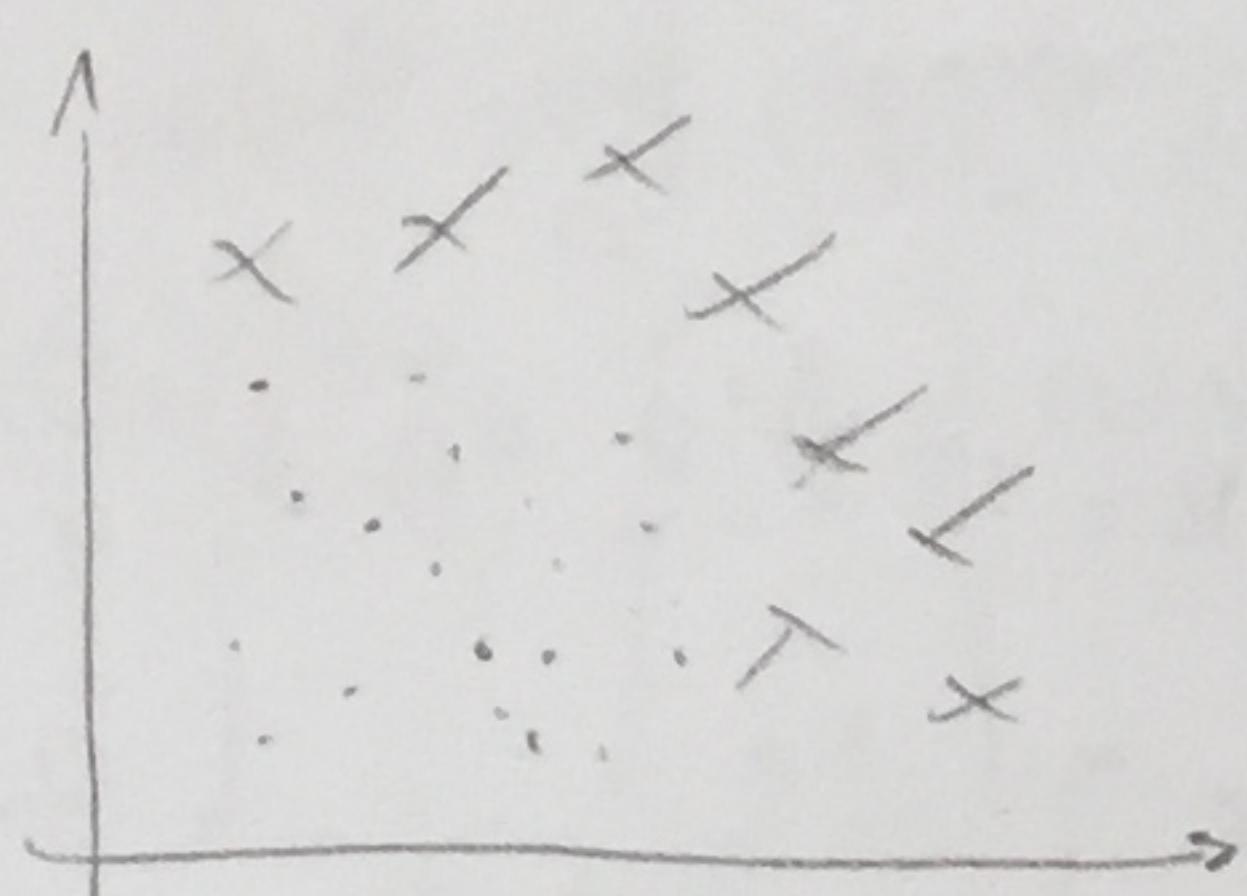


Non-Linear Hypothesis:

could use $g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \dots)$

$x_1 = \text{size}$

$x_2 = \# \text{ bedrooms}$

$x_3 = \# \text{ floors}$

$x_4 = \text{age}$

\vdots

x_{100}

Non-Linear

$x_1^2, x_1 x_2, x_1 x_3, x_1 x_4, \dots, x_1 x_{100}$

$x_2^2, x_2 x_3, x_2 x_4, \dots, x_2 x_{100}$

\downarrow
 $\approx 5000 \text{ features!}$

• # features $O(n^2) \approx \frac{n^2}{2}$

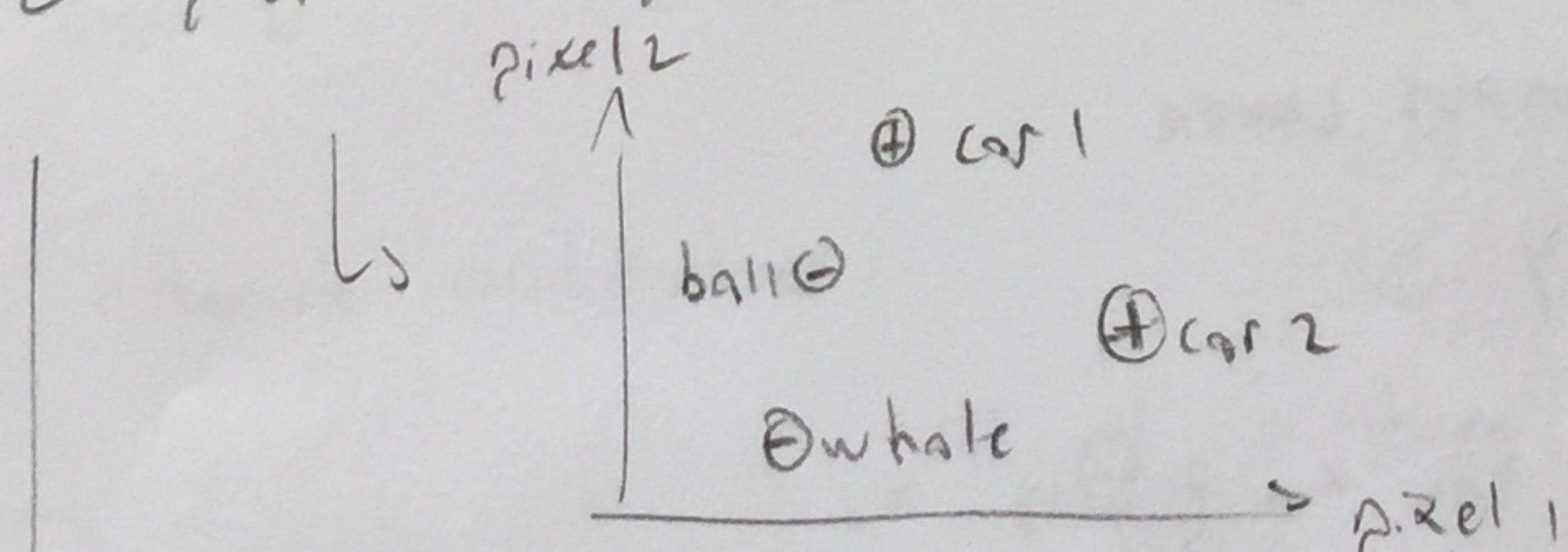
⇒ computationally ~~I.~~. Could just include $x_1^2, x_2^2, \dots, x_{100}^2$, but doesn't cover all decision boundary forms.

↳ additionally, including $x_1^2 x_2 \dots \Rightarrow$ # feat. $O(n^3)$!

⇒ for many problems, n is large.

e.g.: picture of car → matrix of pixel intensity values

2 pixs, 2 pixel locations (same location on each pix)



↳ 50x50 pixel images → $n = 2500$!

