

⇒ Weeks - 4 : Neural networks representation

→ Why do we need neural networks?

* Says you have a complex supervised learning classification problem.

You can use logistic feature if it has 1-2 features BUT suff. if you have 100 features. It WILL NOT work. Here's why :

eg of Supp. m = training samples = 100
if $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{100} \end{bmatrix}$ ignoring x_0

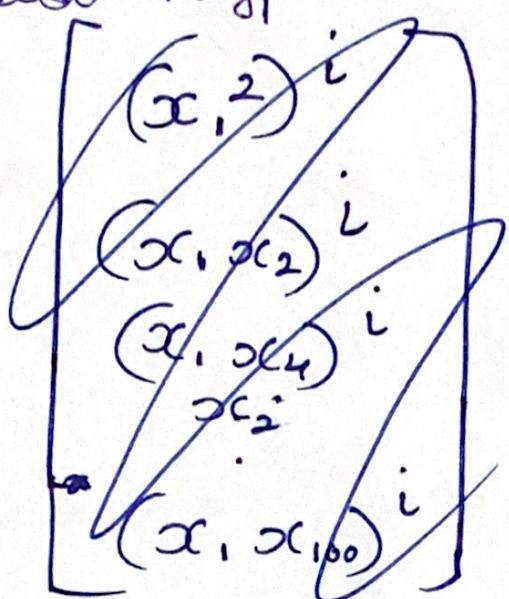
eg : Supp. m = no. of features = 100

then $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{100} \end{bmatrix}$

then where $x^i = \begin{bmatrix} x_0^i \\ x_1^i \\ \vdots \\ x_{100}^i \end{bmatrix}$

* Supp. to prevent underfitting , we increase our features & include all the quadratic terms (second order).
 in our hypothesis then :

~~this~~ ~~set~~



for n=100

$$x^i = \begin{bmatrix} (x_1^2)^i \\ x_1 x_2^i \\ \vdots \\ x_1 x_{100}^i \\ ((x_2)^2)^i \\ \vdots \\ (x_2 x_{100})^i \\ \vdots \\ ((x_{100})^2)^i \end{bmatrix}$$

Then our hypothesis will be :

$$h = g(\theta^T x)$$

$$= g(\theta_0 + \theta_1 x_1^2 + \theta_2 x_1 x_2 + \dots)$$

→ So as you can see how many features are added in our hypothesis.

* The exact way to calculate how many features for all polynomial terms is the combination function with repetition:

$$\binom{m+r-1}{r} = \frac{m+r-1}{r!}$$

$$= \frac{(m+r-1)!}{(r-1)! m!}$$

Refer NPTEL vid. for the formula

→ For 100 features (m), if we wanted to make them quadratic ($\binom{m+2-1}{2}$) we would get

$$\frac{(100+2-1)!}{(2-1)! (100-1)!} = 5050$$

resulting new features! We can approximate this as

→ We

→ Suppose for 100 features, we wanted

→ We can approximate the growth of the no. of new features we get with all quadratic terms with $O(m^2/2)$.

→ As our decision boundary gets more & more complex, the no. of inputs in our hypothesis will increase which would take

Now, suppose we want to include all the cubic terms in our hypothesis, no. of

features $\propto = \frac{100+3-1}{13 \cdot 100-1} = 171700$, the

features grow asymptotically at $O(n^3)$.

These are very steep growths,

so as the no. of our features increase, the no. of quadratic or cubic features^{in our hypothesis} increases very rapidly & becomes quickly impractical.

So, as our hypothesis becomes more & more complex

→ Neural networks offer an alternate way to perform machine learning when we have complex hypothesis with many features.

As ~~this~~ we ^{wouldn't} need

to "add" more ~~for~~ inputs in our hypothesis just to ^{satisfy} increase its^(hypothesis) complexity.

⇒ Neurons & the Brain

* Neural nets are limited imitations of how our brain works! They've had a big recent resurgence because of advances in computer hardware.



⇒ There is evidence that the brain uses only one "learning algorithm" for all its different functions.

→ Evidence for this hypothesis:

* Auditory cortex in brain: Responsible for hearing in humans

- If you cut the wiring from the ear to the auditory cortex, it re-route optic nerve to the auditory cortex.

- The auditory cortex learns to see !!!!!

* Somatosensory cortex: Responsible for touch processing

- If you rewire optic nerve to the somatosensory cortex, then it learns to see.

* The principle is called "neuroplasticity" & has many examples & experimental evidence.

→ ~~e.g.~~ More examples:

* Brainport: Seeing with your ~~tongue~~ tongue!

