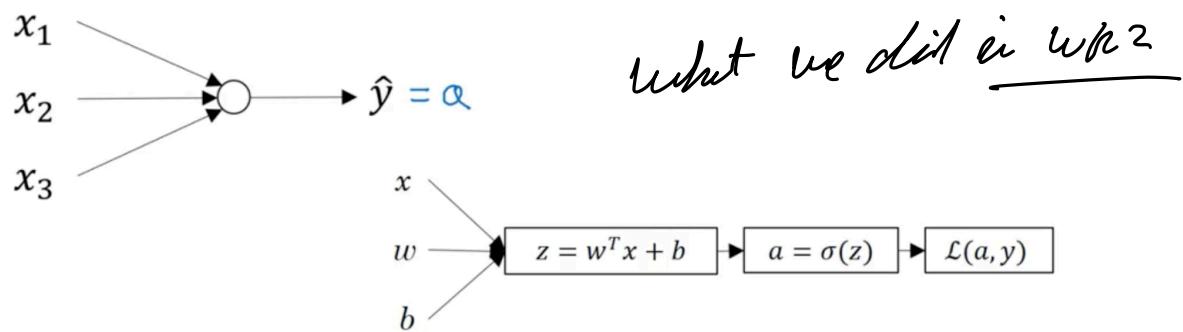


# Overview of Neural Networks

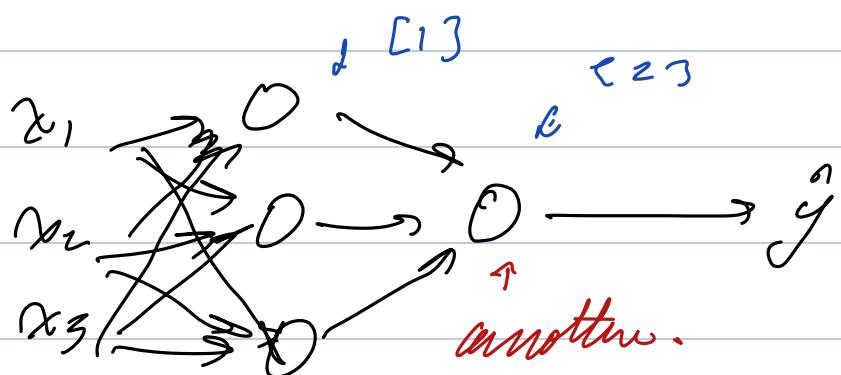
What is a Neural Network?

What is a Neural Network?



Features (Inputs)  $x$ , Parameters:  $w$  and  $\sigma$

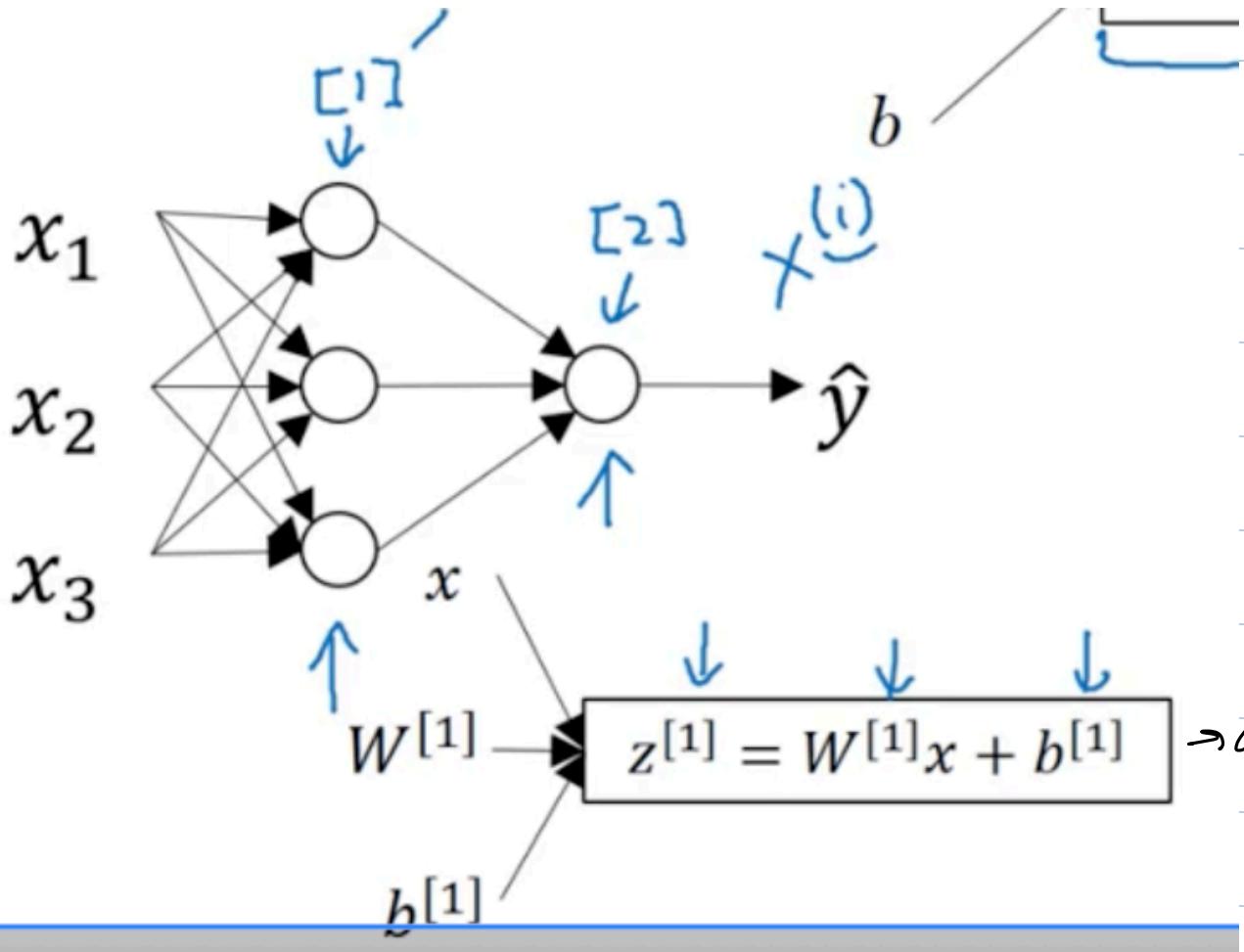
Neural Network:



" $z$ " like  
calculation

and " $a$ " like calculation

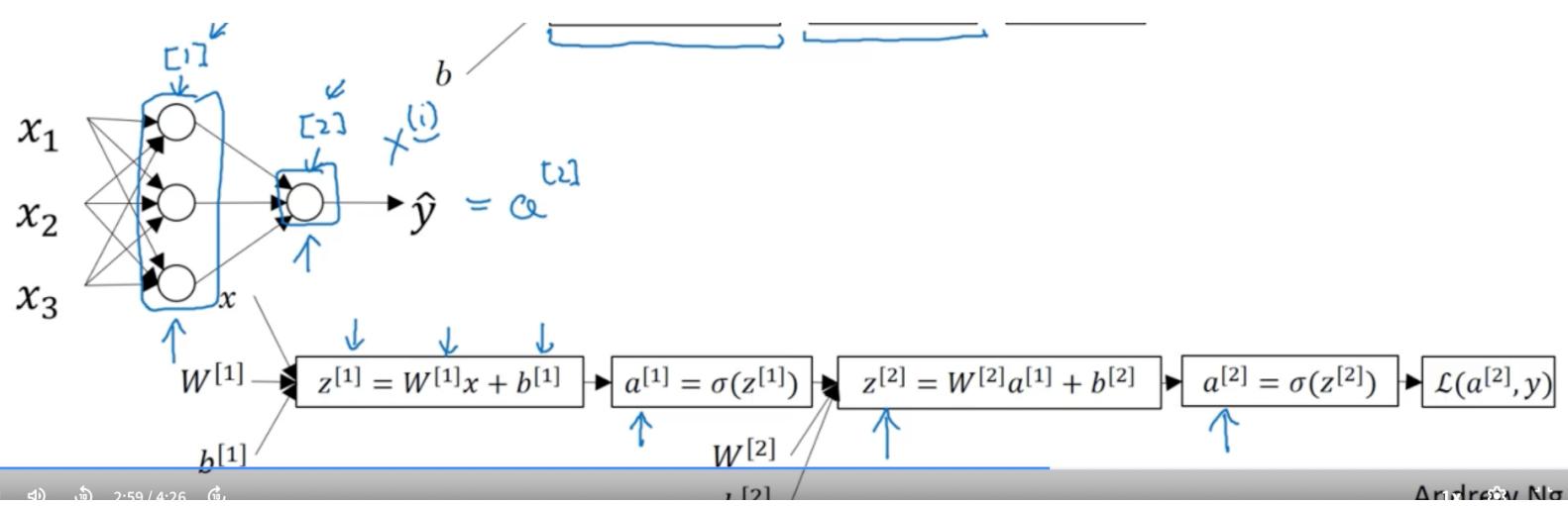
$$w^{[1]} \times z^{[1]} = w^{[1]} x + b^{[1]}$$



$[1] \rightarrow$  features associated to layer 1  
 $[2] \rightarrow$  layer 2

(i)  $\rightarrow$  individual training example.  $x^{(i)}$

$W^{[1]}$   $\rightarrow$  layer 1.



In logistic regression  $\rightarrow$  2 calls  $\Rightarrow$  a cost  
in neural networks it's a lot of them negative calculations.

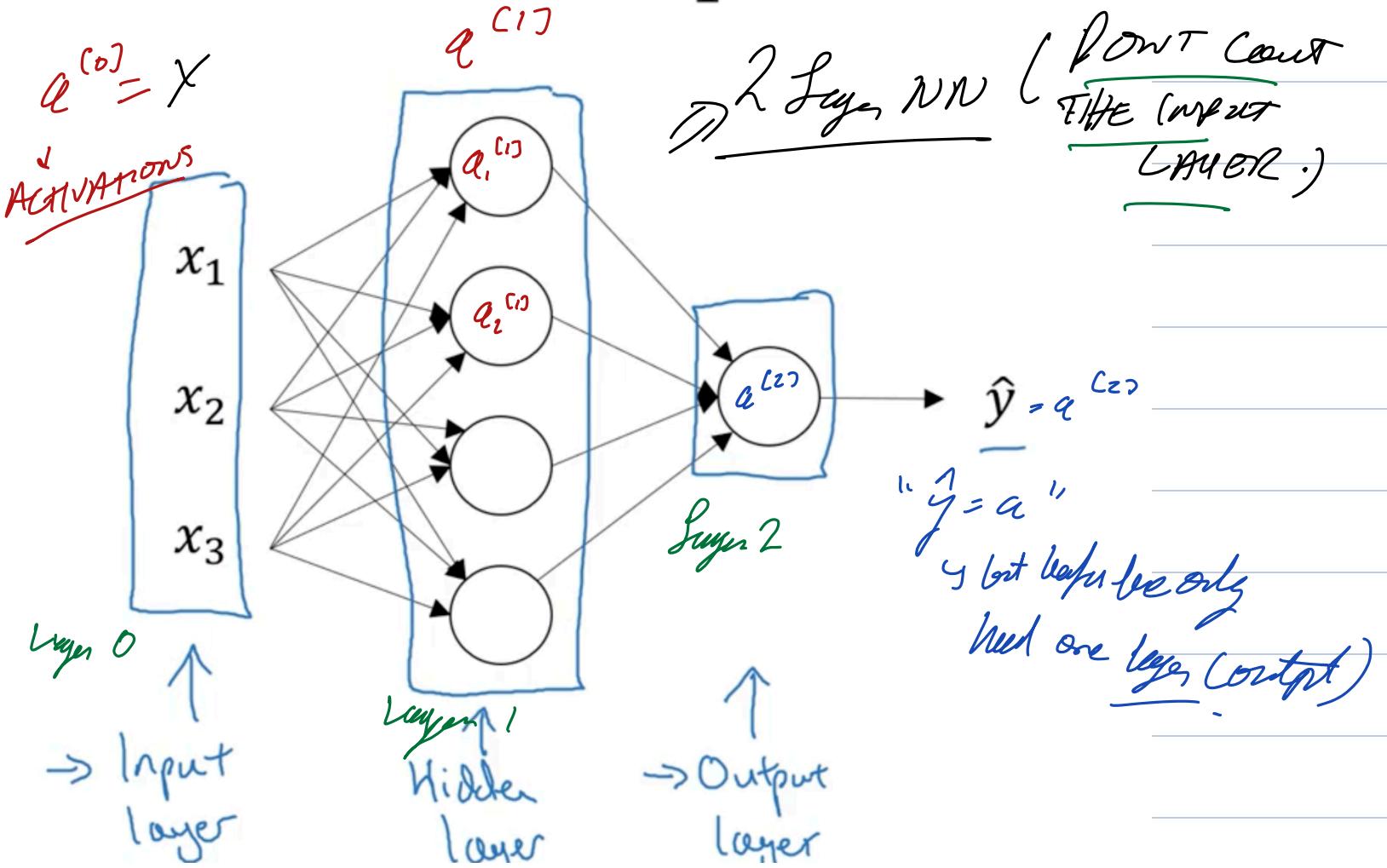
All lots of backwards calculations here as well for the gradient

$$d\omega^{(2)} \quad da^{(2)}$$

Think of logistic regression model and repeating it over all  
over layers

Neural Network representations and what they mean.