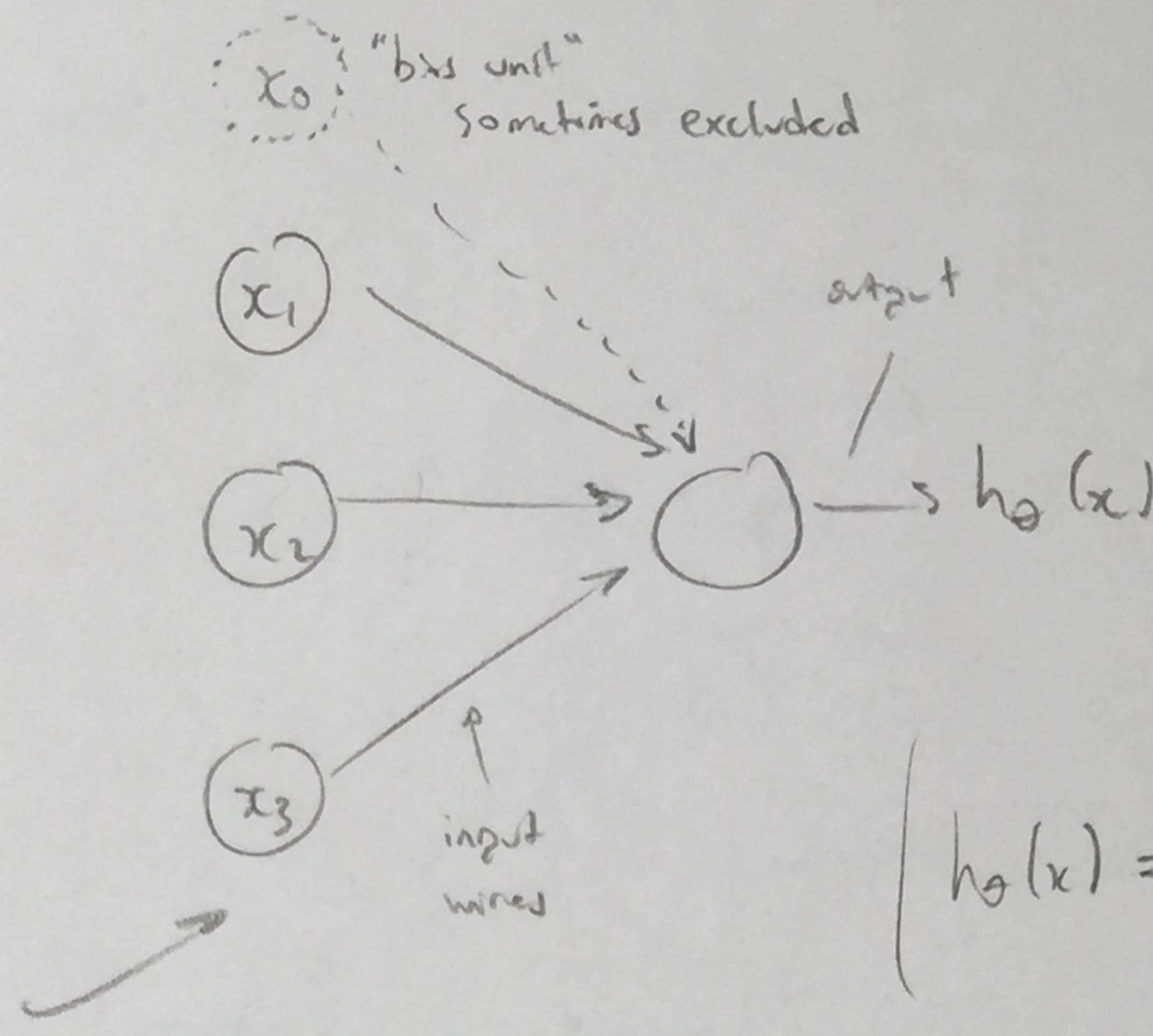
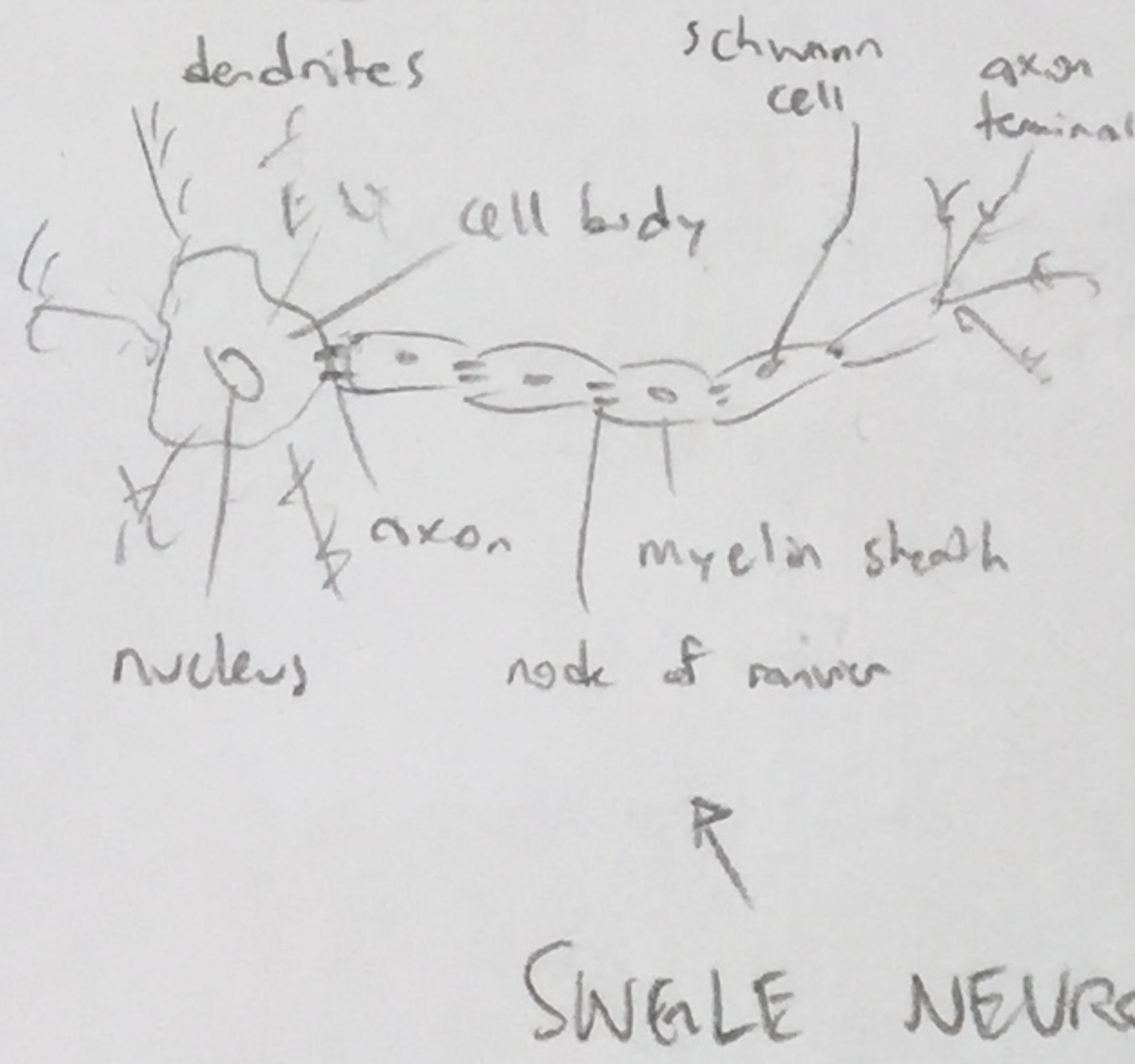
$$x = \begin{bmatrix} \text{pixel 1} \\ \text{pixel 2} \\ \vdots \\ \text{pixel 2500} \end{bmatrix}$$

↳ quadratic features: $x_i \times x_j \approx 3 \text{ MM}$
features! $\left(\frac{n^2}{2}\right)$

\Rightarrow Neural networks \rightarrow better way to learn complex non-linear hypothesis than using massive # of features!

\hookrightarrow plug any sensor into the brain \rightarrow brain will learn how to interpret this data!