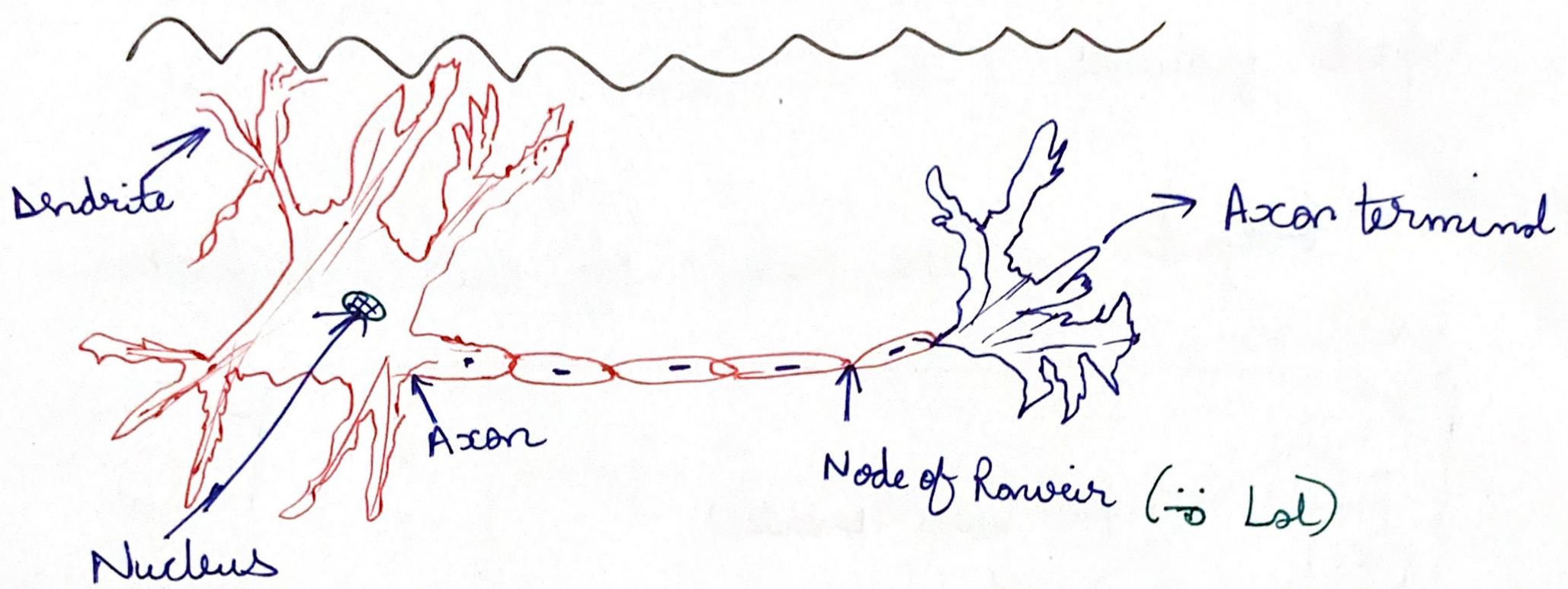
* Human echolocation: Blind people trained to interpret sound & echo (sonar).



⇒ Model representation



→ What does a neuron look like?



→ Working:

* Neuron gets one or more inputs through dendrites.

* Does some processing

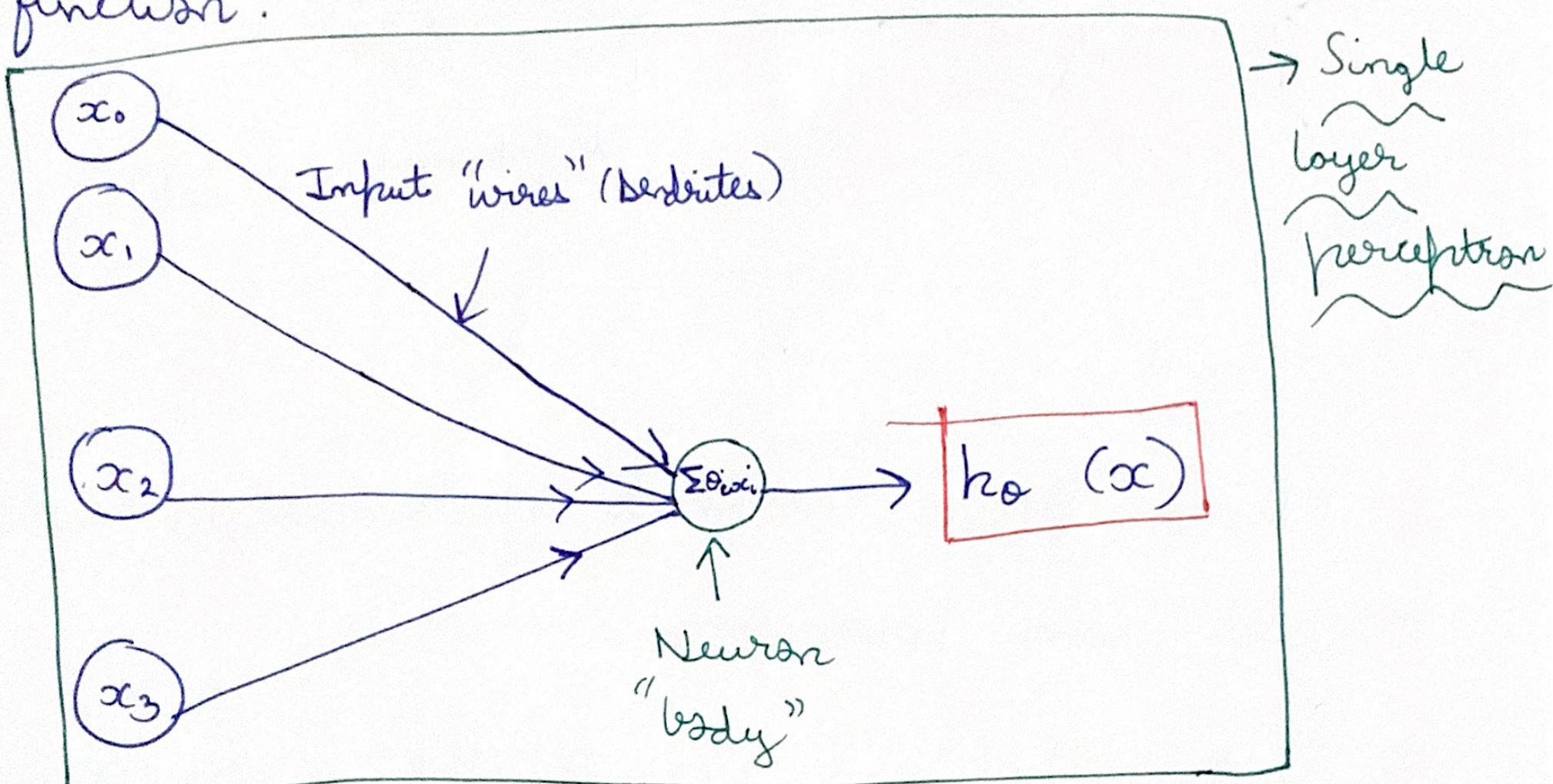
* Sends output down axon to other neuron(s).

→ Neurons communicate through electric spikes.



→ Artificial neural network - Representation of a neuron

- * In our model, our dendrites are like the input features ($x_1 \dots x_m$).
- * The output is the result of our hypothesis function.



