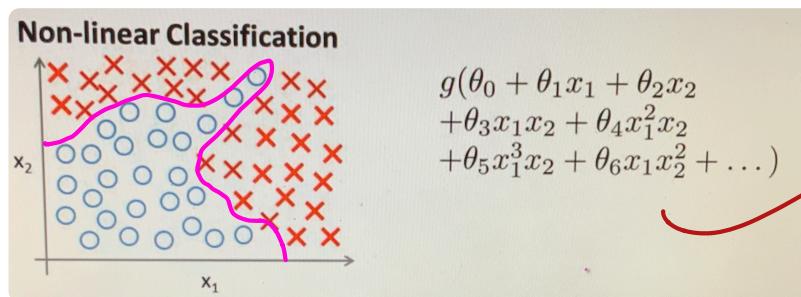


Stanford ML

Representing Neural Networks

- Why linear and logistic regression aren't enough.

Ex: Non - Linear Classification



- Easy to create a non-linear hypothesis with only a small number of features

- If $n \sim 100$, you end up with ≈ 5000 features to $O(n^2)$
- Complexity: $O(n^2/2)$
- Can lead to **OVERTFITTING** via too many features

Computer Vision

Computer Vision: Car detection

Cars		
Not a car		

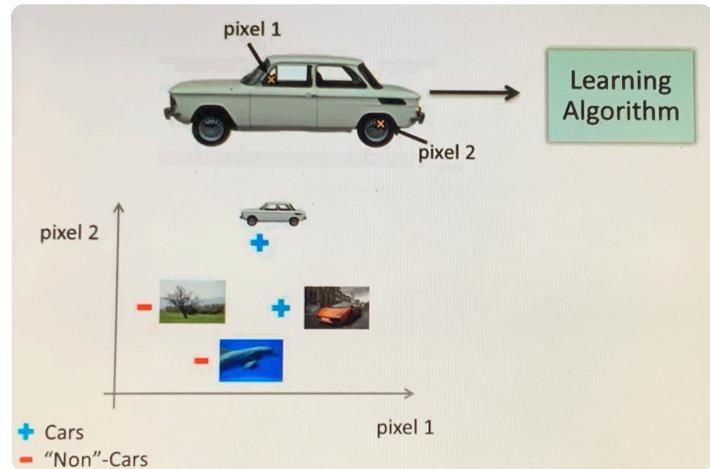
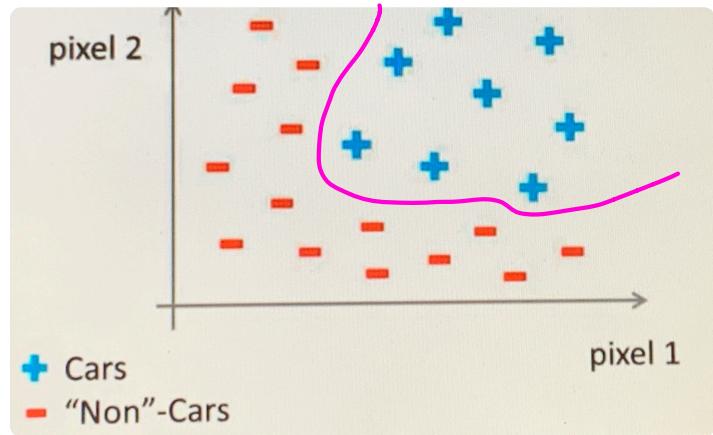
Testing:

What is this?