Similar Hidden Unit

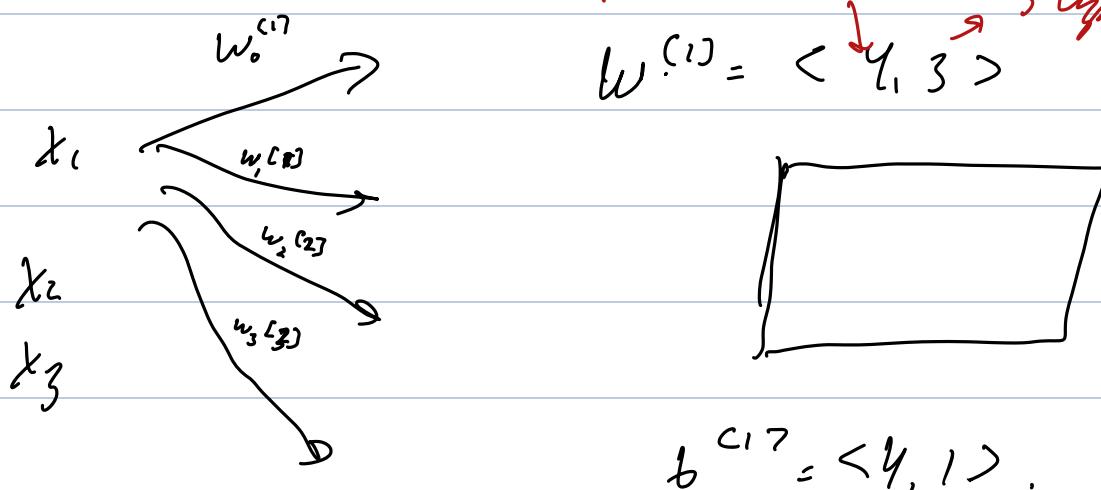


Input layer has weights and target output  $\hat{y}$ .  
 So hidden layers, they we don't know what they should  
 be in the training set.

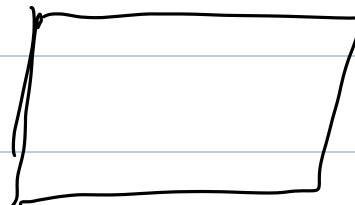
INPUT LAYER passes on  $x$  to hidden layers on the activation

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_n^{(1)} \end{bmatrix}$$

Hidden layers will have  $w^{(l)}, b^{(l)}$  associated with each layer.



$y$  nodes /  $y$  hidden units  
 $w^{(1)} = \langle 4, 3 \rangle$  → 3 input features

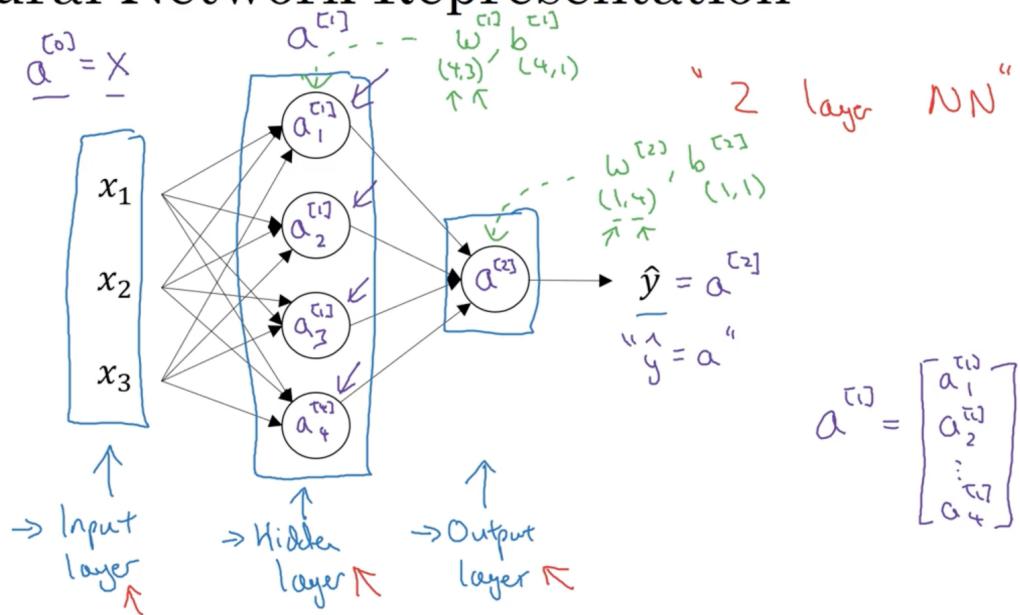


annnotated to the specific layer.

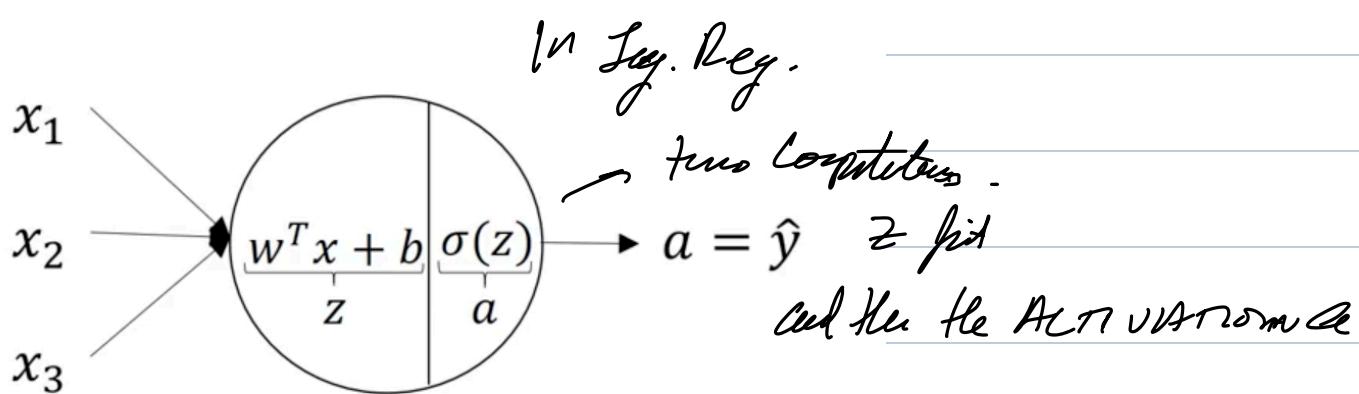
$$w^{(2)}, b^{(2)} \Rightarrow \langle 1, 4 \rangle \quad b^{(2)} \Rightarrow \langle 1, 1 \rangle$$

output      hidden layer by  $y$   
by 1              layers

## Neural Network Representation



## Computing A Neural Network's Outputs

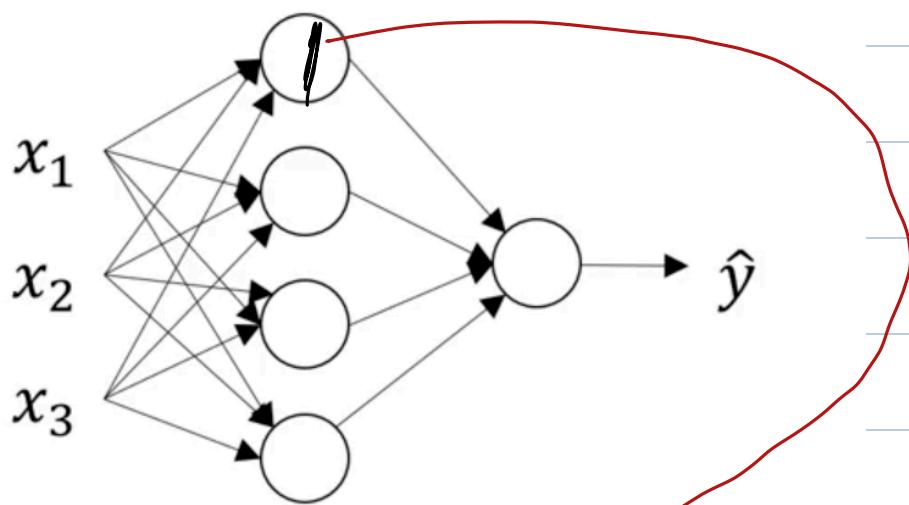


$$z = w^T x + b$$

$$a = \sigma(z)$$

/

A neural network just does this  $\times$  bunch of trials.