- * $\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$; In our this model our x_0 input node is sometimes called the "bias unit". It is always equal to 1.

- * In this neural nets, we use the some logistic function as in classification :
$$g(z) = \frac{1}{1+e^{-z}}$$
. In neural networks, we sometimes call it a sigmoid (logistic) activation function.

* The parameter $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$ we discussed before are called "weights" in the neural networks model.

* A simplistic representation:

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow [] \rightarrow h_{\theta}(x)$$

\Rightarrow Representing a neural net : Neural net is just a group of neurons "strung" together.

* Our inputs nodes (layer 1) go into another node (layer 2), & are output as the hypothesis function.

* Layer-1 : Input layer

Layer-3 : Output layer

We can have intermediate layers of nodes b/w the input & output layers called the "hidden layer".

Notation for neurons:

Well

* $a_i^s \equiv$ activation of unit i in layer s

$\rightarrow a_1^2 \equiv$ activation of unit 1 in layer 2

\rightarrow Activation : By activation, we mean the value which is computed as output by the node

* $\theta^s \equiv$ matrix of weights controlling function mapping from layer(s) to layer ($s+1$)

\rightarrow Dimension of θ^s :

If network has (s_j) units in layer(j)

$s(s_{j+1})$ units in layer ($j+1$), then θ^s will

be of dimension $(s_{j+1}) \times (s_j + 1)$ ~~(s_j)~~

The ' $+1$ ' comes from the addition in θ^s of the "bias nodes"

$x_0 \theta_0^s$. In other words, the output

nodes will not include the bias nodes

→ Looking at the δ matrix

* Column length \rightarrow no. of units in the foll. layer

* Row length \rightarrow no. of units in the current layer + 1
(because we have to map the bias unit)

→ What are the computations which occur?

* We have to calculate the activation for each node. (The $g(\delta^T x)$ term)

⊗

* That activation depends on :

1) the inputs to the node (from the prev. layer)

2) the parameters associated with the node

(from the δ vector associated with that layer)

⇒ Note: Every input/activation goes to every node in following layer.

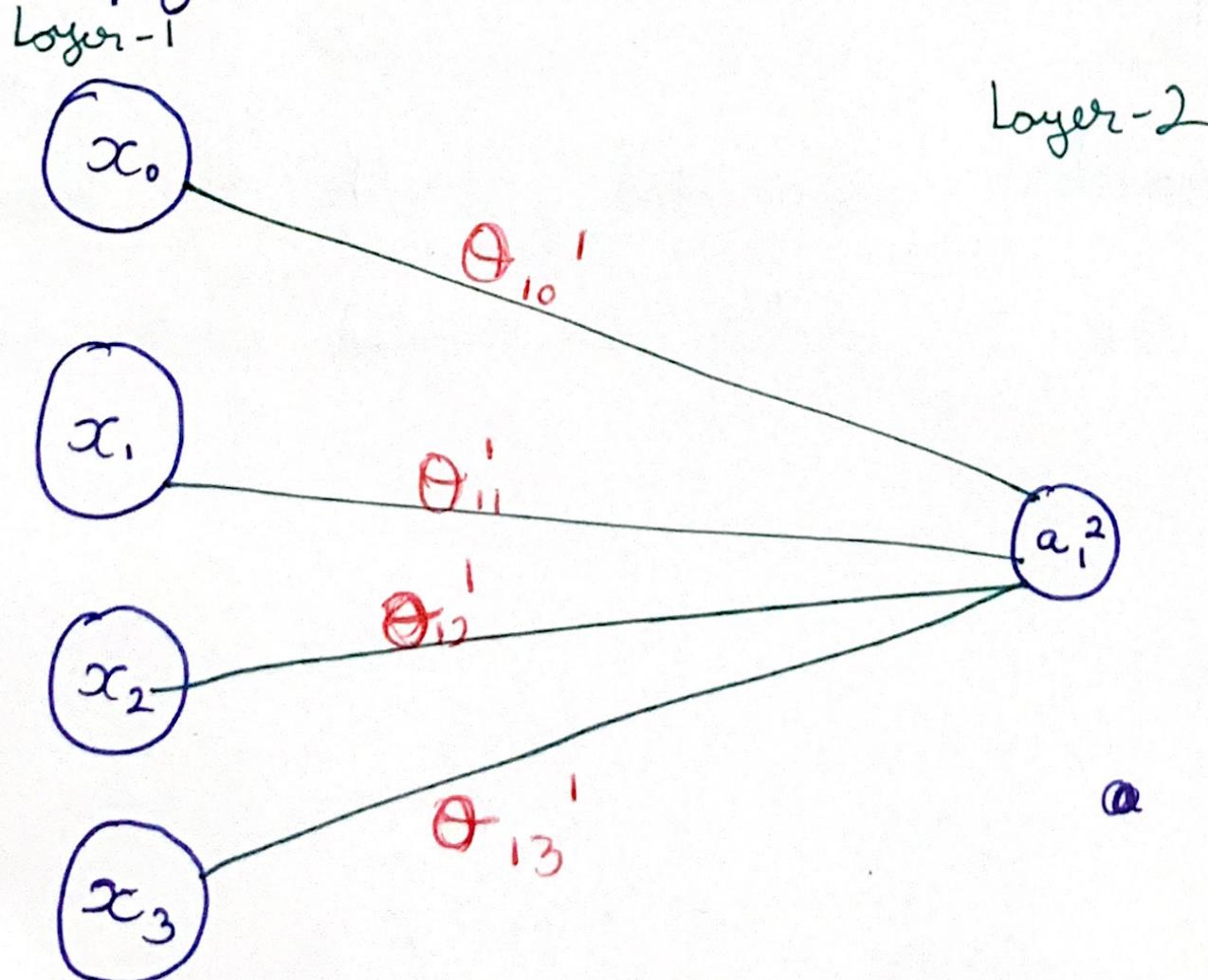
$$\rightarrow \theta_{ab}^c$$

* a → Ranges from 1 to the no. of units in layer (c+1)

* b → Ranges from 0 to the no. of units in layer (c)

* c → layer which you're moving FROM.