Model Representation 1:



$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

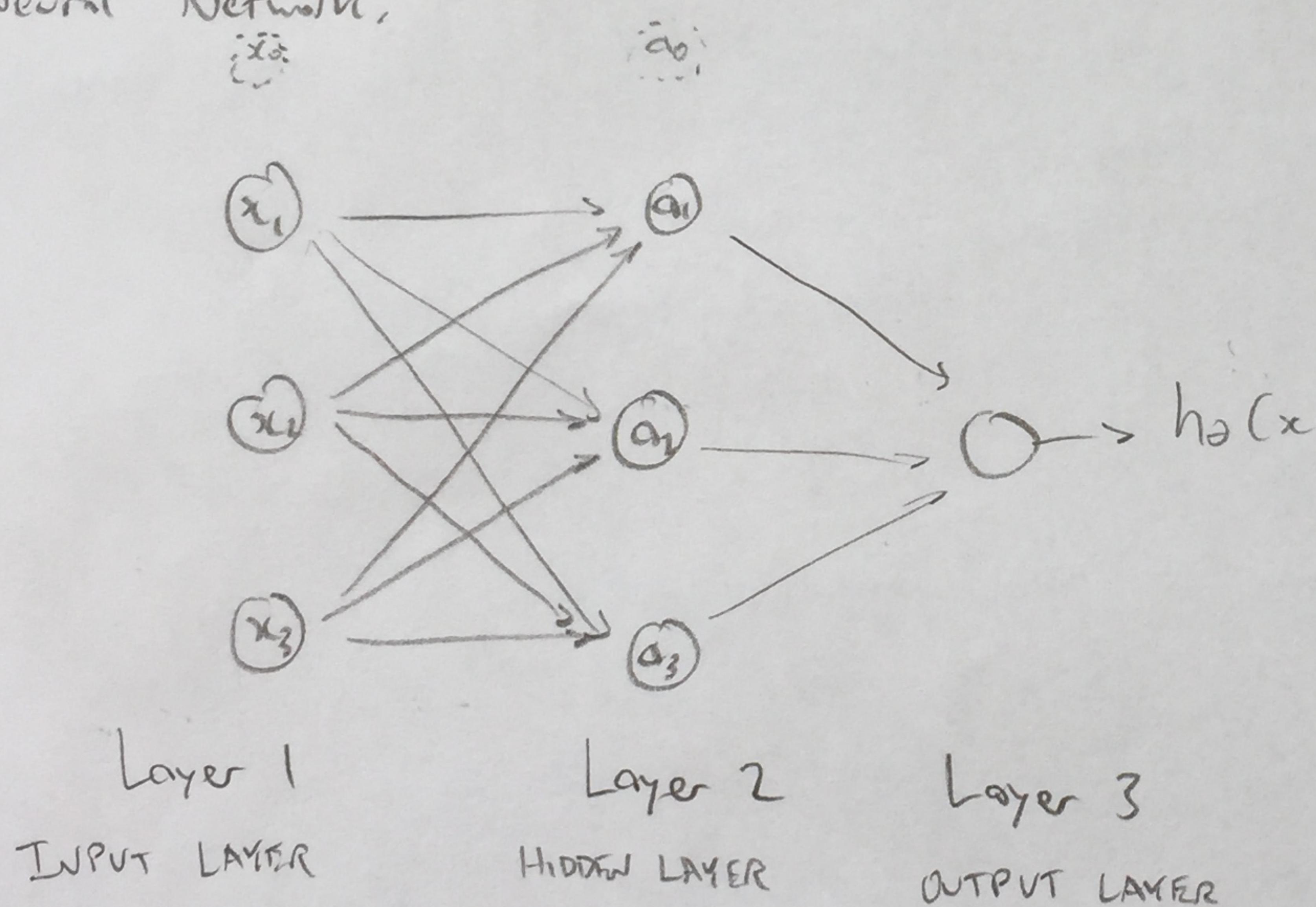
n=3 here

Parameters
= "weights"

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

$g(z)$: Sigmoid (logistical) activation function

Neural Network:



$a_i^{(j)}$ \rightarrow activation of unit i in layer j

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

x

y

$$a_1^{(2)} = g(\theta_{00}^{(1)} x_0 + \theta_{10}^{(1)} x_1 + \theta_{20}^{(1)} x_2 + \theta_{30}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{01}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{21}^{(1)} x_2 + \theta_{31}^{(1)} x_3)$$

$$a_3^{(2)} = \dots$$

$$h_\theta(x) = a_1^{(3)} = g(\theta_{00}^{(2)} a_0^{(2)} + \theta_{10}^{(2)} a_1^{(2)} + \theta_{20}^{(2)} a_2^{(2)} + \theta_{30}^{(2)} a_3^{(2)})$$

If network has s_j units in layer j , s_{j+1} units in layer $j+1$ $\rightarrow \theta^j \in \mathbb{R}^{s_{j+1} \times (s_j + 1)}$

Model Representation II

$$a_1^{(2)} = g(z_1^{(2)}) \quad \text{layer 2 (hidden layer)}$$

where $z_1^{(2)} = \Theta_{1,0}^{(1)}x_0 + \Theta_{1,1}^{(1)}x_1 + \Theta_{1,2}^{(1)}x_2 + \Theta_{1,3}^{(1)}x_3$

$$a_2^{(2)} = g(z_2^{(2)})$$

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

if $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$ $z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \Rightarrow z^{(2)} = \Theta^{(1)}x \in \mathbb{R}^3$

$$a^{(2)} = g(z^{(2)}) \in \mathbb{R}^3$$

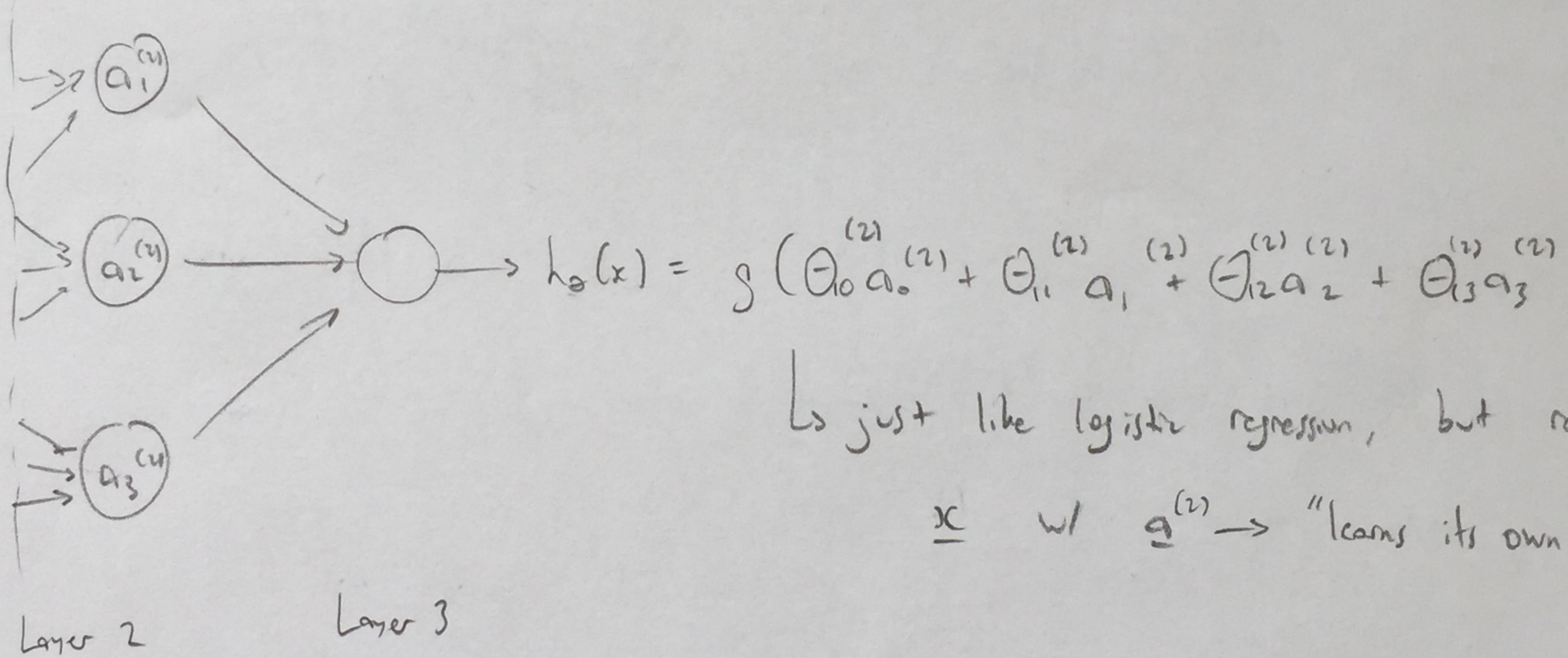
$$\Rightarrow \text{add } a_0^{(2)} = 1 \quad (x_0) \rightarrow a^{(2)} \in \mathbb{R}^4$$