- Identifying what is a car involves giving the

algorithm a labelled data set for training
the classifier

- By comparing pixel intensities between images like on the right, we end up with something like



Therefore we need
a non-linear hypothesis
in order to classify
images into cars and not
cars

USING 50×50 images

$$\Rightarrow 2500 \text{ pixels} \Rightarrow n = 2500$$

$$\text{OR } n = 7500 \text{ (RGB)}$$

$$x = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel } 2500 \text{ intensity} \end{bmatrix} \quad \therefore \text{in order to create a quadratic } h_0(x) \dots$$

Quadratic Features
 $(x_i \times x_j)$

$$\approx 3 \times 10^6$$

WAT
TOO
BIG

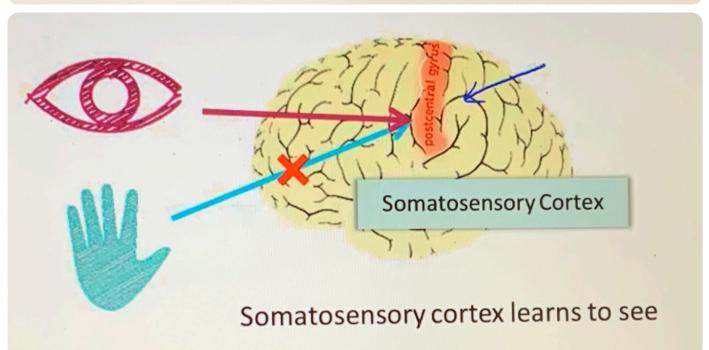
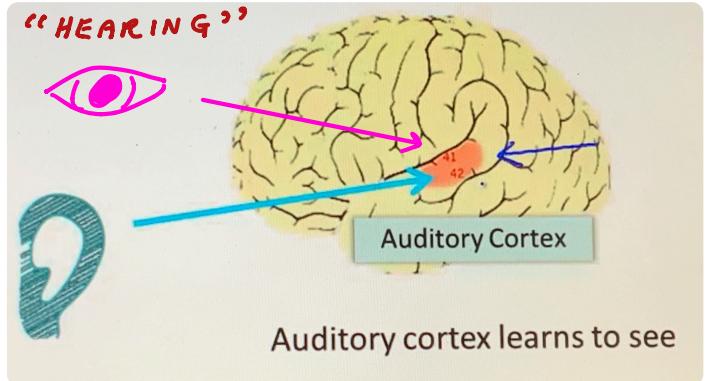
Neurons and the Brain

- ORIGIN: Algorithms that try mimic the brain
- The "one learning algorithm" hypothesis
- If you re-wire the signal from the ears to the auditory cortex to now come from the EYE

=> The auditory cortex

LEARNs TO SEE

- The same happens with your somatosensory cortex, responsible for the sensation of touch



NEURO - REWIRING EXPERIMENTS

EXAMPLES

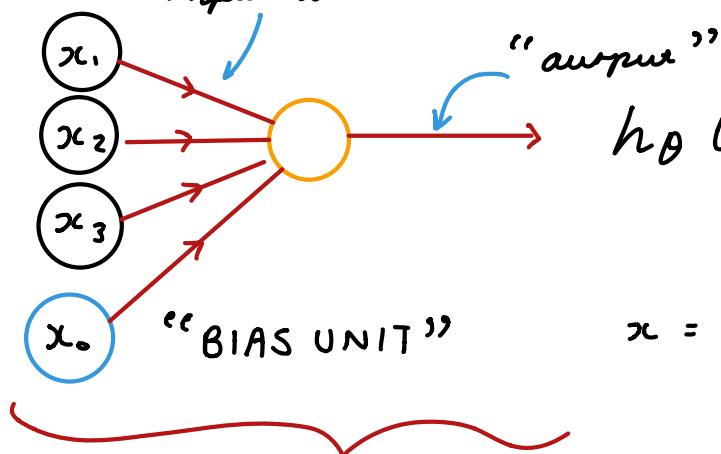


We wish to emulate how the brain does this, what algorithm does it use?