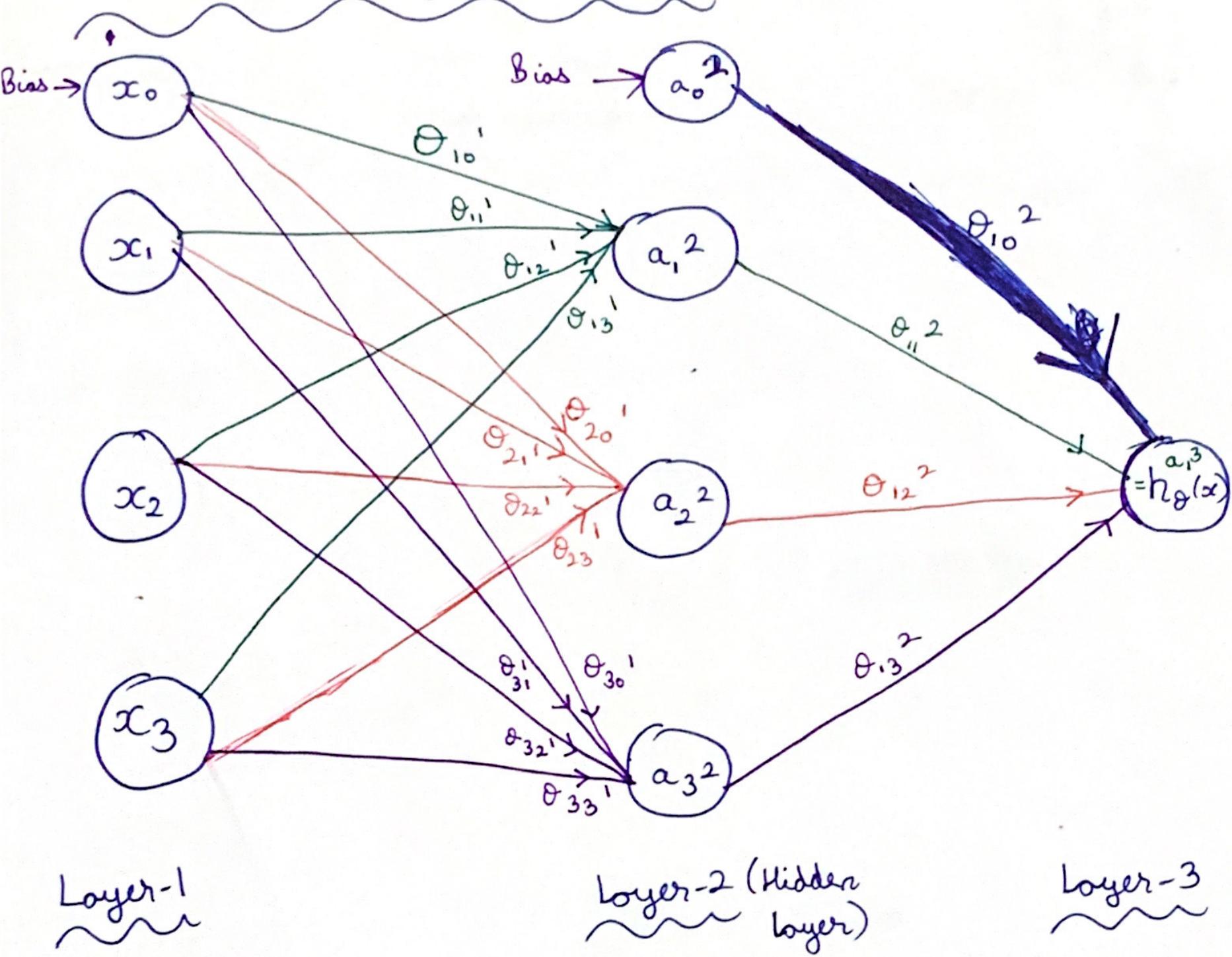
e.g.: This foll. single layered perceptron will be helpful:



$$\text{Activation of 1st unit} = a_1^2 = g(\theta_{10}^1 x_0 + \theta_{11}^1 x_1 + \theta_{12}^1 x_2 + \theta_{13}^1 x_3)$$

at Layer 2

⇒ An example of neural net



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

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

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

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

$$* \theta^1 = \begin{bmatrix} \theta_{10}^1 & \theta_{11}^1 & \theta_{12}^1 & \theta_{13}^1 \\ \theta_{20}^1 & \theta_{21}^1 & \theta_{22}^1 & \theta_{23}^1 \\ \theta_{30}^1 & \theta_{31}^1 & \theta_{32}^1 & \theta_{33}^1 \end{bmatrix} \rightarrow \begin{array}{l} \text{Maps layer-1 to layer-2} \\ \downarrow \quad \downarrow \\ (3 \text{ input nodes}) \quad (3 \text{ activation nodes}) \\ + 1 \text{ bias } (x_0) \end{array}$$

$(3) \times (3+1)$

$$\theta_2 = \begin{bmatrix} \theta_{10}^2 \\ \theta_{11}^2 \\ \theta_{12}^2 \\ \theta_{13}^2 \end{bmatrix} \quad 1 \times (3+1)$$

→ Maps layer 2 → layer 3

↓
(3 input nodes)
+ 1 bias (a_0^2)

↓
1 activation node

* we can also set $x = a'$

$$\therefore x = a' = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix}$$



~~Also~~: $a^{\delta^{-1}} = \begin{bmatrix} a_0^{\delta^{-1}} \\ a_1^{\delta^{-1}} \\ \vdots \\ a_m^{\delta^{-1}} \end{bmatrix}$ ~~(m+1) × 1~~