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

$$h_\theta(x) = a^{(3)} = g(z^{(3)})$$

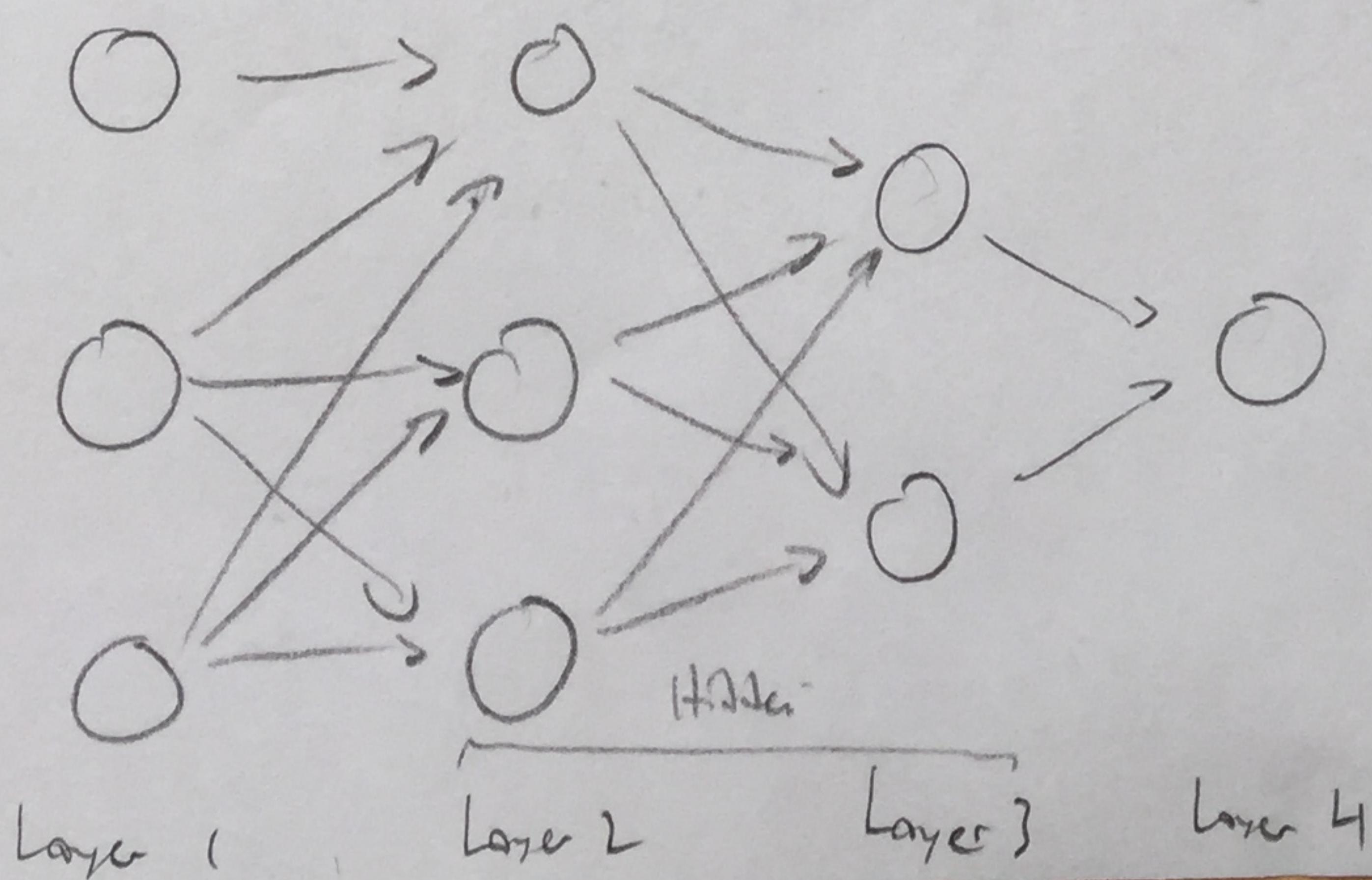
Vectorized implementation
of forward propagation

Note :



↳ just like logistic regression, but replaces x w/ $a^{(2)}$ → "learns its own features!"

Could have alternate architectures:



Vector Summary:

$$a_1^{(1)} = g(\Theta_{1,0}^{(1)}x_0 + \Theta_{1,1}^{(1)}x_1 + \Theta_{1,2}^{(1)}x_2 + \Theta_{1,3}^{(1)}x_3)$$

$$a_2^{(1)} = g(\Theta_{2,0}^{(1)}x_0 + \Theta_{2,1}^{(1)}x_1 + \Theta_{2,2}^{(1)}x_2 + \Theta_{2,3}^{(1)}x_3)$$

$$a_3^{(1)} = g(\dots)$$

$$h_\theta(x) = a_1^{(3)} = g(\Theta_{1,0}^{(2)}a_0^{(2)} + \Theta_{1,1}^{(2)}a_1^{(2)} + \Theta_{1,2}^{(2)}a_2^{(2)} + \Theta_{1,3}^{(2)}a_3^{(2)})$$

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

$$a_2^{(2)} = g(z_2^{(2)}) \Rightarrow z_k^{(2)} = \Theta_{k,0}^{(1)}x_0 + \Theta_{k,1}^{(1)}x_1 + \dots + \Theta_{k,n}^{(1)}x_n$$

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

$$\Rightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix} \xrightarrow{[(n+1) \times 1]}$$

$$\Rightarrow \text{if } x = a^{(1)} \text{ (first layer)} \rightarrow z^{(j)} = \Theta^{(j-1)} a^{(j-1)}$$

$$\xleftarrow{\quad} \quad \quad \downarrow \quad \quad [s_j \times 1] \quad [s_j \times (n+1)] \quad \text{where } s_j \# \text{ activation nodes}$$