Zoom into one node of hidden layer

fit to leg region, does 2 steps of computation

$\downarrow$  first layer

$$z_i^{(1)} = w_i^{(1)T} x + b_i^{(1)} \quad a_i^{(1)} = \sigma(z_i^{(1)})$$

$\downarrow$

from  
node

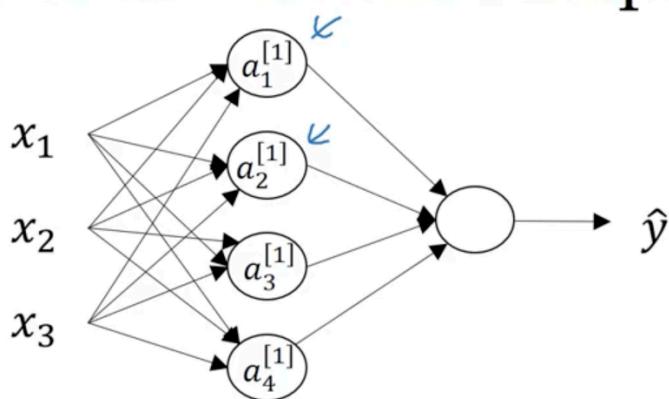
$a_i^{(1)} \rightarrow$  layer  
 $\downarrow$  node in layer.

2nd node in hidden layers.

$$z_2^{(1)} = w_2^{(2)T} x + b_2^{(1)}$$

$$a_2^{(1)} = \sigma(z_2^{(1)})$$

## Neural Network Representation



$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

$$z_2^{(1)} = w_2^{(1)T} x + b_2^{(1)}, \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

$$z_3^{(1)} = w_3^{(1)T} x + b_3^{(1)}, \quad a_3^{(1)} = \sigma(z_3^{(1)})$$

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}, \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

vectorize the equations.

$w^{(1)}$

$4 \times 3$ .

$w_{vec} \in \mathbb{R}^{3 \times 4}$

now  $\in \mathbb{R}^{4 \times 3}$

$b^{(1)}$

$$y = \underbrace{\begin{bmatrix} -w_1^{(1)\top} \\ -w_2^{(1)\top} \\ -w_3^{(1)\top} \\ -w_4^{(1)\top} \end{bmatrix}}_3 + \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}}_{3 \times 1} = \begin{bmatrix} w_1^{(1)\top} x + b \\ w_2^{(1)\top} x + b \\ w_3^{(1)\top} x + b \\ w_4^{(1)\top} x + b \end{bmatrix}$$

$\downarrow$  for each hidden node

$\langle y, 3 \rangle$  each of those nodes

compute by previous  
computations differently.