Also, $\theta^{\delta^{-1}} = \begin{bmatrix} \theta_{1,0}^{\delta^{-1}} & \theta_{1,1}^{\delta^{-1}} & \theta_{1,2}^{\delta^{-1}} & \dots & \theta_{1,\delta_g-1}^{\delta^{-1}} \\ \theta_{2,0}^{\delta^{-1}} & \theta_{2,1}^{\delta^{-1}} & \theta_{2,2}^{\delta^{-1}} & \dots & \theta_{2,\delta_g-1}^{\delta^{-1}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \theta_{\delta_g,0}^{\delta^{-1}} & \theta_{\delta_g,1}^{\delta^{-1}} & \theta_{\delta_g,2}^{\delta^{-1}} & \dots & \theta_{\delta_g,\delta_g-1}^{\delta^{-1}} \end{bmatrix}$

$$\begin{aligned} \rightarrow \text{size}(\theta^{\delta^{-1}}) &= (\delta_g) \times (\delta_{g-1} + 1) \\ &= (\delta_g) \times (m+1) \end{aligned}$$

We can write $\theta^{\delta^{-1}} = a$

$$\theta^{\delta^{-1}} = a_0^{\delta^{-1}} A + a_1^{\delta^{-1}} A^2 + \dots + a_{\delta_g-1}^{\delta^{-1}} A^{\delta_g-1}$$

\Rightarrow Vectorized implementation of a neural net

- * Define a new variable z_k^s that encompasses the parameters inside our g function.

If we replace the variable z for all the parameters, we get in our previous example :

$$\begin{aligned}a_1^2 &= g(z_1^2) \\a_2^2 &= g(z_2^2) \\a_3^2 &= g(z_3^2)\end{aligned}$$

\therefore For layer $j = 2$ & node unit $= k$, the variable z will be :

$$z_k^2 = \theta_{k,0}^T x_0 + \theta_{k,1}^T x_1 + \dots + \theta_{k,m}^T x_m$$

- * Vectorizing x & z^s :

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix}, z^s = \begin{bmatrix} z_1^s \\ z_2^s \\ \vdots \\ z_m^s \end{bmatrix}$$

Vector of z values from the s th layer

We can rewrite the eqⁿ as :

$$z^s = \theta^{s-1} a^{s-1} ; z^s \rightarrow [\Delta_s \times 1] \text{ height vector}$$

* Now we can get a vector of our activation nodes for layer s as follows :

$$a^s = g(z^s) ; \text{ here our func}^n g \text{ can be applied element wise to our vector } z^s$$

* We can then add a bias unit to layer s after we've computed a^s . This will be element a_0^s 's will be equal to 1.

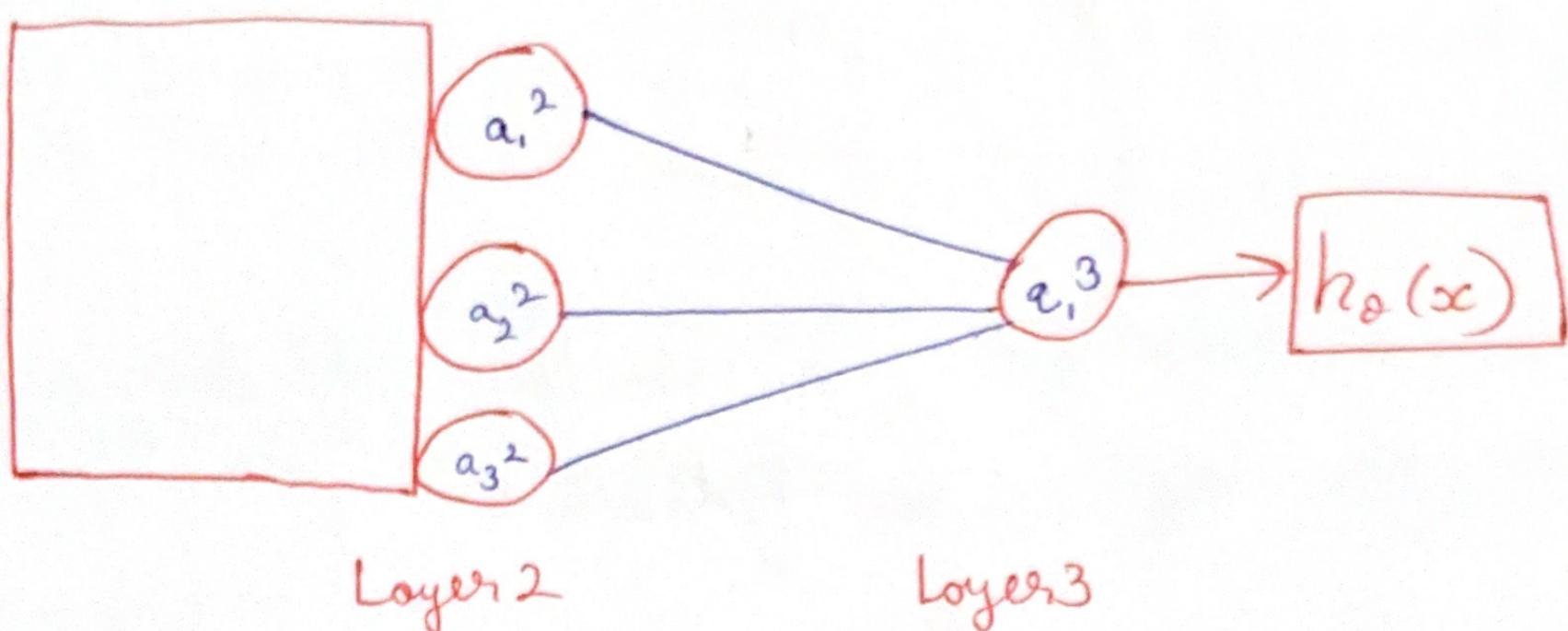
* Continuing this sequentially, layer by layer we'll finally we can finally calculate our hypothesis funcⁿ as

$$h_{\theta}(x) = a^{s+1} = g(z^{s+1}) = g(\theta^s a^s)$$

(Assuming our network has $(s+1)$ layers in total)

This process is called forward propagation. (input + hidden + output)

⇒ Neural network learning its own features



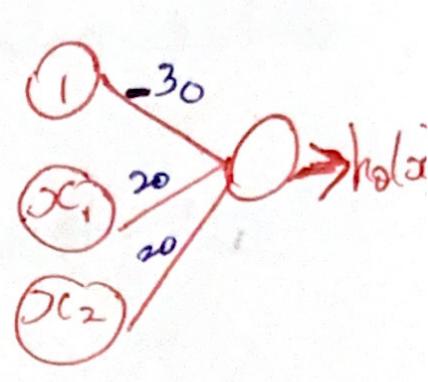
* Notice, that in this last step, between Layer-3 & Layer-2) we are doing exactly the same thing as we did in logistic regression!