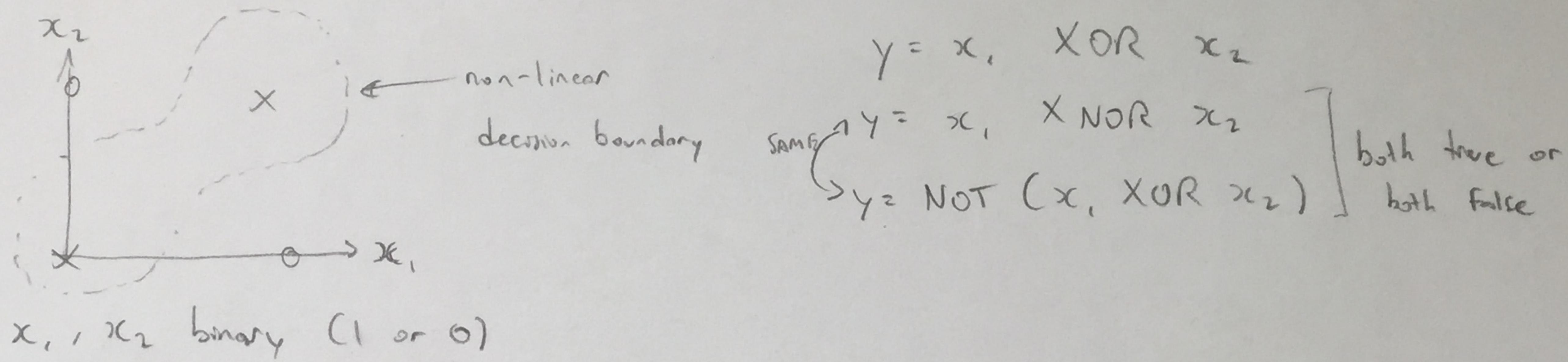
→ add bias unit → $a_0^{(j+1)}$ added to top of $a^{(j+1)}$ to make $[(n+1) \times 1]$

$$\Rightarrow z^{(j+1)} = \Theta^{(j)} a^{(j)}$$

$$\hookrightarrow \text{final } \Theta^{(j)} [1 \times (n+1)] \text{ times } a^{(j)} [(n+1) \times 1] \rightarrow [1 \times 1]$$

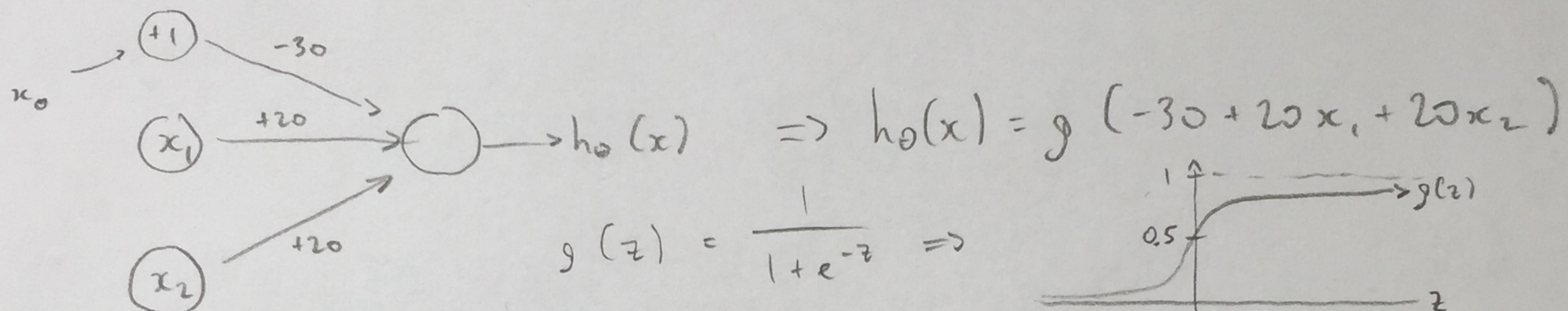
$$\Rightarrow h_\theta(x) = a^{(j+1)} = g(z^{(j+1)}) \leftarrow \text{same as logistical regression!}$$

Non-Linear Classification:



⇒ Neural Example: AND:

$$x_1, x_2 \in [0,1], y = x_1 \text{ AND } x_2:$$



$$\text{if } x_1 = x_2 = 0 \rightarrow g(-30) \approx 0$$

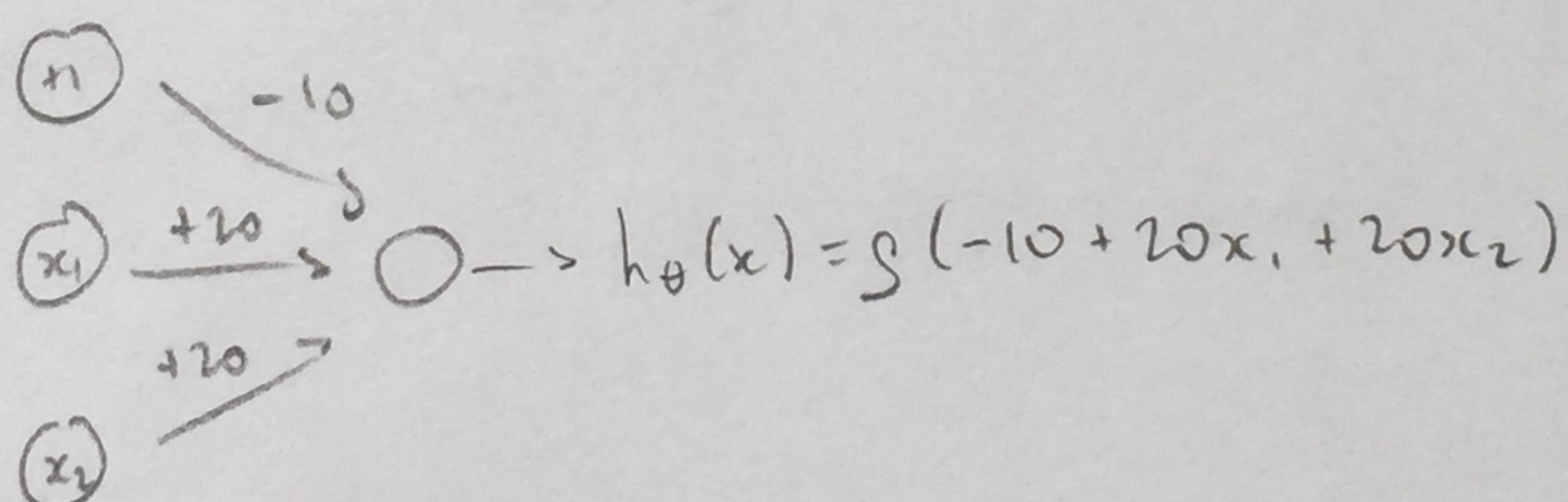
$$\text{if } x_1 = 0 \& x_2 = 1 \rightarrow g(-10) \approx 0$$

$$\text{if } x_1 = 1 \& x_2 = 0 \rightarrow g(10) \approx 1$$

$$\text{if } x_1 = 1 \& x_2 = 1 \rightarrow g(10) \approx 1$$

Logical AND function!

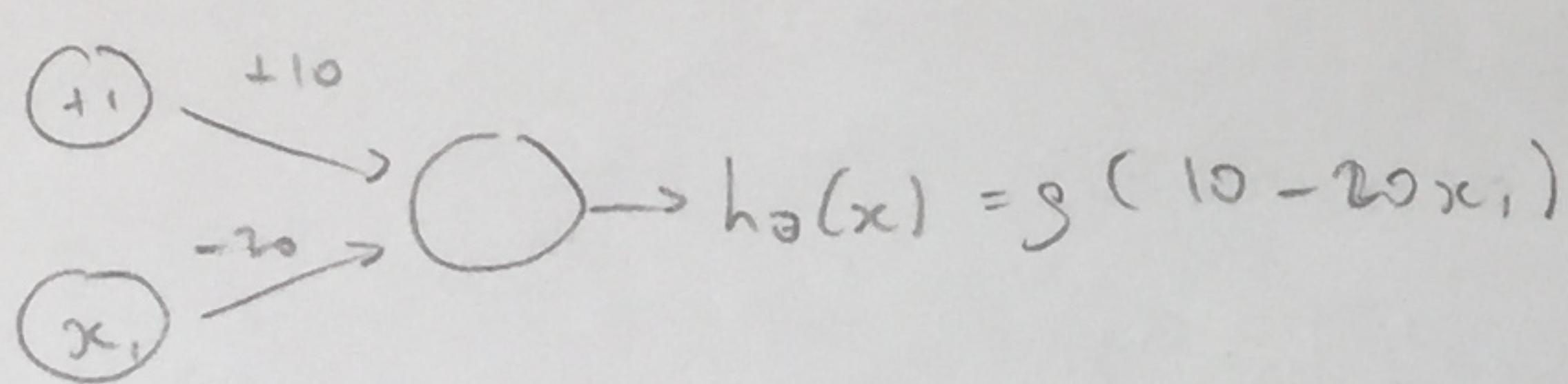
ex 2)