$$z^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$$

apply  $\sigma(z)$

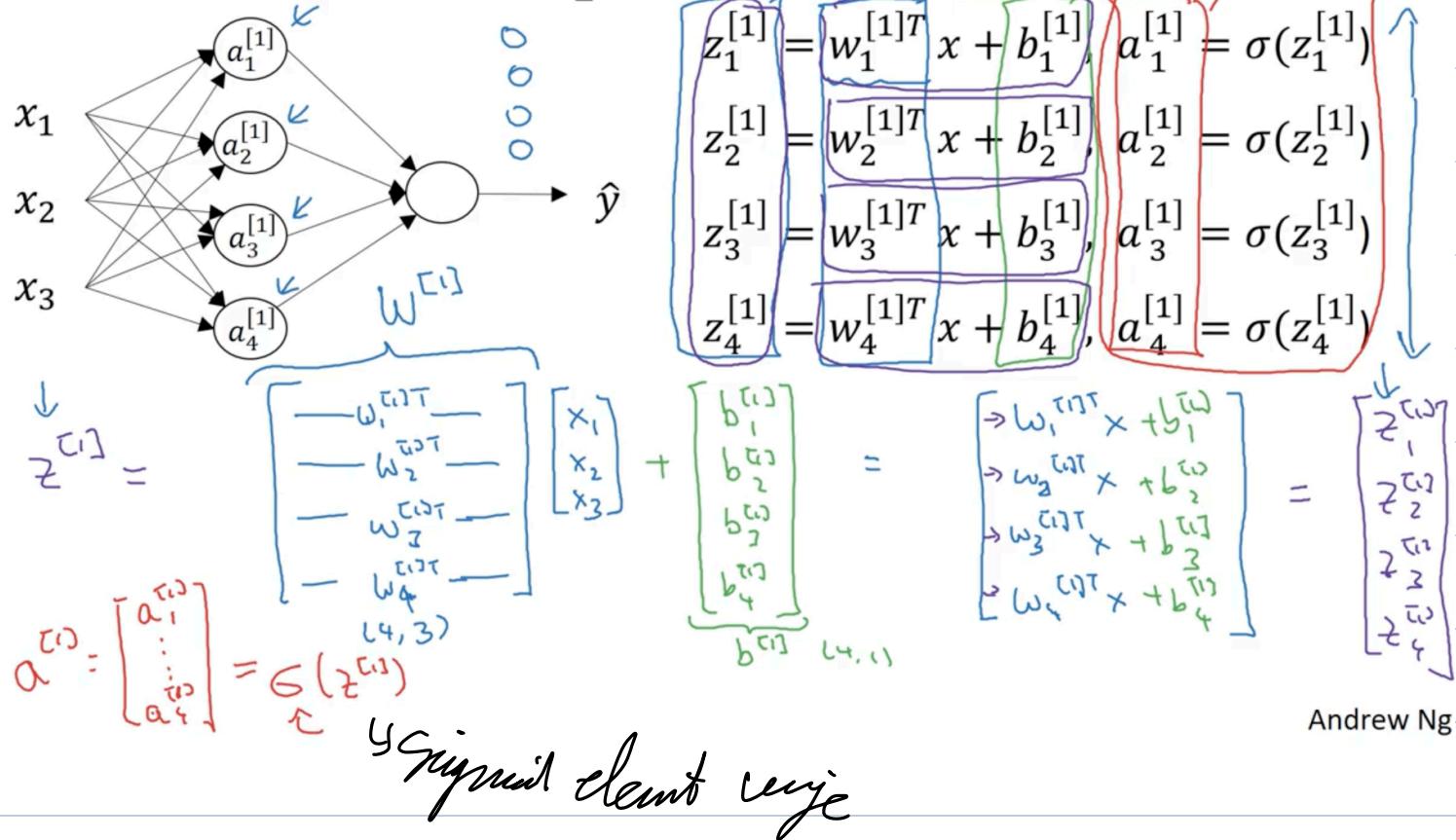
to get ACTIVATIONS

$\downarrow$  nodes scale vertically to get  $z^{(1)}$

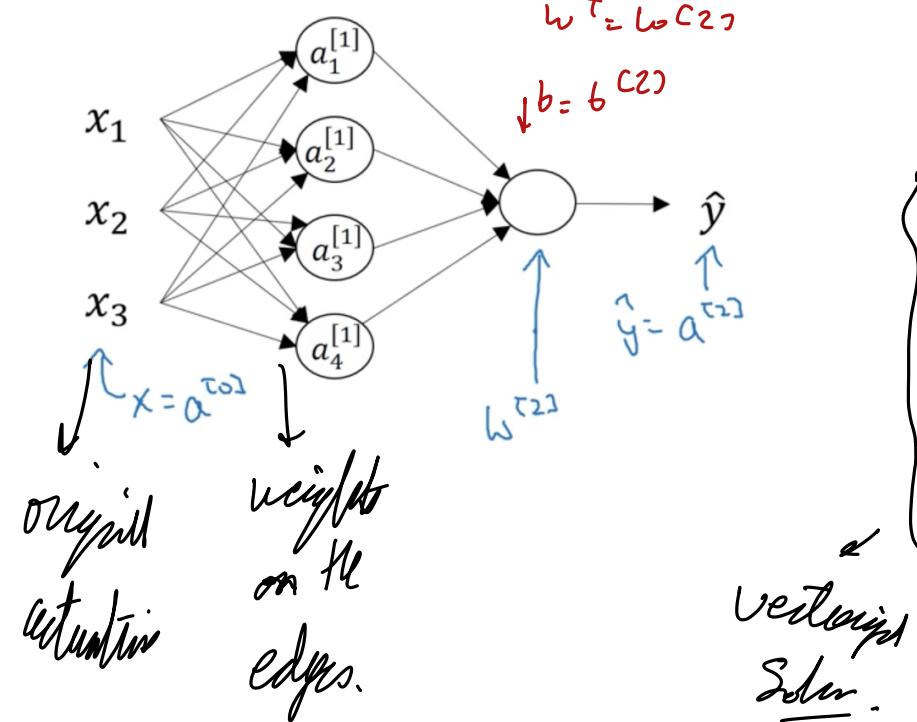
$\downarrow$

$\downarrow$

# Neural Network Representation



# Neural Network Representation learning



Given input  $x$ :

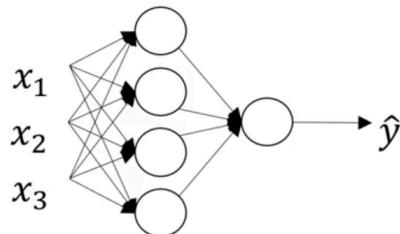
$$\begin{aligned} &\rightarrow z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} \\ &\rightarrow a^{[1]} = \sigma(z^{[1]}) \\ &\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ &\rightarrow a^{[2]} = \sigma(z^{[2]}) \end{aligned}$$

Andrew Ng

Really just a bunch of logistic regressions with different activations and weights.

## Vectorize Across Multiple Examples

Vectorizing across multiple examples



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

for  $i=1 \text{ to } n$

$$z^{[1](i)} = w^{(1)}x^{(i)} + b^{(1)}$$

$$x \xrightarrow{\quad} a^{(2)} = \hat{y}$$

input  
feature  
vector  $x$

$m$  examples.

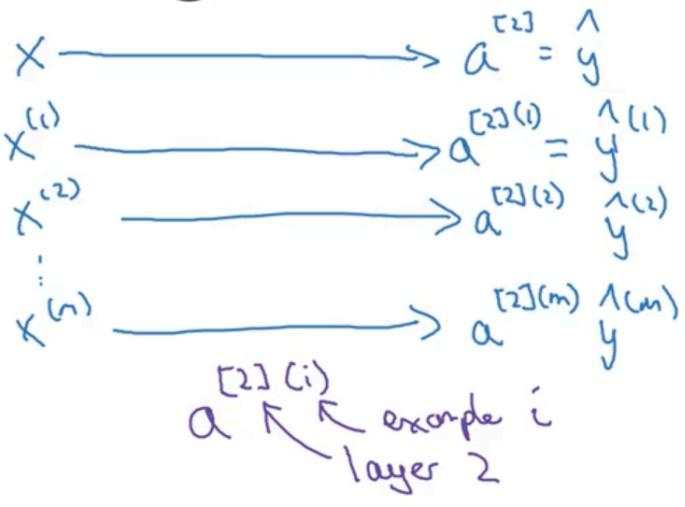
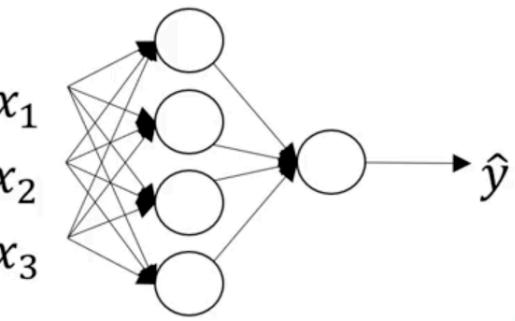
$$x^{(1)} \rightarrow a^{(2)(1)} \hat{y}^{(1)}$$

$$x^{(2)} \rightarrow a^{(2)(2)} \hat{y}^{(2)}$$

$$\vdots \qquad \qquad \qquad a^{(2)(m)} \hat{y}^{(m)}$$

$$x^{(m)} \rightarrow a^{(2)(m)} \hat{y}^{(m)}$$

$a^{(2)(i)}$  example  $i$   
layer 2



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right\}$$

for  $i = 1$  to  $m$ ,

$$z^{[1](i)} = w^{[1]} x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

↓  
repeat to

Andrey

Vetorize all this

for  $i = 1$  to  $m$ :

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

full

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

$$(n_{\infty}, m) \rightarrow X = A^{(\infty)}$$

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

$$Z^{[1]} = \begin{bmatrix} | & | \\ z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](M)} \\ | & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & & | \\ a^{[1](1)} & & a^{[1](2)} & \dots & a^{[1](M)} \\ | & & | \end{bmatrix}$$

*Neurons*  $\downarrow$  *vertically across different nodes of activation / neurons*

*horizontally across different examples*

← → TRAINING EXAMPLES  
[ HIDDEN LAYER UNITS  
OR NEURONS

# JUSTIFICATION OF THE VECTORIZED IMPLEMENTATION

$$z^{(1)(1)} = w^{(1)} x^{(1)} + b^{(1)} \quad \xrightarrow{\text{to simplify}}$$

$$z^{(1)(2)} = w^{(1)} x^{(2)} + b^{(1)} \quad \xrightarrow{\text{to}} \dots$$

$$z^{(1)(3)} = w^{(1)} x^{(3)} + b^{(1)} \quad \dots$$

$$w^{(1)} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

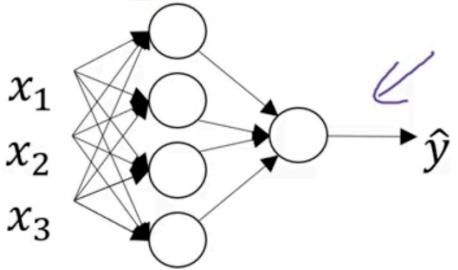
$$w^{(1)} x^{(2)} = \begin{bmatrix} : \\ : \end{bmatrix}$$

$$w^{(1)} x^{(3)} = \begin{bmatrix} : \\ : \end{bmatrix}$$

$$w^{(1)} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ | & | & | \end{bmatrix} = z^{(1)}$$