The only difference is, instead of inputting a input vector, the features are just values calculated by the hidden layer (δ).

* This gives so, we basically input our values & let the hidden layers learn whatever gives the best final result to feed into the output layer **!!!**

\Rightarrow Neural network example : AND function

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow [g(z^2)] \rightarrow h_0(x)$$



x_0 = bias unit = 1

* Set our first theta matrix as:

$$\theta^1 = [-30 \quad 20 \quad 20]$$

* This will cause the output of our hypothesis to only be +ve if both x_1 & x_2 are 1. In other words:

$$h_0(x) = g(-30 + 20x_1 + 20x_2)$$

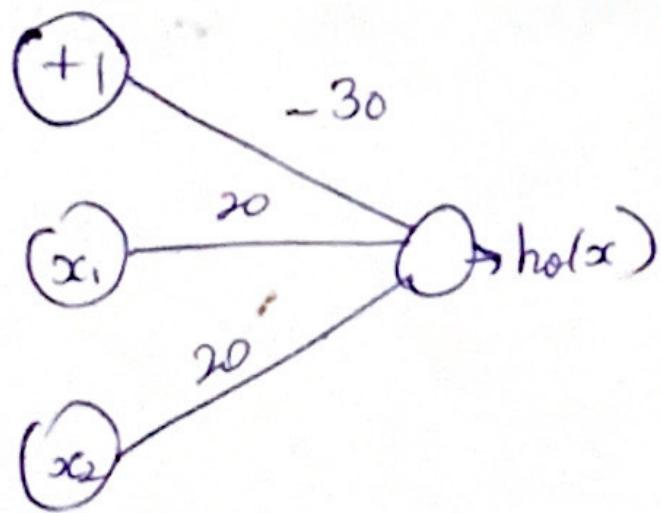
\rightarrow Side note : $g(4) \approx 1$; $g(-4.6) \approx 0$

\rightarrow Table :	x_1	x_2	$h_0(x)$
	0	0	$g(-30) \approx 0$
	0	1	$g(-10) \approx 0$
	1	0	$g(-10) \approx 0$
	1	1	$g(10) \approx 1$

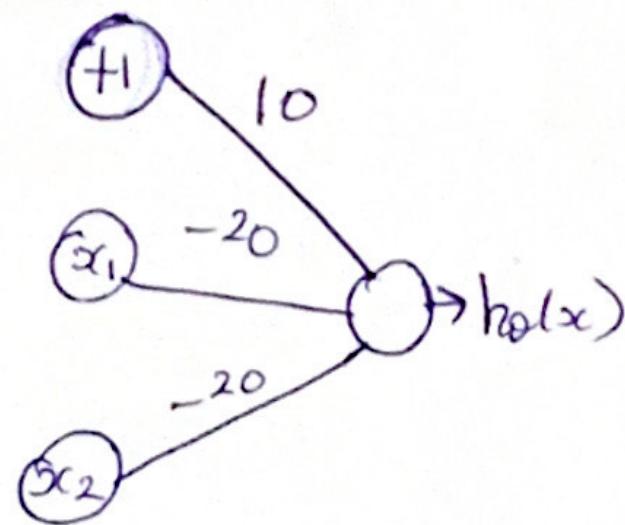
$$\text{NOR} : \theta^1 = [10 \quad -20 \quad -20]$$

$$\text{OR} : \theta^1 = [-10 \quad 20 \quad 20]$$

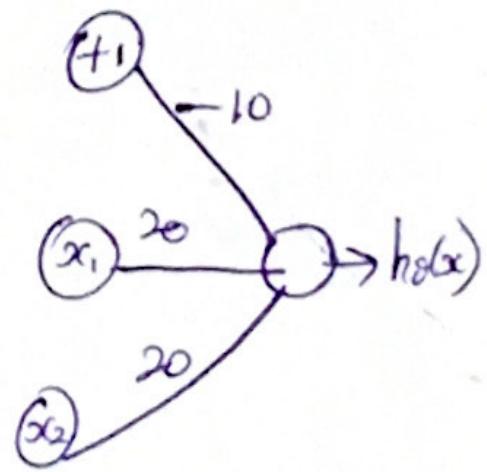
* We combine these 3 into a XOR neural network as shown below:



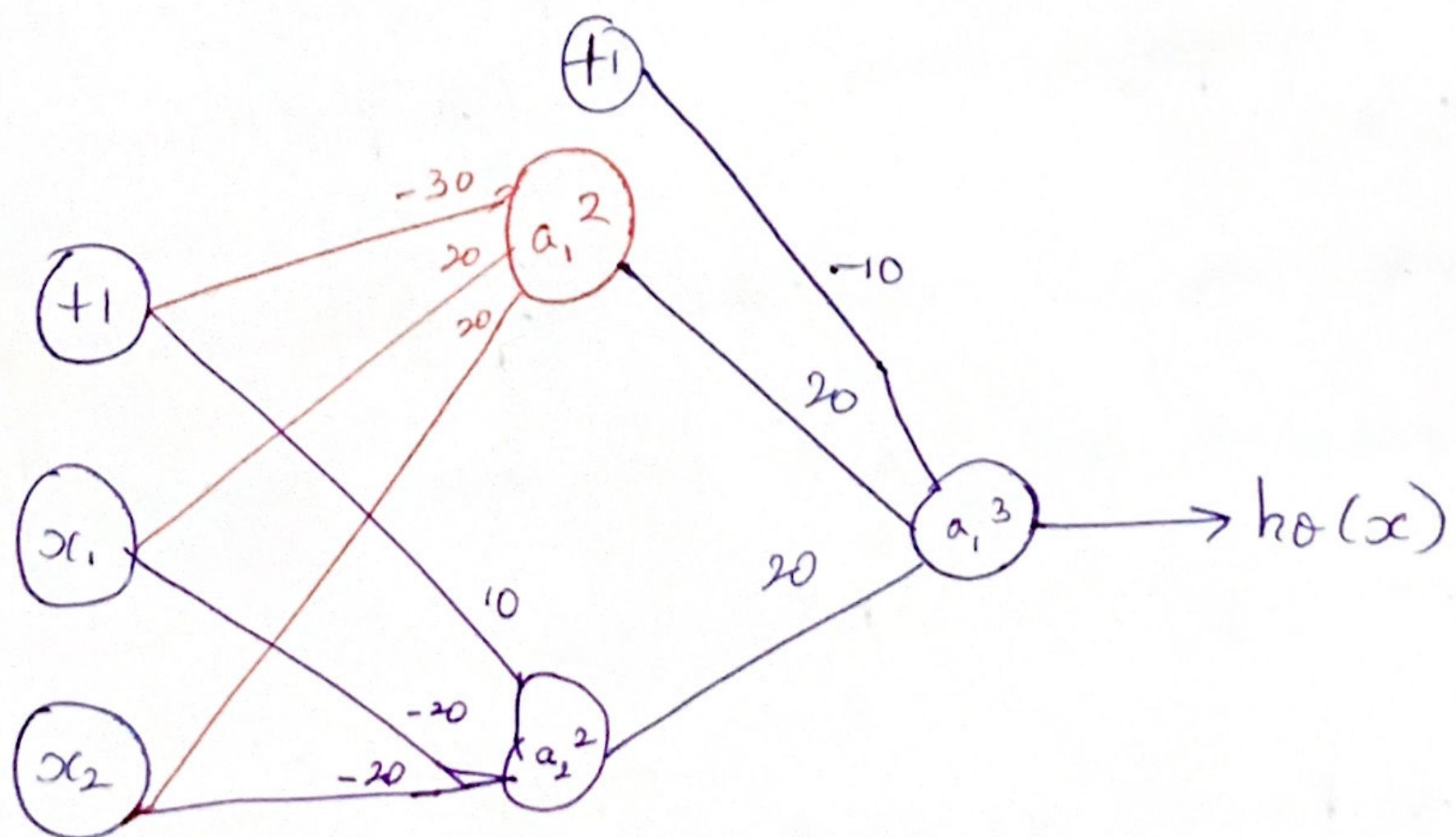
x_1 AND x_2



x_1 XNOR x_2

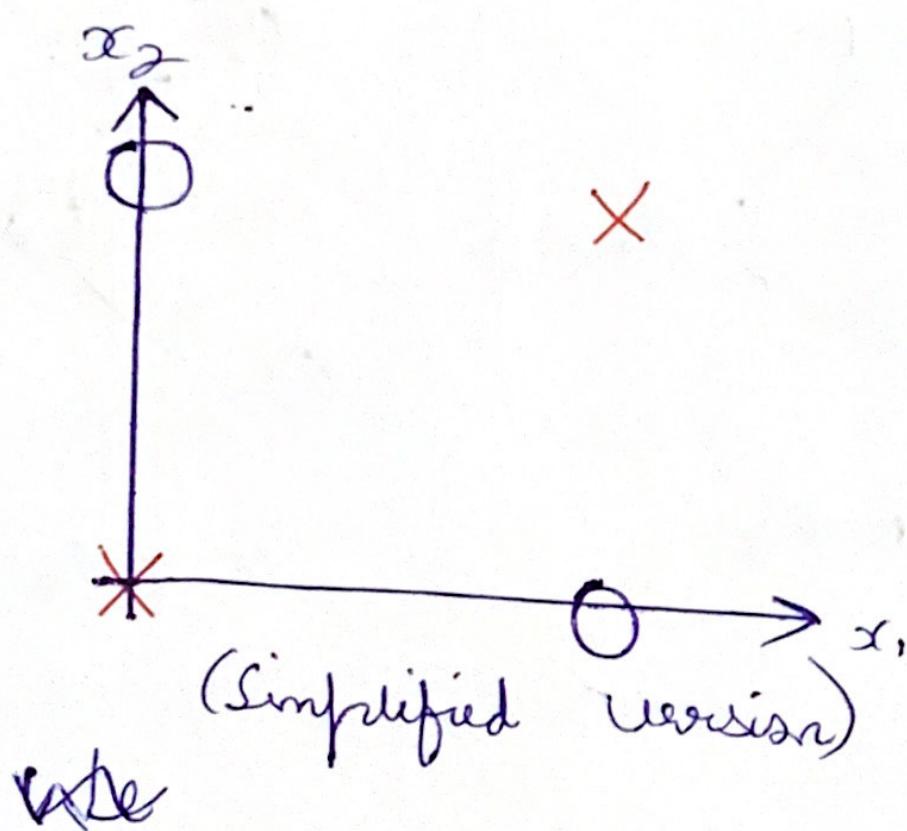


x_1 OR x_2



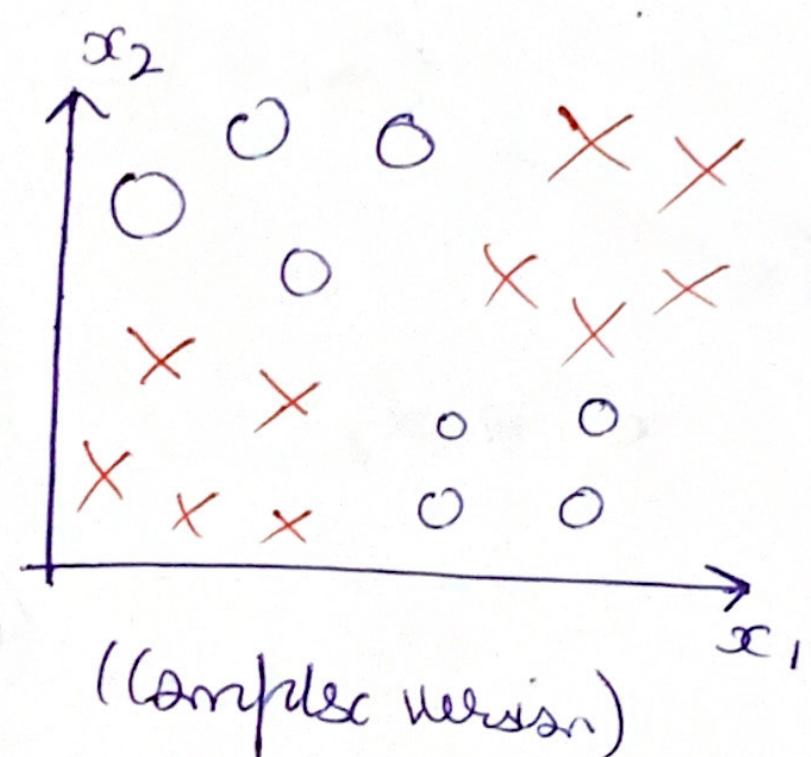
=/ Putting it together: $x_1 \text{ XNOR } x_2$

* Supr. we want to classify the foll. dataset:



(Simplified version)

note



(Complex version)

We can see in the simplified version

We want to find a non linear decision boundary

$$\text{s.t. } y = x_1 \text{ XNOR } x_2$$

* We know $x_1 \text{ XNOR } x_2 =$

$$x_1 \cdot x_2 + \overline{x_1} \cdot \overline{x_2}$$

$\underbrace{\quad\quad\quad}_{\text{AND}}$ $\underbrace{\quad\quad\quad}_{\text{NOR}}$

OR

for
→ ~~the~~ the prev. sample,

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

$$; a^1 = \text{oc} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

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

$$a^2 = g(\Theta^1 \text{oc})$$

$$a^3 = g(\Theta^2 a^2)$$

$$h_\Theta(\text{oc}) = a^3$$

And there we have the XNOR operator
using 2 hidden layers!