Model Representation

NEURON MODEL: Logistic unit

"input units"



SINGLE NEURON (SIGMOID)

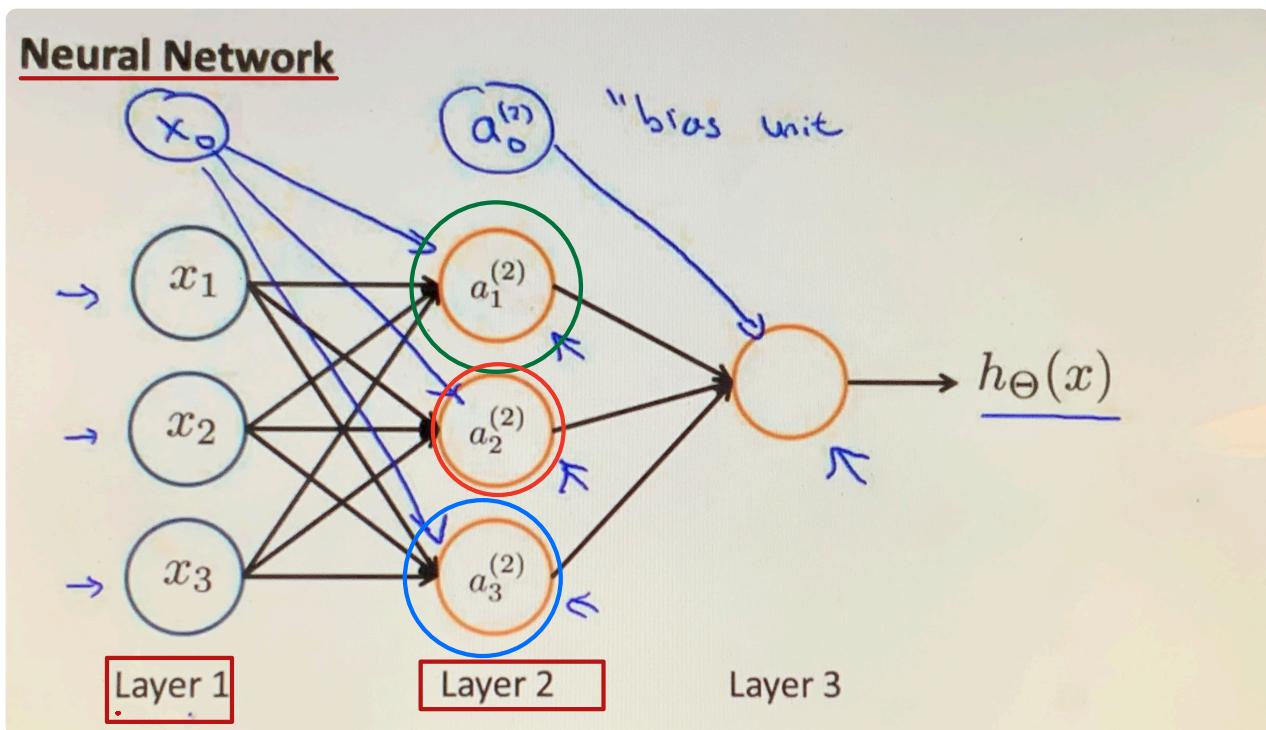
ACTIVATION FUNCTION

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

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

WEIGHTS / PARAMETERS

- A neural network can simply be characterized as a collection of these neurons:



INPUT

x

"3 units" "3 units"

HIDDEN

$a^{(2)}$

OUTPUT

y

$$\therefore \theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

NOMENCLATURE:

- * $a_i^{(j)}$ = ACTIVATION of unit i in layer j
- * $\Theta^{(j)}$ = matrix of weights connecting neurons
mapping from layer $j \rightarrow j+1$
- $a_1^{(2)} = g \left\{ \Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 \right\}$
- $a_2^{(2)} = g \left\{ \Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3 \right\}$
- $a_3^{(2)} = g \left\{ \Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3 \right\}$
- $h_\theta(x) = a_1^{(3)} = g \left\{ \Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right\}$

- In general, if a network has s_j units in layer j ,
 s_{j+1} units in layer $j+1$, then

$$\Theta^{(j)} \in \mathbb{R}^{s_{j+1} \times (s_j + 1)}$$

Layer 1: $s_1 = 3$

3

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

Layer 2: $s_2 = 3$

3

$$\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

Layer 3: $s_3 = 7$

3

Representing

the

Hypothesis

VECTORIZATION $a_1^{(3)} = g \left\{ \Theta$

- Let $a_1^{(2)} = g(z_1^{(2)})$ where

$$z_1^{(2)} = \Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 = \begin{bmatrix} \Theta_{1j}^{(1)} \end{bmatrix}^T x$$