x_1	x_2	$h_0(x)$
0	0	0
0	1	1
1	0	1
1	1	1

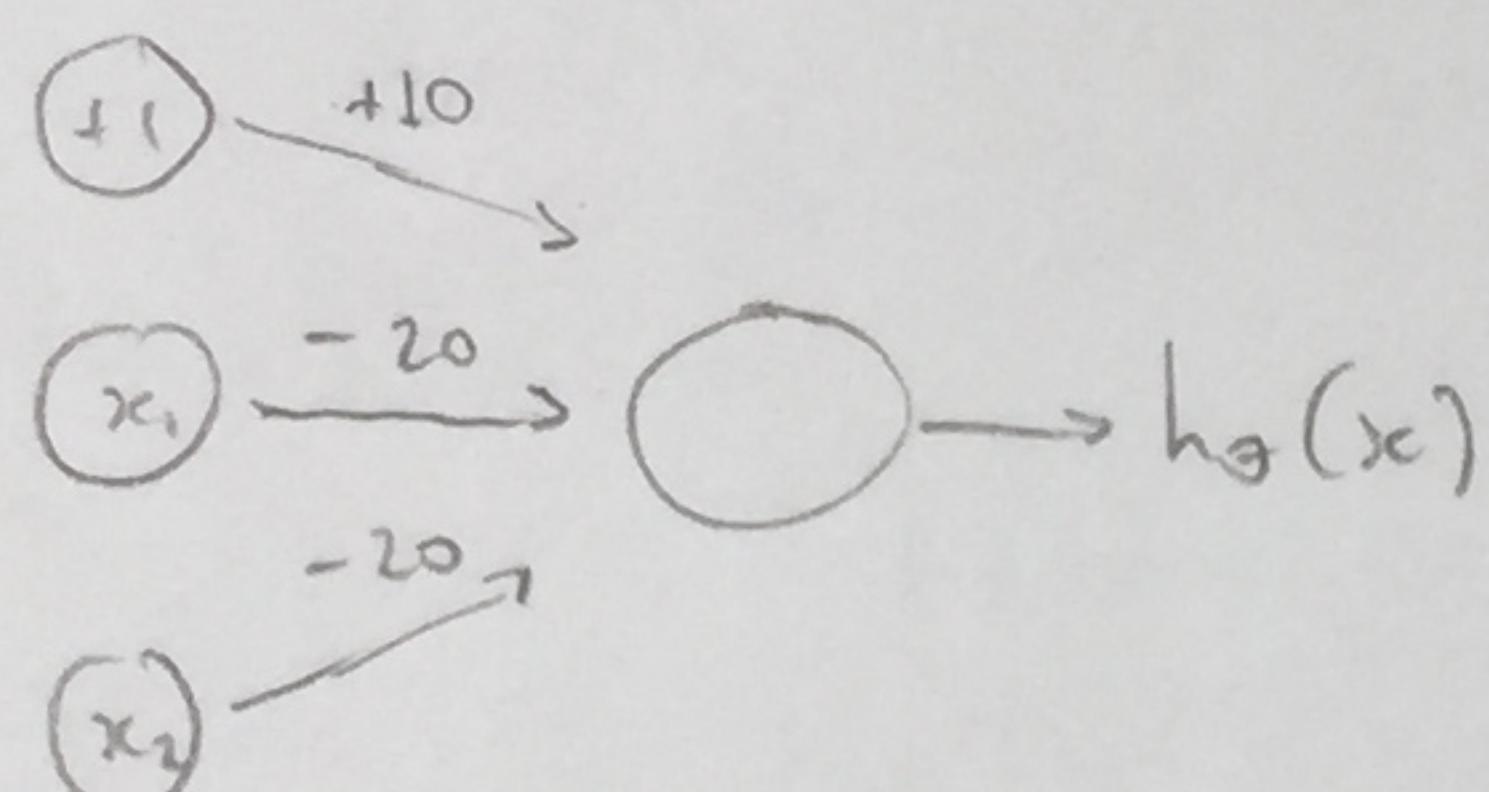
$x_1 \text{ OR } x_2$ logic

\Rightarrow Negation ($\text{NOT } x_1$):