$\lambda + b^{(1)} k + b^{(1)} c + b^{(1)}$

# Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$\underline{A^{[1]}} = \begin{bmatrix} | & | & | \\ a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

for  $i = 1$  to  $m$

$$\rightarrow z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$\rightarrow a^{[1](i)} = \sigma(z^{[1](i)})$$

$$\rightarrow z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$\rightarrow a^{[2](i)} = \sigma(z^{[2](i)})$$

$\downarrow \text{if } \sigma$

$$\left. \begin{array}{l} Z^{[1]} = W^{[1]}X + b^{[1]} \\ A^{[1]} = \sigma(Z^{[1]}) \\ Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} = \sigma(Z^{[2]}) \end{array} \right\} \text{Vectorized}$$

$$Z^{[1]} = W^{[1]}A^{[0]} + b^{[1]}$$

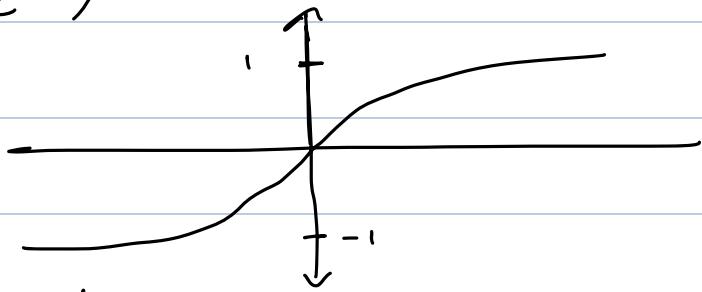
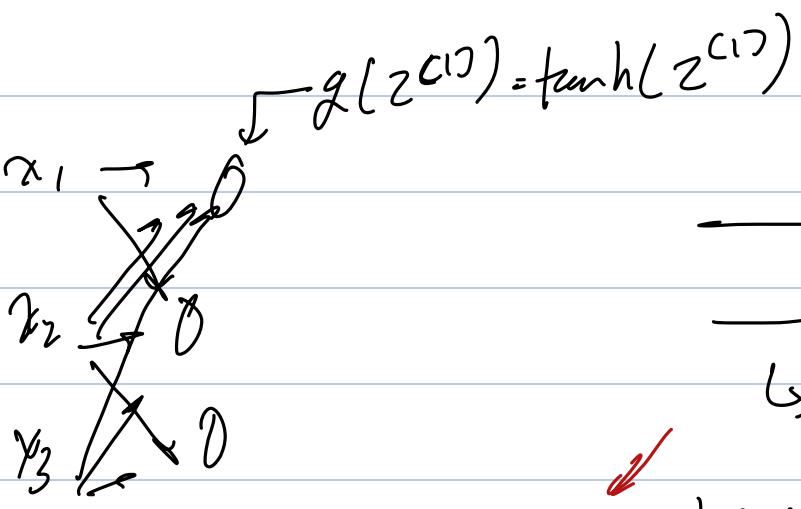
$\rightarrow$  so ~~the~~ diff layers are  
essentially doing the  
same thing to do.

ACTIVATION FUNCTIONS  $\rightarrow$  for hidden and output  
layers.

Better than sigmoid

$$A^{[1]} = \sigma(Z^{[0]})$$





$$a = \tanh(z)$$

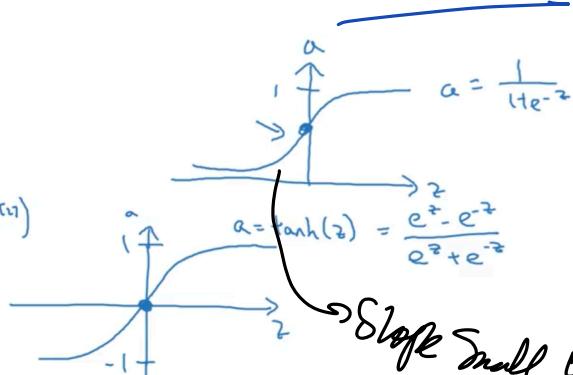
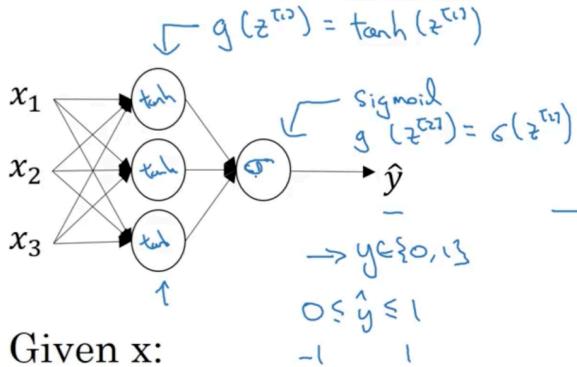
then

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

$\partial z$  which  
is what we do to stabilize to has flat  
kind of effect.

Sigmoid could be better for binary classification  
where  $y \in \{0, 1\}$ . For the output  
 $0 \leq \hat{y} \leq 1$

### Activation functions



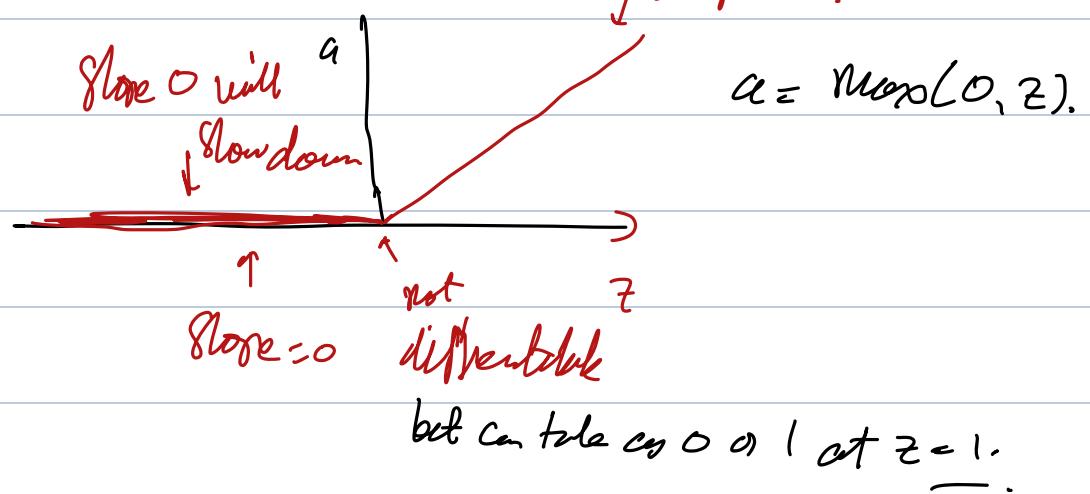
Slope small when  $z$  is large  
or small, so

slows down

gradient descent

Activation function can do different at different layers.

Rectified



ReLU For

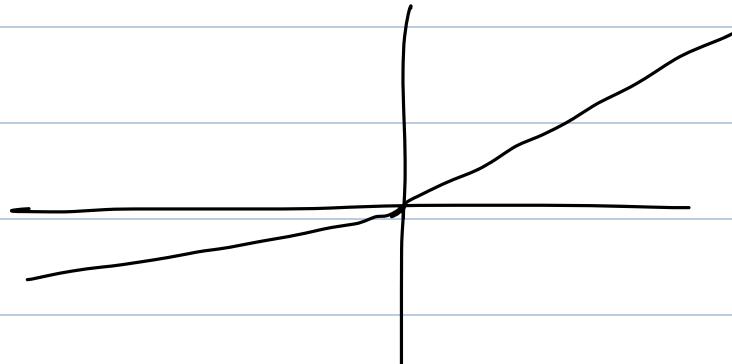
Rectified  
Linear Unit

if Binary classifier  $\rightarrow$  Sigmoid,

Else very ReLU is a rule closer to hidden layer.

Leaky ReLU

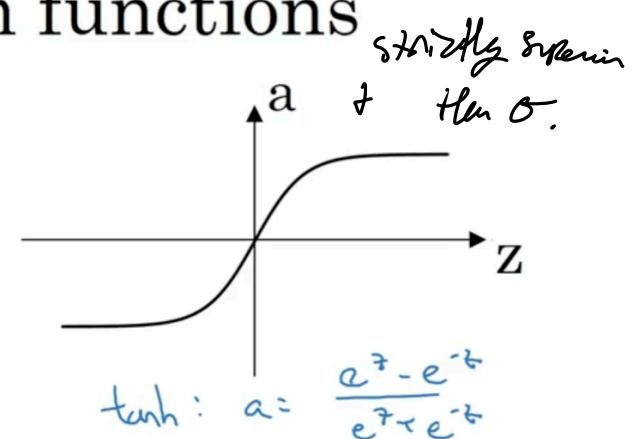
also works  
better



## Pros and cons of activation functions

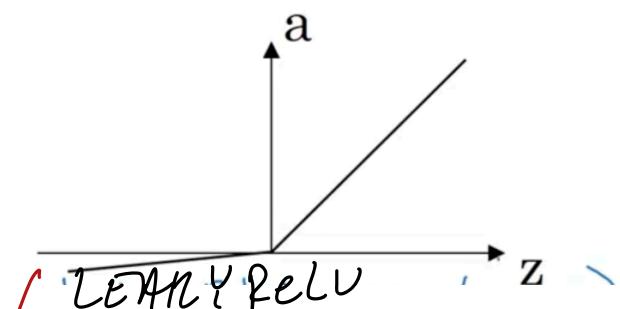
$\rightarrow$  Never use  
EXCEPT for binary  
logits to do binary  
classification.

$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



DEFAULT.  
. Not smooth  
pls

$$\text{ReLU } a = \max(0, z)$$



$$\alpha = \text{max}(0.01z, z)$$

could by.

$$0.01$$

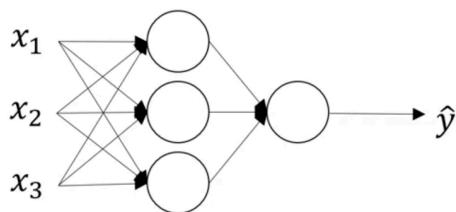
Role by this  
is also a  
parameter

## WHY DOES NEURAL NETWORK NEED ACTIVATION FN

Why not just use the ReLU

WHY DOES IT NEED A NON-LINEAR TRANSFER FN

Activation function



Given  $x$ :

$$\rightarrow z^{[1]} = W^{[1]}x + b^{[1]}$$

$$\rightarrow a^{[1]} = g^{[1]}(z^{[1]}) \stackrel{C^1}{\geq}$$

$$\rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$\rightarrow a^{[2]} = g^{[2]}(z^{[2]})$$

$g(z) = z$   
"linear activation  
function"  
identity?

why not do this?

$$\text{if } a^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$\text{will then } a^{[2]} = w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]}$$

$$a^{[2]} = \underbrace{w^{[2]}w^{[1]}}_{w'} x + \underbrace{w^{[2]}b^{[1]} + b^{[2]}}_{b'}$$

$$= \mathbf{w}' \mathbf{x} + b'$$

So  $\therefore$  no matter how many hidden layers, it's essentially one big linear function.

Same computation of 2 lin fun = 1 lin fun.

So you need non-lin fun.

Where would you use lin fun?

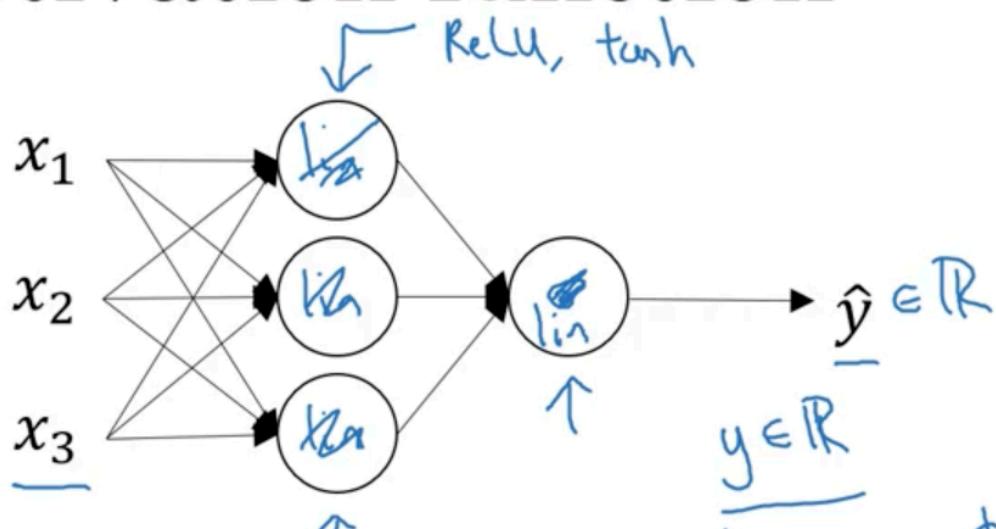
↳ On Regression Problem, where  $y$  is a real #

$y \in \mathbb{R} \rightarrow$  house prices  $0 \rightarrow \$1,000,000,$

$\therefore$  it might be possible to have lin act the output

layer

## Activation function



So linear function activation function is not that good.

