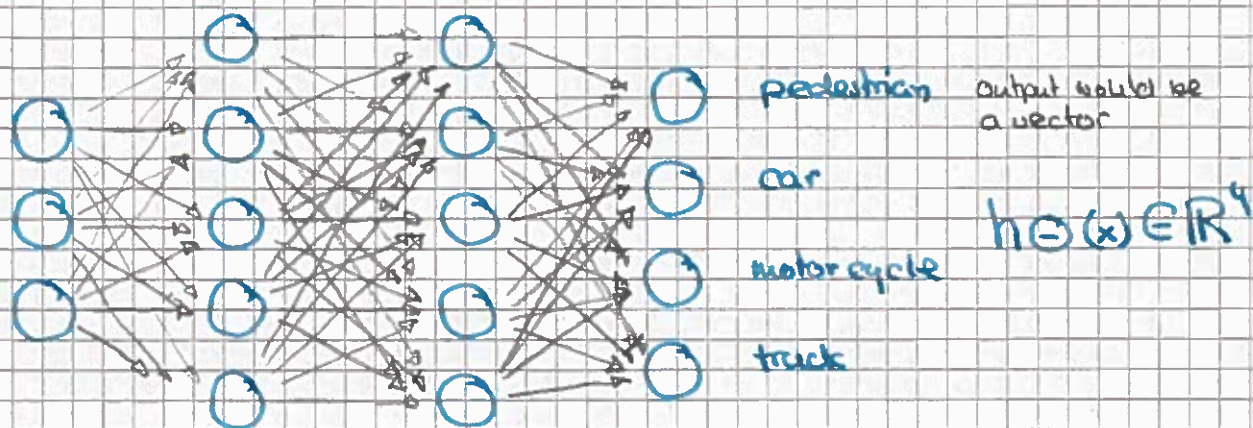
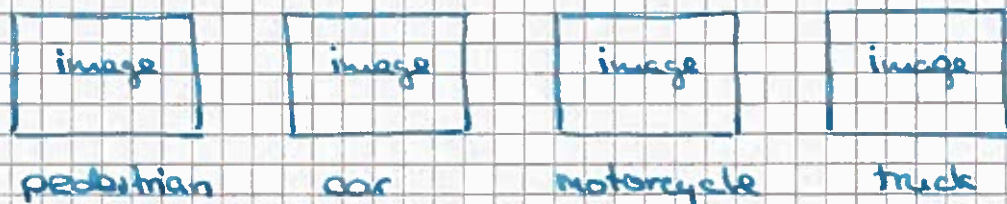


x_1	x_2	$q_1^{(1)}$	$q_2^{(1)}$	$h\Theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

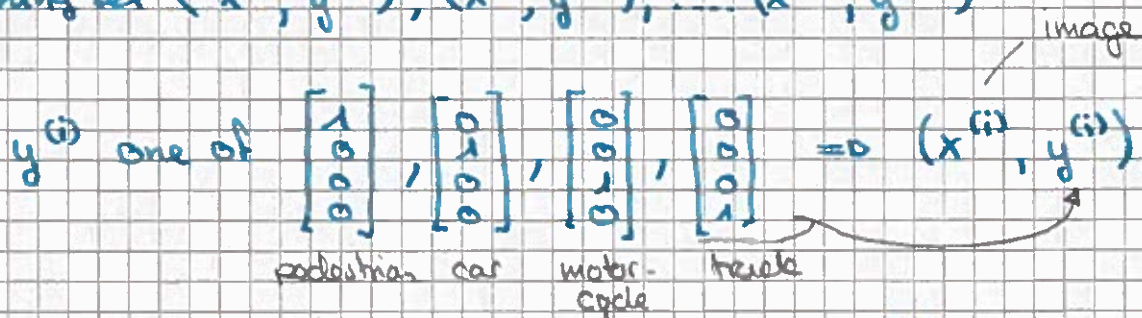
Multiclass Classification

18

Multiple output units: One vs. all



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



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

\mathbb{R}^4 4 dim. vectors