

① MOTIVATIONS.

1.1. Non-linear hypothesis:

Non-linear Classification:

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

This method works well when we have only 2 features (x_1, x_2)

But for many ML problems would have a lot more features!

if we have $n=100$ features, you end up with ≈ 5000 features (quadratic features!) $\approx \frac{n^2}{2}$. Because of many features, it may overfit the training set.→ if you add cubic polynomials, there would be $\approx 170\,000$ features. ($O(n^3)$)ex/ if we have 50×50 pixel img, there are $n=2500$ (7500 if RGB) features.Quadratic features ≈ 3 million, features.

if you try to classify imgs only by looking at pixel values, probably you will overfit! cannot classify imgs.

1.2. Neurons and the Brain:

NNs: Algorithms that try to mimic the brain.

The "one learning algorithm" hypothesis. Auditory cortex learns to see. } neural-rewiring experiments.
Somatosensory cortex learns to see. }

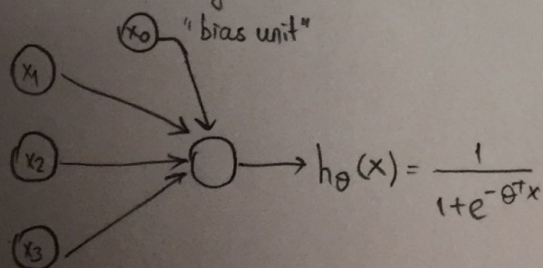
② NEURAL NETWORKS

2.1. Model Representation - 1:

Dendrite — "input wires"

Axon — "output wire"

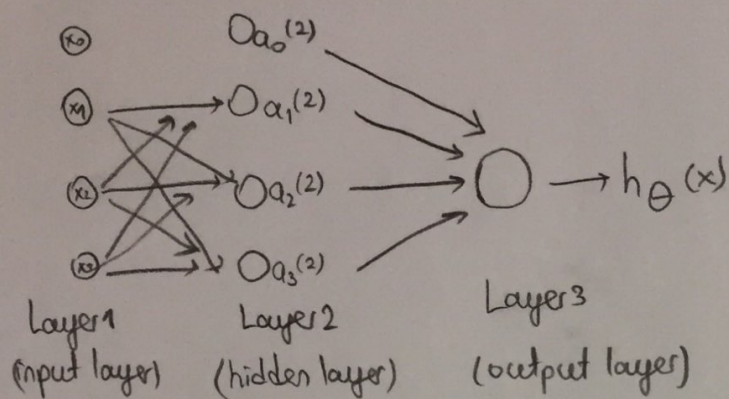
Neural model: Logistic unit



Sigmoid (logistic) activation fn

Weights = θ 's. = $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$ params.

NN:



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

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

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

\vdots

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

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

2.2. Model Representatn -2:

Forward Propagatn: Vectorized implementatn.

$$a_3^{(2)} = g(z_3^{(2)}) \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad z^{(2)} = \Theta^{(1)} \cdot x$$

$\mathbb{R}^3 \quad \mathbb{R}^3$

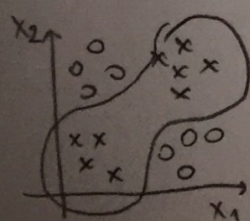


Instead of being constrained to feed the features (x_1, x_2, x_3) to logistic regression, it gets to learn its own features (a_1, a_2, a_3) to feed into logistic regression.

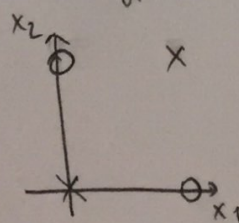
\Rightarrow it can learn some pretty interesting and complex features & can end up with better hypothesis

③ APPLICATIONS.

NN can be used to learn complex non-linear hypothesis.



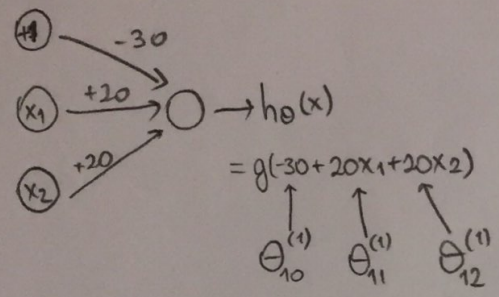
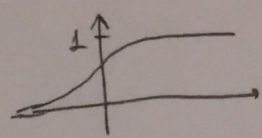
we'd like to learn a non-linear decision boundary



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

$$x_1 \text{ XNOR } x_2$$

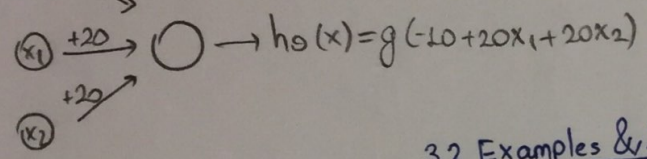
AND $y = x_1 \text{ AND } x_2$



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	$g(10) \approx 1$

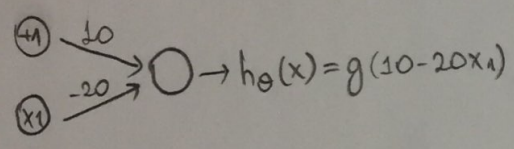
$h_{\theta}(x) \approx x_1 \text{ AND } x_2$

OR: $y_0 = -10$

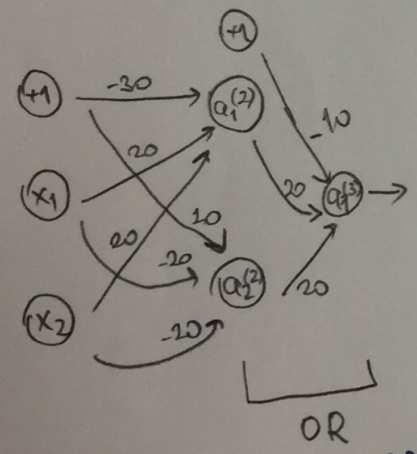


3.2. Examples & Intuits - 2:

NOT:



XNOR:

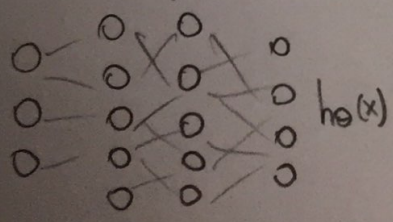


x_1	x_2	$q_1^{(2)}$	$q_2^{(2)}$	$h_{\theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

3.3. Multiclass Classification:

multiple output units: one-vs-all (pedestrian, car, motorcycle, truck)

Want $h_{\theta}(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ when pedestrian, $h_{\theta}(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ when car, $h_{\theta}(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ when motorcycle, $h_{\theta}(x) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ when truck.



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