# DERIVATIVES OF ACTIVATION FUNCS

for Gradient Descent

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

$$g(z) = a . \quad g'(z) = a(1-a)$$

Tanh       $g(z) = \tanh(z)$        $g'(z) = 1 - (\tanh(z))^2$   
 $= \frac{e^z - e^{-z}}{e^z + e^{-z}}$        $g'(z) = 1 - (g(z))^2 .$   
 $g'(z) = 1 - a^2$   
 $g(z) = a .$

ReLU

$$g(z) = \max(0, z) \quad g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Note - if  $z = 0$ ,

for programming ↑ set  $\epsilon$  to 0 or 1

Leaky ReLU

$$g(z) = \max(0.01z, z) \quad g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

$\rightarrow$  for  $z \rightarrow \infty$ , same as  
 ReLU,  
 ReLU

## Gradient Descent for Neural Network with one Hidden Layer

Parameters:  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$ .  $n^x = n^{(0)}$ ,  $n^{(1)}$ ,  $n^{(2)} = 1$   
 $\langle n^{(1)}, n^{(0)} \rangle$   $\langle n^{(1)}, 1 \rangle$   $\langle n^{(2)}, n^{(1)} \rangle$   $\langle n^{(2)}, 1 \rangle$   $\begin{matrix} \text{input features} \\ \text{hidden units} \\ \text{output units} \end{matrix}$

Cost Function for Binary Classification

$$J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) \\ = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i, y_i)$$

Gradient Descent:  $\rightarrow$  Initialize parameters Randomly and not to 0

Repeat {

Compute  $(\hat{y}^{(i)}, i=1 \dots m)$

$$d_w^{(1)} = \frac{dJ}{dw^{(1)}}, d_b^{(1)} = \frac{dJ}{db^{(1)}}, \dots$$

$$w^{(1)} := w^{(1)} - \alpha d_w^{(1)}$$

## FORMULAS FOR COMPUTING PERLATES :

Forward Prop.

$$Z^{(1)} = W^{(1)} X + b^{(1)}$$

$$A^{(1)} = g^{(1)}(Z^{(1)})$$

$$Z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(Z^{(2)})$$

Backprop

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$dZ^{(2)} = A^{(2)} - Y$$

$$dW^{(2)} = \frac{1}{m} dZ^{(1)} A^{(1)T}$$

$$db^{(2)} = \frac{1}{m} \text{Ap. Sum}(dZ^{(2)}, \text{axis} = 1, \text{keepdim})$$

Sum across

axis = 1

(rows sum)

(reduce)

\*Steps

(n,)

$$dZ^{(1)} = W^{(2)T} dZ^{(2)} \quad \begin{matrix} \leftarrow \\ (n^{(1)}, m) \end{matrix} \quad \begin{matrix} \rightarrow \\ \text{element wise} \end{matrix} \quad \begin{matrix} \leftarrow \\ g^{(1)'}(Z^{(1)}) \end{matrix} \quad \begin{matrix} \rightarrow \\ (n^{(1)}, m) \end{matrix}$$

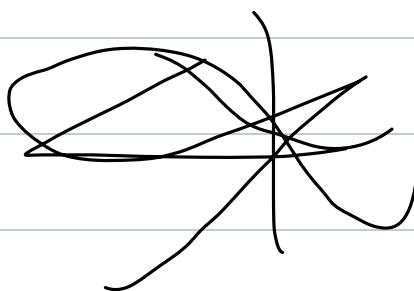
$$dW^{(1)} = \frac{1}{m} dZ^{(1)} X^T$$

$$db^{(1)} = \frac{1}{m} \text{Ap. Sum}(dZ^{(1)}, \text{axis} = 1, \text{keepdim} = \text{true})$$

$$\hookrightarrow (n^{(1)}, 1)$$

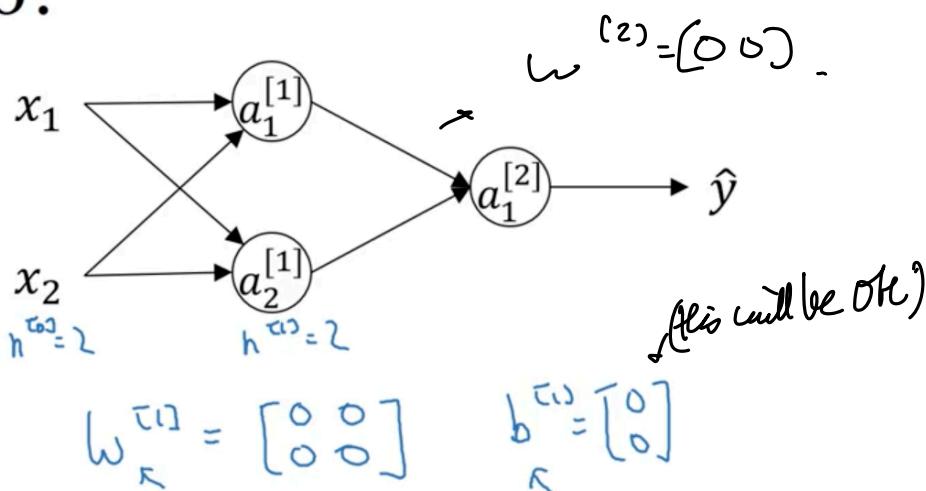
CATCH UP WITH THESE VIDS TO SEE

REDUCTION



## Random Initialization

What happens if you initialize weights to zero?



Therefore,

$$q_1^{(1)} = q_2^{(1)}. \text{ and then } \partial z_1^{(1)} = \partial z_2^{(1)}.$$

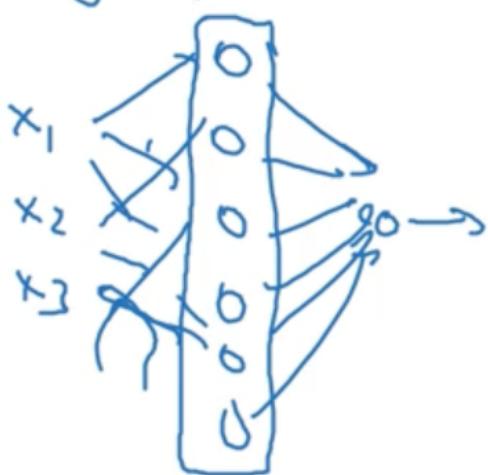
So these two hidden units are symmetric -

$$\partial w = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \therefore w^{(1)} = w^{(1)} - \alpha \partial w.$$

THE 2 HIDDEN UNITS ARE LEARNING

the same thing

Symmetric



→ So they're doing the same thing.

Use Gaussian random.

$$w^{(1)} = \text{np.random.}\tilde{\text{randn}}(2, 2) + 0.01$$

$$b^{(1)} = \text{np.zeros}((2, 1))$$

to a small random value

I don't have problem here.

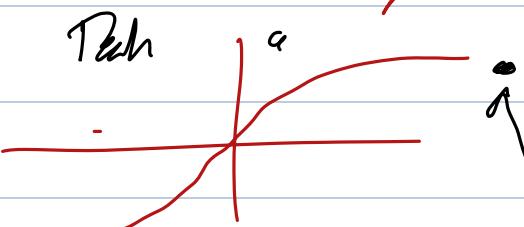
$$w^{(2)} = \text{random}$$

$$b^{(2)} = 0$$

Can be  
better, for deep  
neural network

better

PREFER to start with  
small weights



$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

large weights  $\Rightarrow$  large  $a^{(1)}$

So gradient descent

(or learning) will be  
very slow.