- The same can be done for all activation nodes

$$\therefore \text{let } \underline{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \text{ and } \underline{z}^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

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

$$\Rightarrow \text{Add } \underline{\alpha}_0^{(2)} = 1 \therefore \underline{\alpha}^{(2)} \in \mathbb{R}^n$$

$$\Rightarrow \underline{z}^{(3)} = \Theta^{(2)} \underline{\alpha}^{(2)} \Rightarrow h_{\Theta}(\underline{x}) = \underline{\alpha}^{(3)} = g(\underline{z}^{(3)})$$

- In general : → g applied element-wise

$$\underline{\alpha}^{(i)} = g(\underline{z}^{(i)}) = g(\Theta^{(i-1)} \underline{\alpha}^{(i-1)})$$

FORWARD PROPAGATION

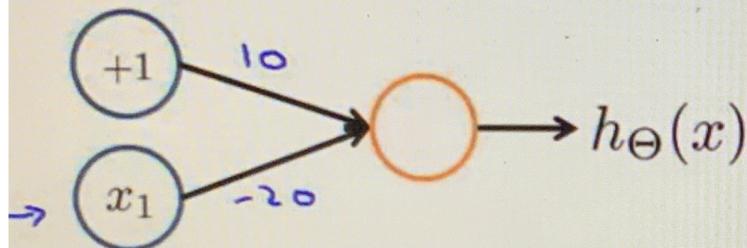
$$\underline{\alpha}_1^{(3)} = g(\Theta^{(2)} \underline{\alpha}^{(2)})$$

- Effectively, multiple logistic regressions where the algorithm ends towards the most OPTIMAL FEATURES before the final output layer - performs logistic regression with those $\underline{\alpha}^{(N-1)}$ for N layers.

Negation:

NOT x_1

NOT x_1



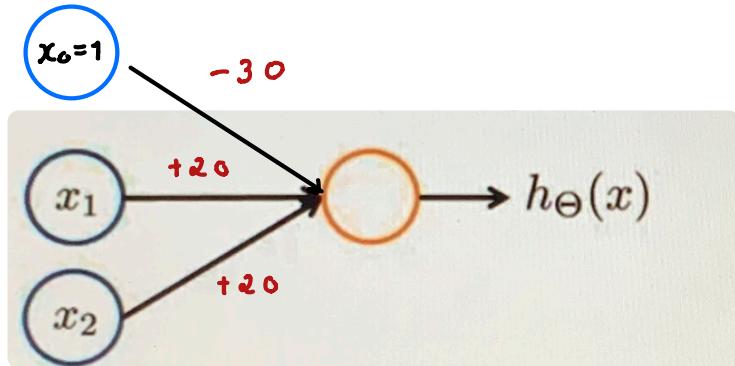
x_1	$h_{\Theta}(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

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

- Simply place 'large' negative weight in front of the feature / variable you wish to negate

AND GATE EXAMPLE

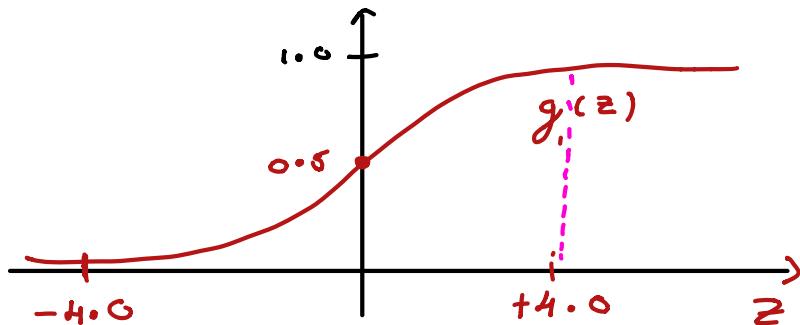
* $x_1, x_2 \in \{0, 1\}$ * $y = x_1 \text{ AND } x_2$



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	$g(10) \approx 1$

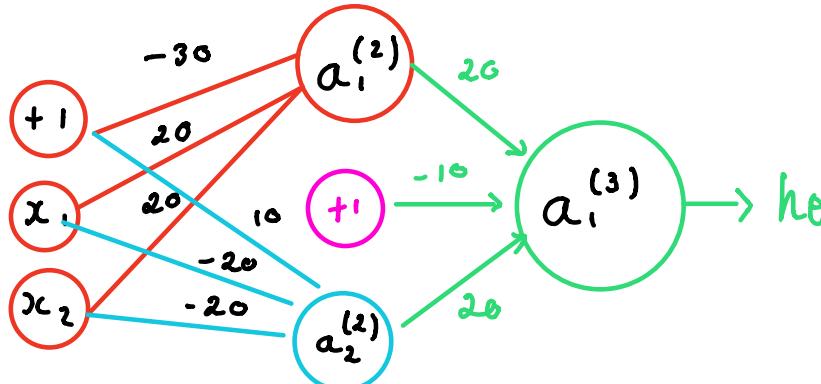
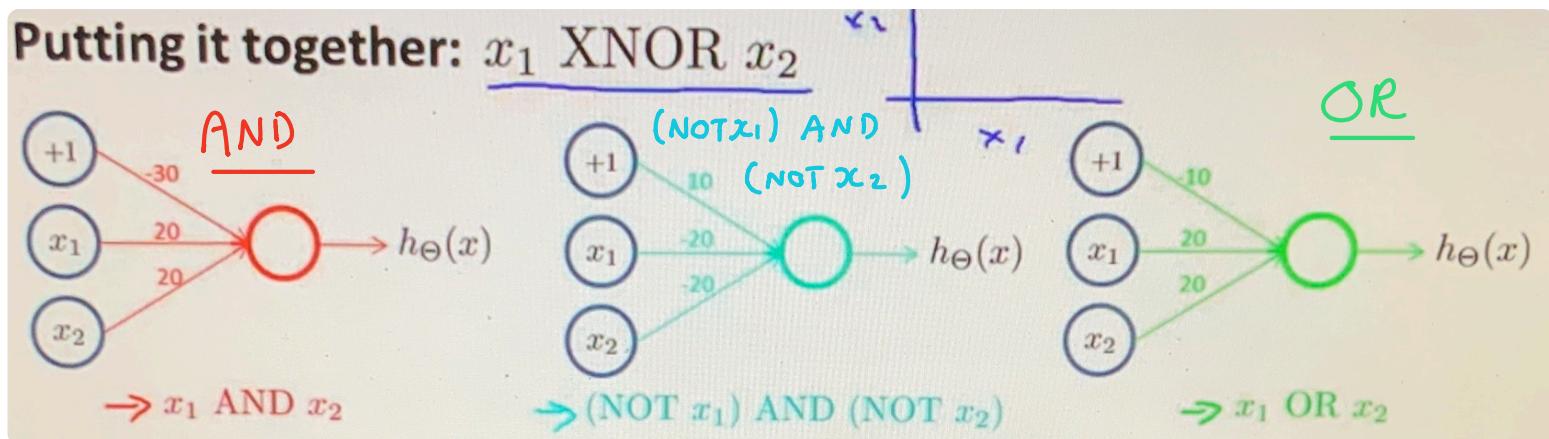
- Assigning weights of $\Theta = (-30, +20, +20)$

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



$$\therefore h_\Theta \approx x_1 \text{ AND } x_2$$

- If you set $\Theta_0 : -30 \rightarrow -10$
 $h_\Theta : \text{AND} \rightarrow \text{OR}$

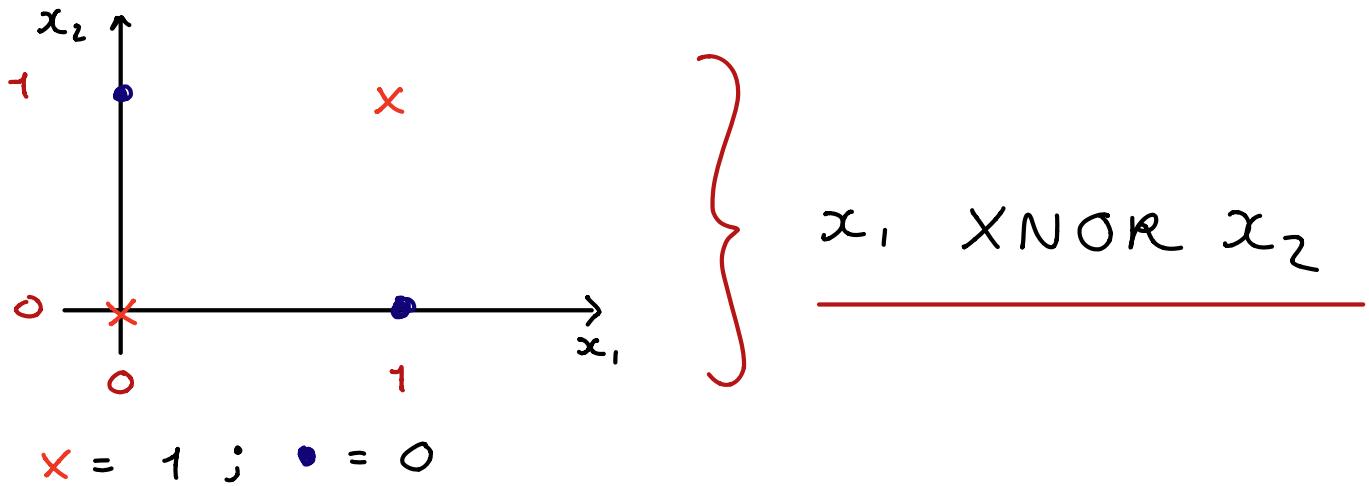


x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\Theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

$$\Theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix};$$

$$\Theta^{(2)} = \begin{bmatrix} -10 & 20 & 20 \end{bmatrix};$$

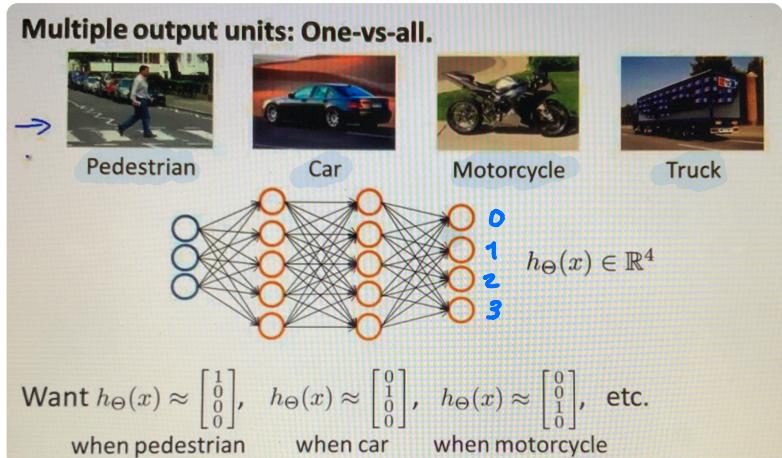
- This neural network composed of simple neurons has created a non linear hypothesis for us to classify this data



ANN: Multiclass Classification

- Extension of one vs. all logistic regression algorithm
- The output can now be one of 4 outcomes

$\therefore h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ when pedest.



$$\therefore h_0(x) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \text{when etc. car}$$

- Training set: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$
where $y^{(i)}$ is one of:

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

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

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

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

$\} y^{(i)}$

person

car

bike

truck

↳ Rather than $y = \{1, 2, 3, 4\}$ as in LOG RECr.

- We want

$h_0 \propto y$