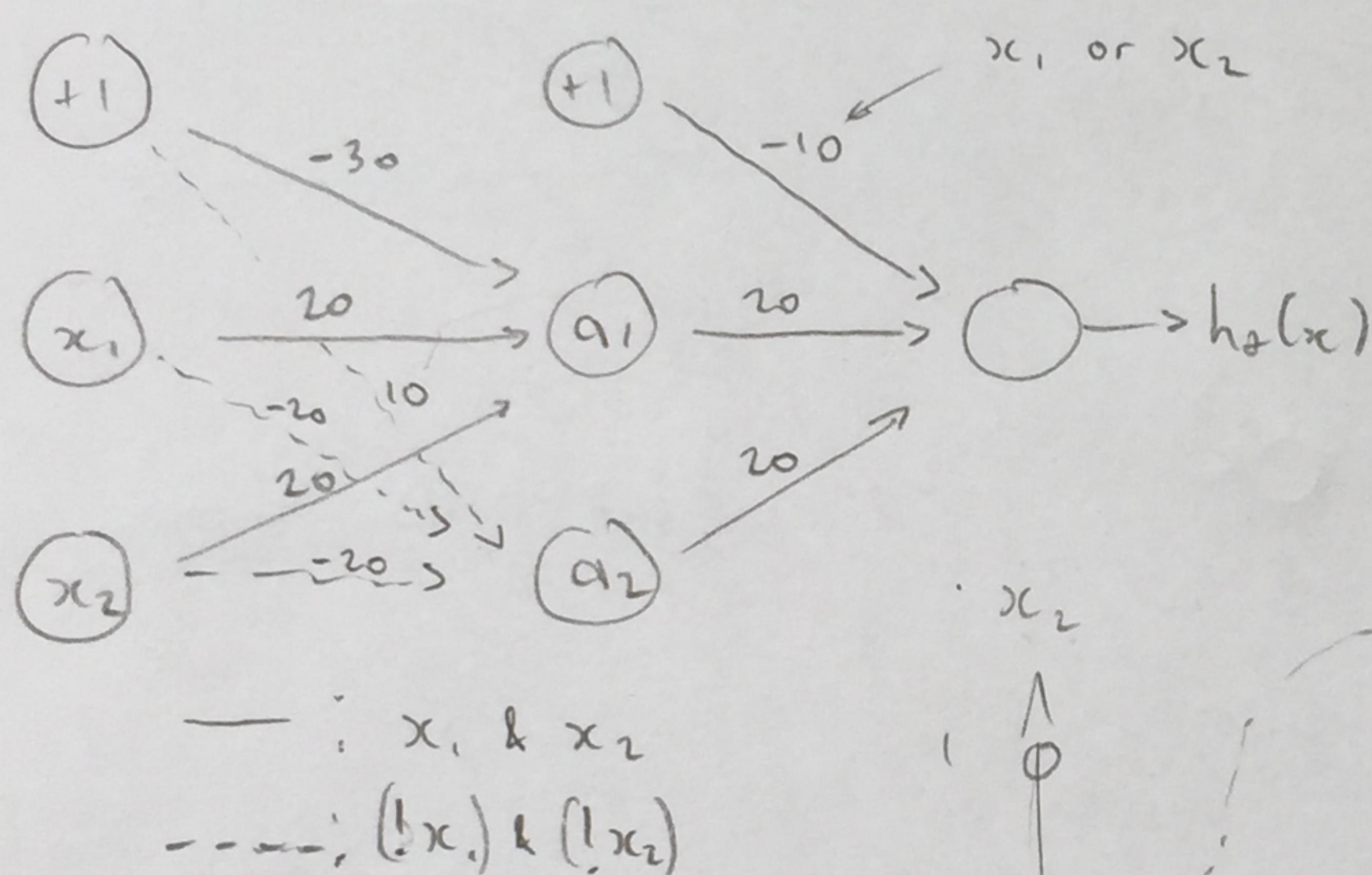
x_1	$h_\theta(x)$
0	1
1	0

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

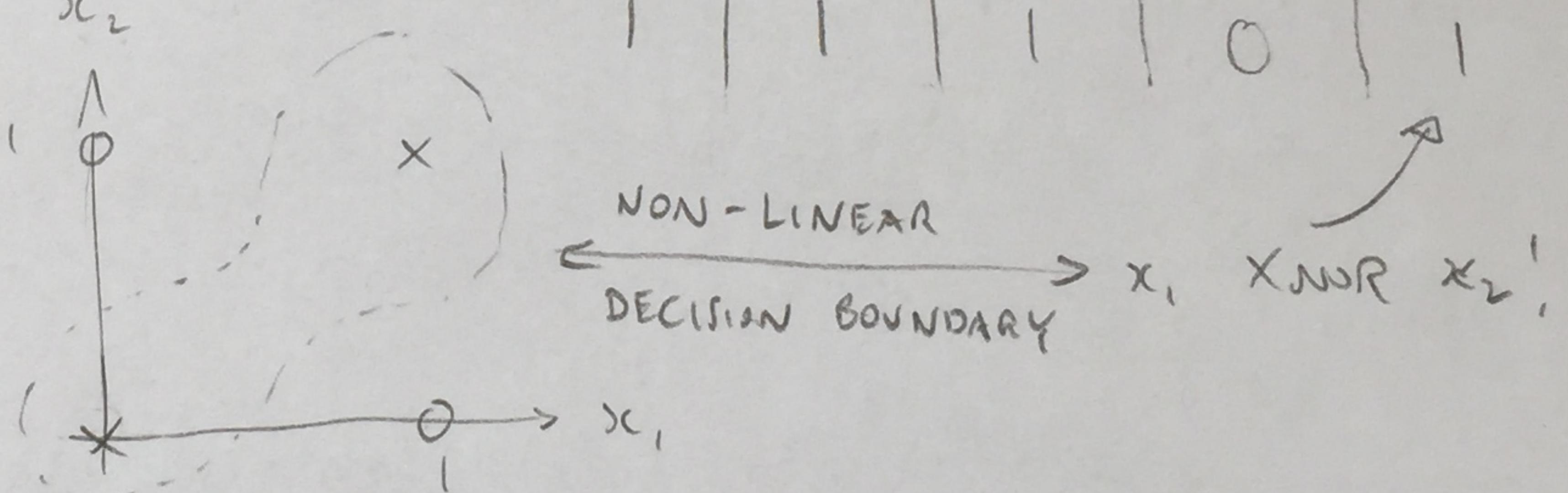


x_1	x_2	$h_\theta(x)$
0	0	1
0	1	0
1	0	0
1	1	0

To create $x_1 \text{ XNOR } x_2$ (both 0 or both 1):



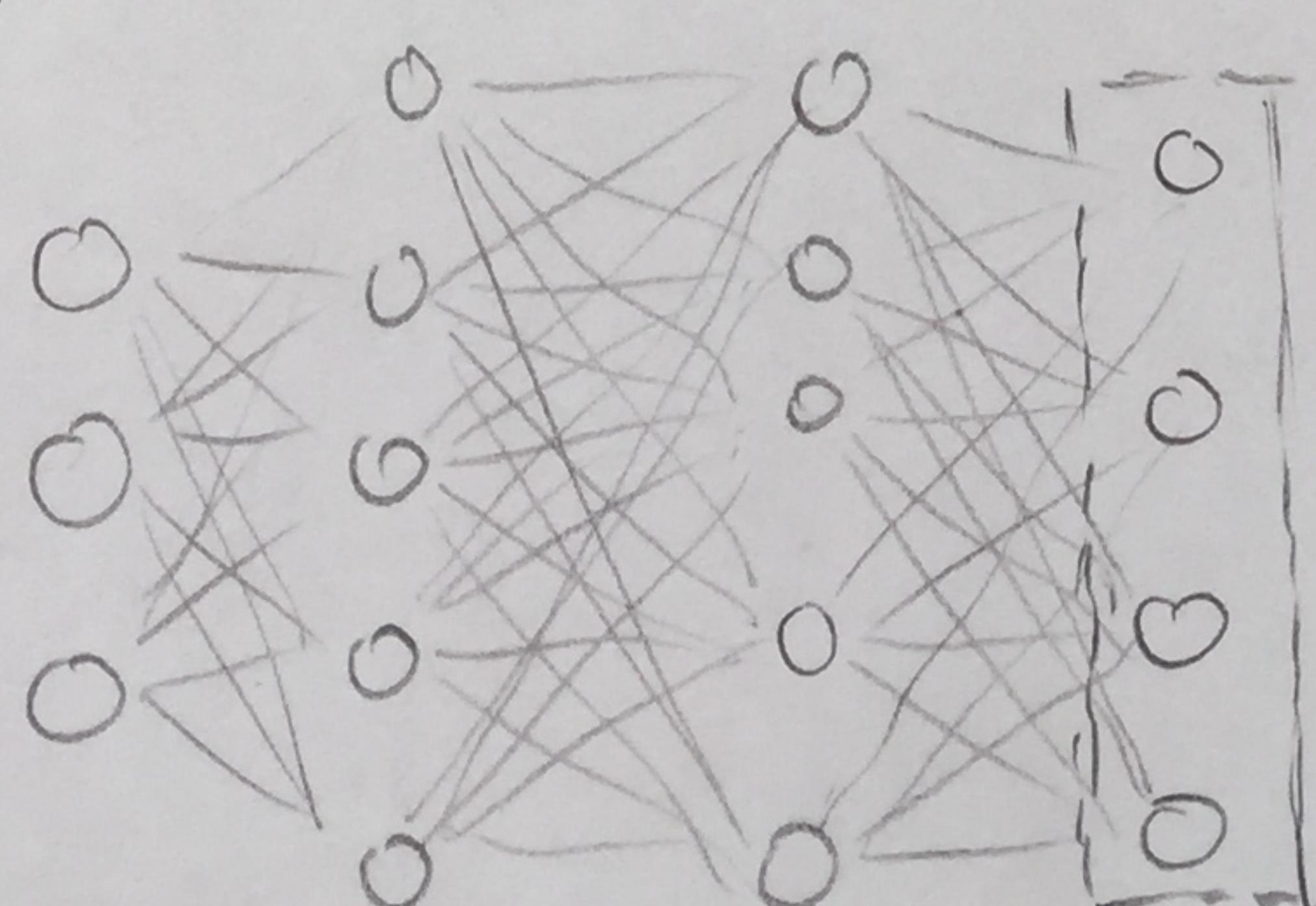
x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1



Multiclass Classification:

e.g. rather than 0/1, 0/1 \rightarrow digits from 0-9

e.g.: recognize car, pedestrian, motorcycle or truck:

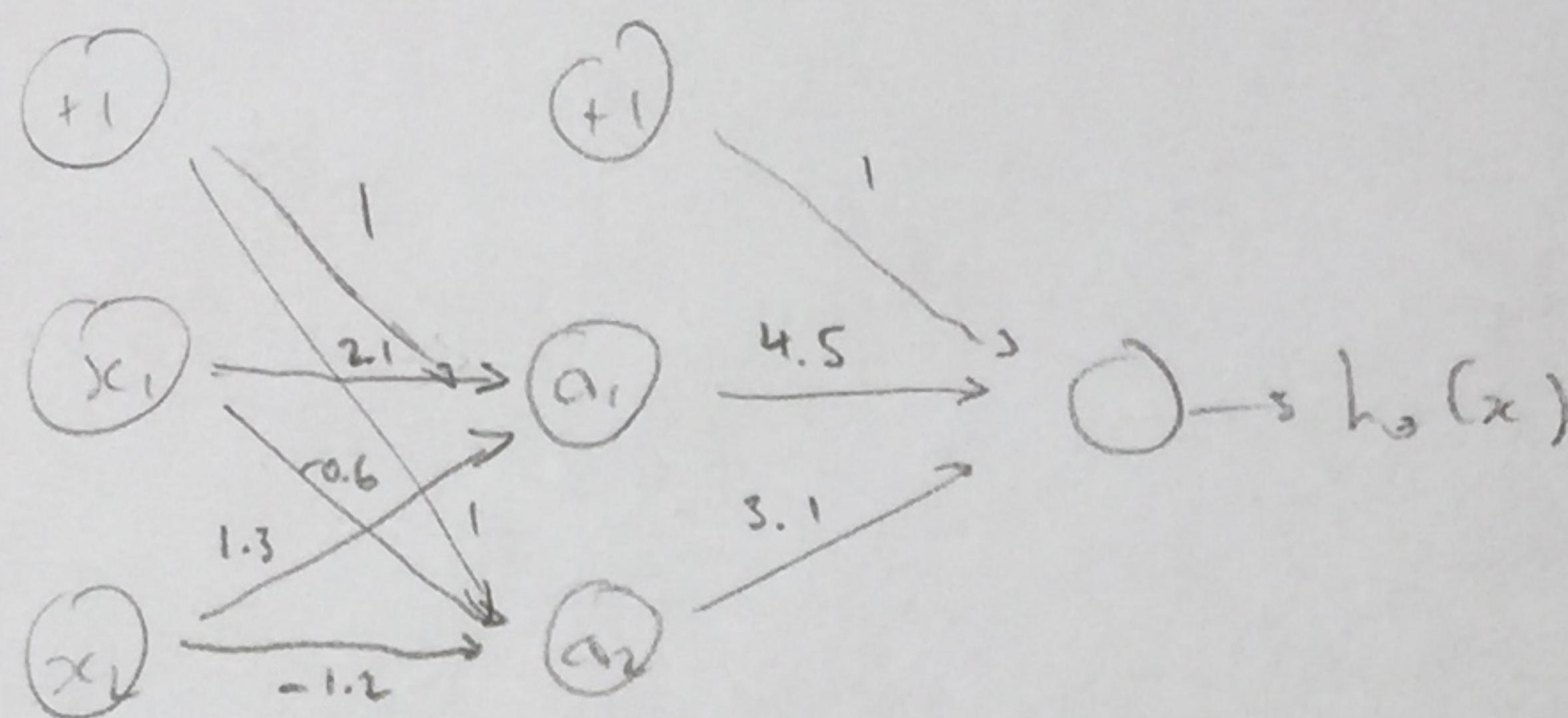


4 logistic regression classifiers:

outputs $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \dots$, etc.

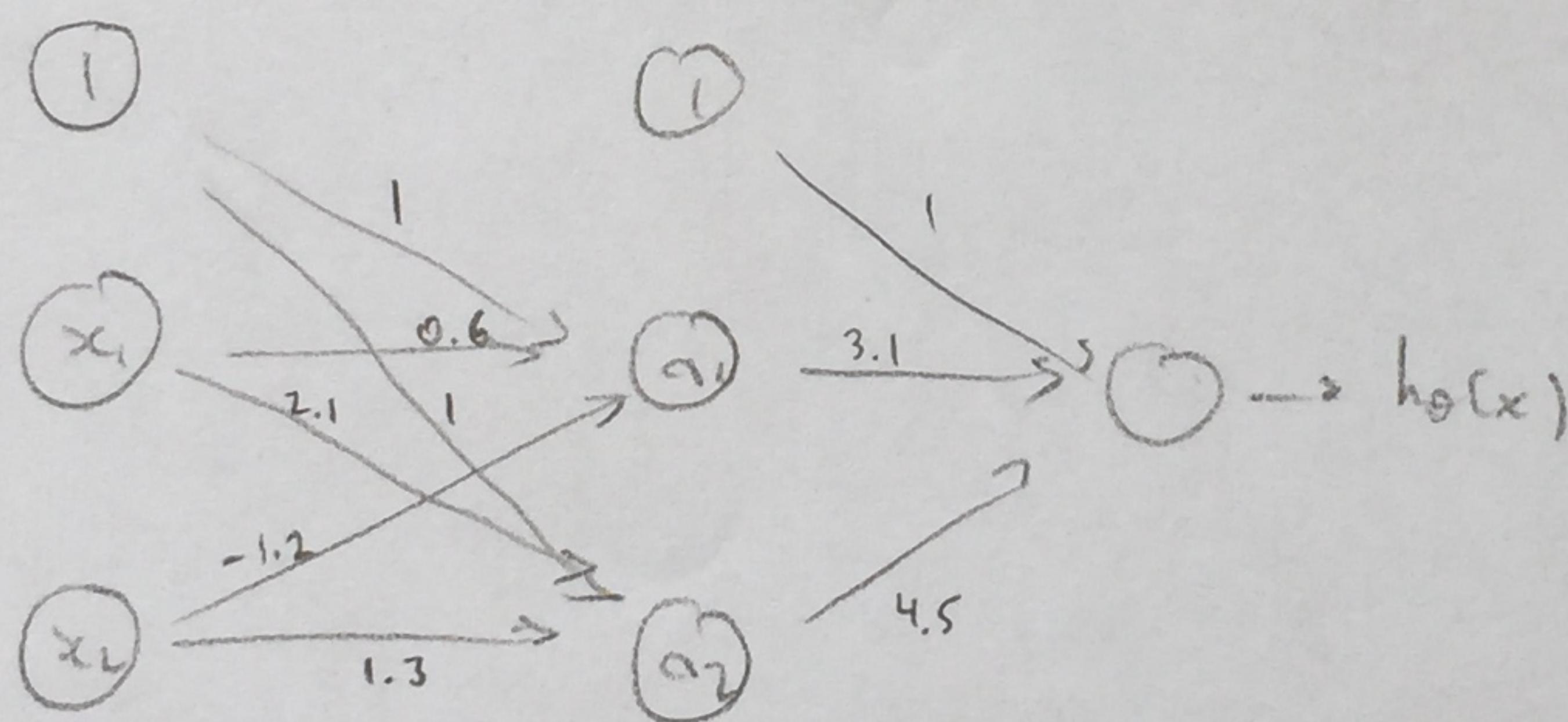
- Represent training set: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(n)}, y^{(n)})$
- $y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ (4 classes)
- $h_\theta(x^{(i)}) \approx y^{(i)}$ (hopefully)

QUIZ:



$$\left. \begin{aligned} a_1 &= g(1 + 2.1x_1 + 1.3x_2) \\ a_2 &= g(1 + 0.6x_1 - 1.2x_2) \end{aligned} \right\} \rightarrow h_\theta(x) = g(1 + 4.5a_1 + 3.1a_2)$$

VS.



$$\left. \begin{aligned} a_1 &= g(1 + 0.6x_1 - 1.2x_2) \\ a_2 &= g(1 + 2.1x_1 + 1.3x_2) \end{aligned} \right\} \rightarrow h_\theta(x) = g(1 + 3.1a_1 + 4.5a_2)$$

SAME