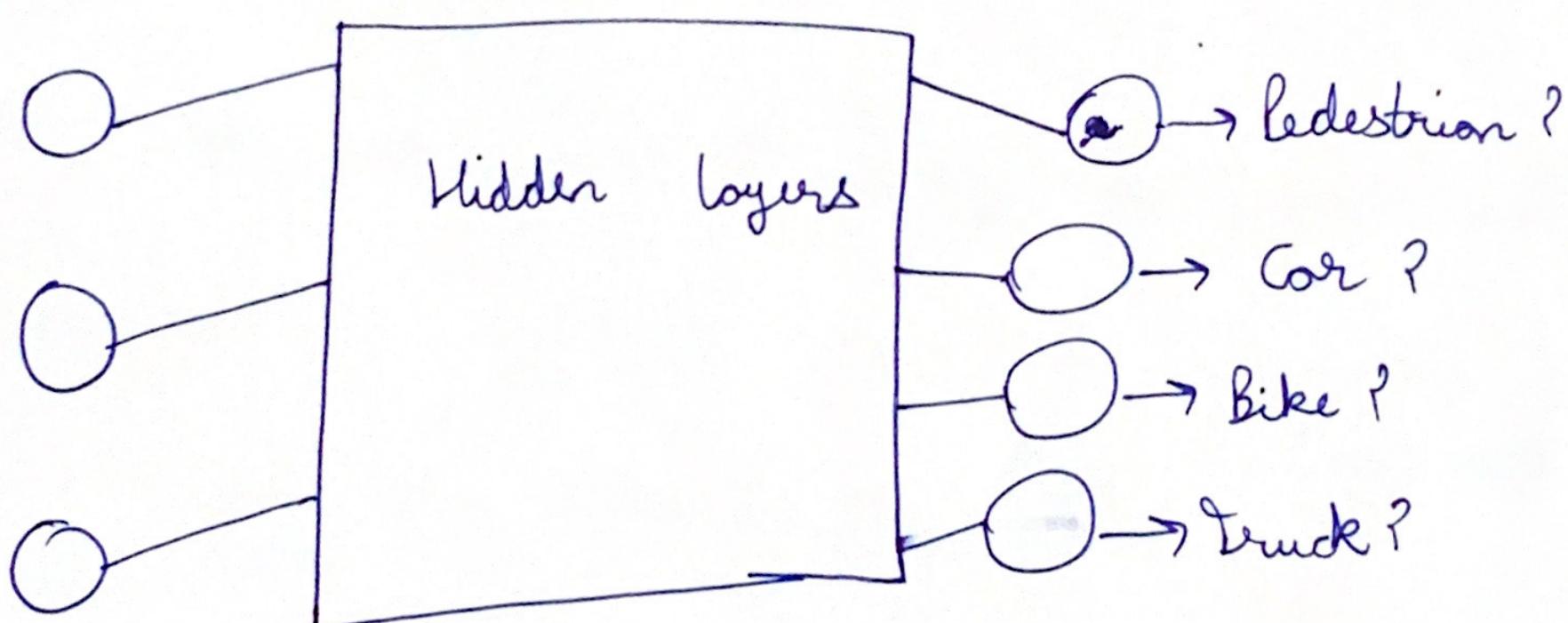
Imagine if we don't even have to
"prespecify" the Θ 's matrices before hand.

That's the power of neural nets! 🎉

\Rightarrow Multiclass classification

* To classify data into multiple classes, we let our hypothesis function return a vector of values.

\Rightarrow Multiple output units : one vs all



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

Want

$$h_{\theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \text{ when pedestrian}$$

$$h_{\theta}(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \text{ when car}$$

$$h_{\theta}(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \text{ when bike}$$

, etc.

We can define our set of resulting classes as y :

$$y^i = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

- * Each y^i represents a diff. image corresponding to either a car, pedestrian, truck, or bike.

The ~~are~~ inner layers, each provide us with ~~new~~ some new information which leads to our final hypothesis function.

- * The setup looks like:

