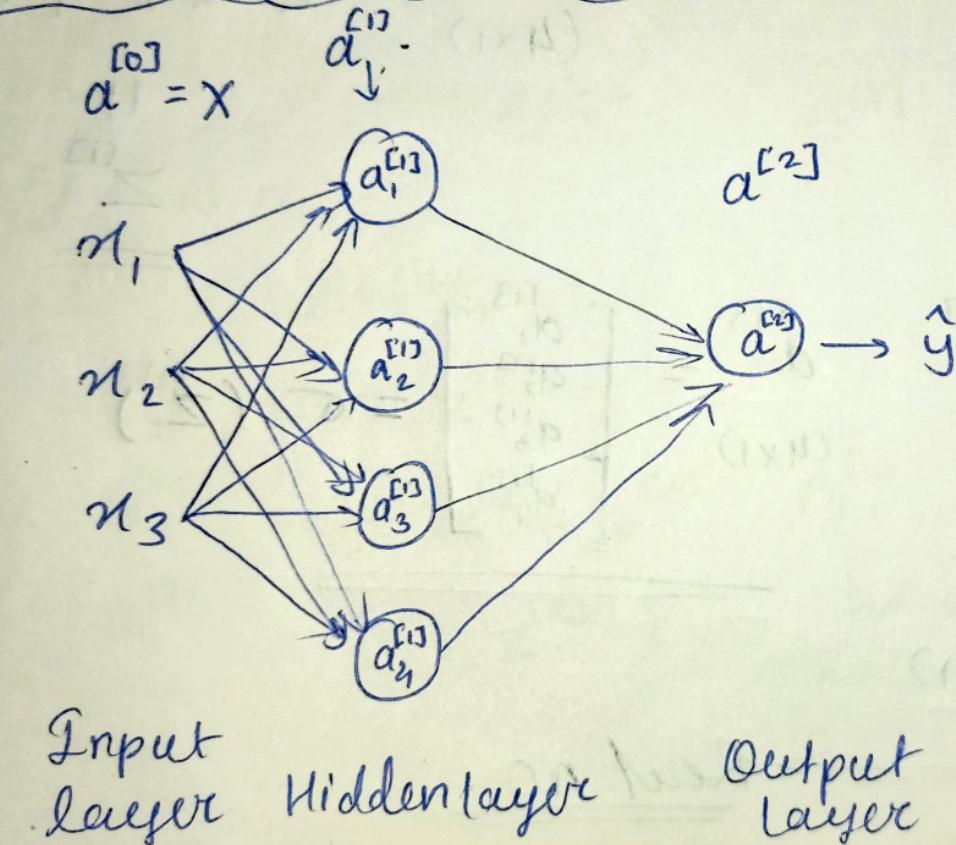


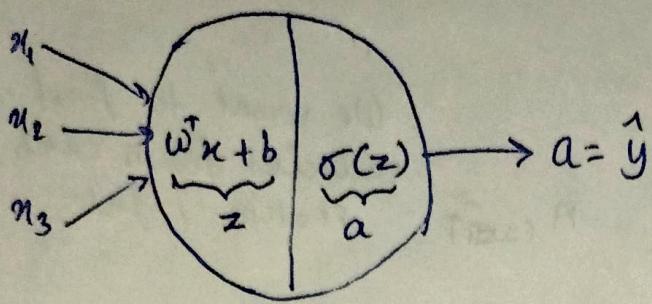
Week 3

* Neural Network Representation.



This is also called
a 2 layer NN (we don't
count input layer)

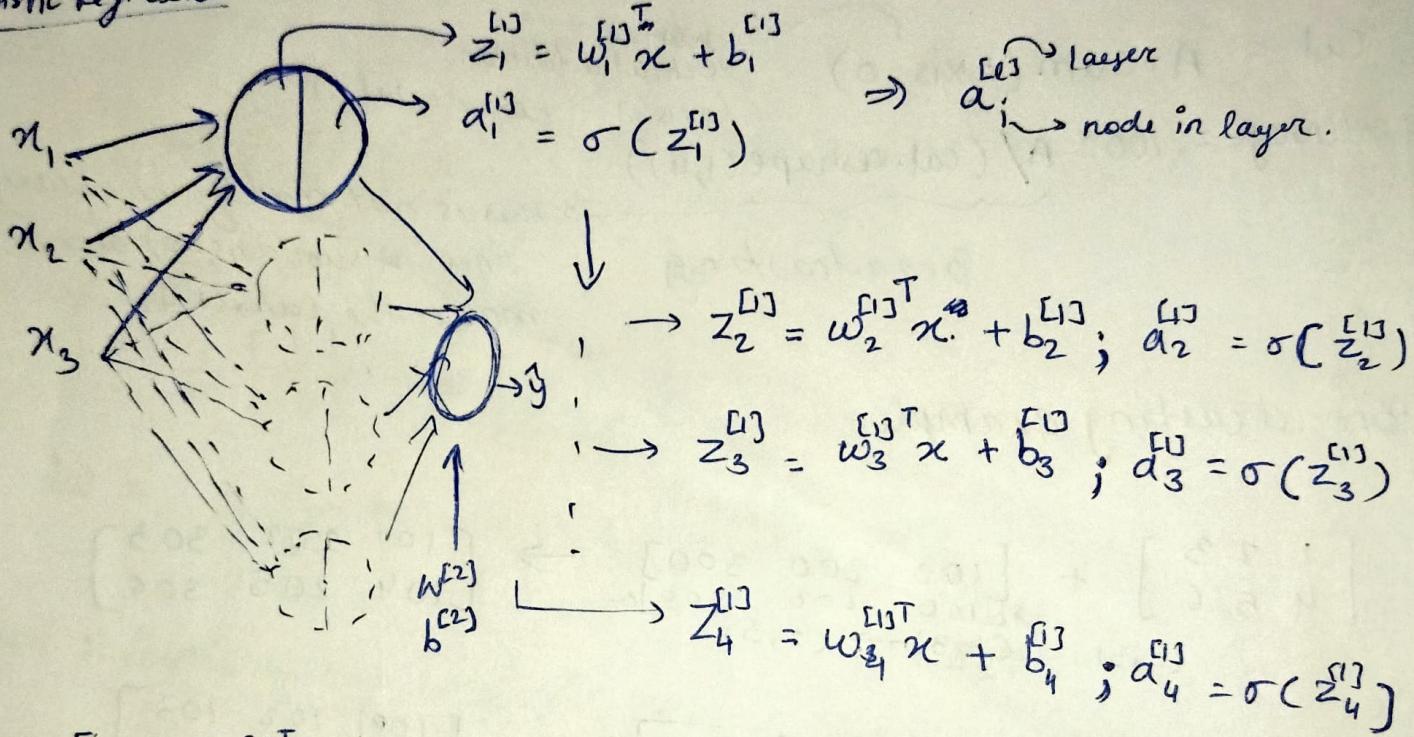
$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \\ a_5^{[1]} \end{bmatrix}$$



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

$$a = \sigma(z).$$

logistic Regression



$$\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{(3 \times 1)} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}_{(4 \times 1)} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}_{(4 \times 1)}$$

because there are 3 inputs.

Output layer

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

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

$$\begin{bmatrix} w^{[2]T} \\ b^{[2]} \end{bmatrix} \rightarrow (1 \times 4)$$

$$\rightarrow (1 \times 1) \quad z^{[2]} = (1 \times 1)$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}_{(4 \times 1)} = \sigma(z)$$

$\Sigma^{[1]}$

∴ In summary :- Given Input X :

$$\left. \begin{array}{l} \rightarrow z^{[1]} = w^{[1]} a^{[0]} + b^{[1]} \\ (4 \times 1) \quad (4 \times 3) (3 \times 1) \end{array} \right\}$$

$$\rightarrow a^{[1]} = \sigma(z^{[1]}) \quad (4 \times 1)$$

$$\rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \quad (1 \times 1) \quad (1 \times 4) \quad (4 \times 1) \quad (1 \times 1)$$

$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \quad \xrightarrow{\text{to}} \hat{y} \quad (1 \times 1)$$

This computes the result for 1 training example.
 $\therefore \hat{y}$ here is for the 1st training example.

few lines
of code

$$X \longrightarrow a^{[2]} = \hat{y}$$

$$X^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$

$$X^{(2)} \longrightarrow a^{2} = \hat{y}^{(2)}$$

$$\vdots$$

$$X^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$$

} m training examples

example no.

no. of features
= no. of nodes
in input layer

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

no. of examples

We want to vectorize :-

for $i = 1$ to m :

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

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

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

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

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

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

no. of examples

no. of nodes
in the layer

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

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

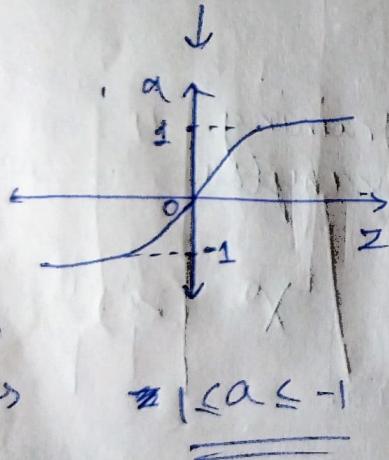
$$Z^{[2]} = W^{[2]}(A^{[1]}) + b^{[2]}$$

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

* Activation functions.

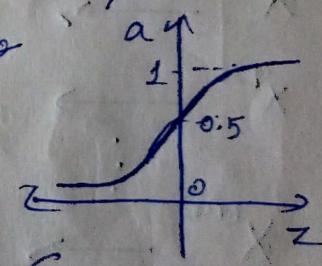
$a = \sigma(z)$
Sigmoid ✓ this is an activation function → this could be different

In fact, $\tanh(z)$ function works better than $\sigma(z)$.



$$\hookrightarrow -1 \leq a \leq 1$$

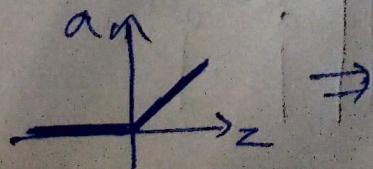
for 1st hidden layer or any hidden layer except maybe the last to output layer.



$$\hookrightarrow 0 \leq a \leq 1$$

→ Problem for both of them is that as the value of z gets really high or really low, the slope tends to zero and the gradient descent starts getting slower and slower.

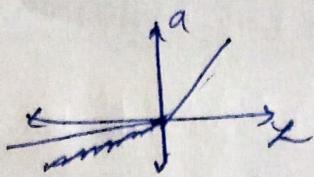
Hence we use ~~ReLU~~ ReLU function (Rectified Linear unit)



$$a = \max(0, z)$$

↪ used most commonly.

Leaky ReLU \rightarrow

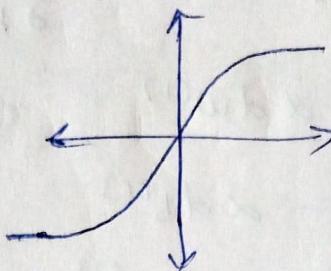


-ve z values do not lead to 'zero' a .

$$a = \max(0.01z, z)$$

* Slopes/derivatives.

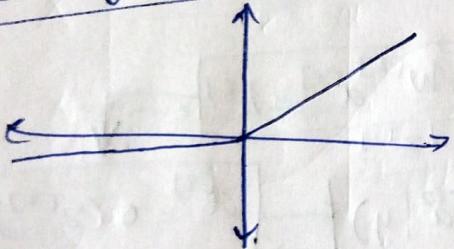
$$\begin{aligned}g(z) &= \tanh h(z) \\&= \frac{e^z - e^{-z}}{e^z + e^{-z}}\end{aligned}$$



$$g'(z) = 1 - (\tanh(z))^2$$

$$\hookrightarrow g'(z) = 1 - a^2$$

Leaky ReLU



$$g(z) = \max(0, z)$$

$$\begin{aligned}g'(z) &= \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \\ \text{undefined if } z = 0 \end{cases} \\&\quad z = 0.000\dots 0,\end{aligned}$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

* Gradient Descents for neural networks.

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$, $n_x = n^{[0]}$, $n^{[1]}$, $n^{[2]} = 1$
 $(n^{[1]}, n^{[0]}) (n^{[1]}, 1) (n^{[1]}, n^{[2]}) (n^{[2]}, 1)$
input features hidden units output units

Cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(g_i, y_i)$

Gradient Descent

Repeat {

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

$$d w^{[1]} = \frac{\partial J}{\partial w^{[1]}} , d b^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$$

$$w^{[1]} = w^{[1]} - \alpha d w^{[1]}$$

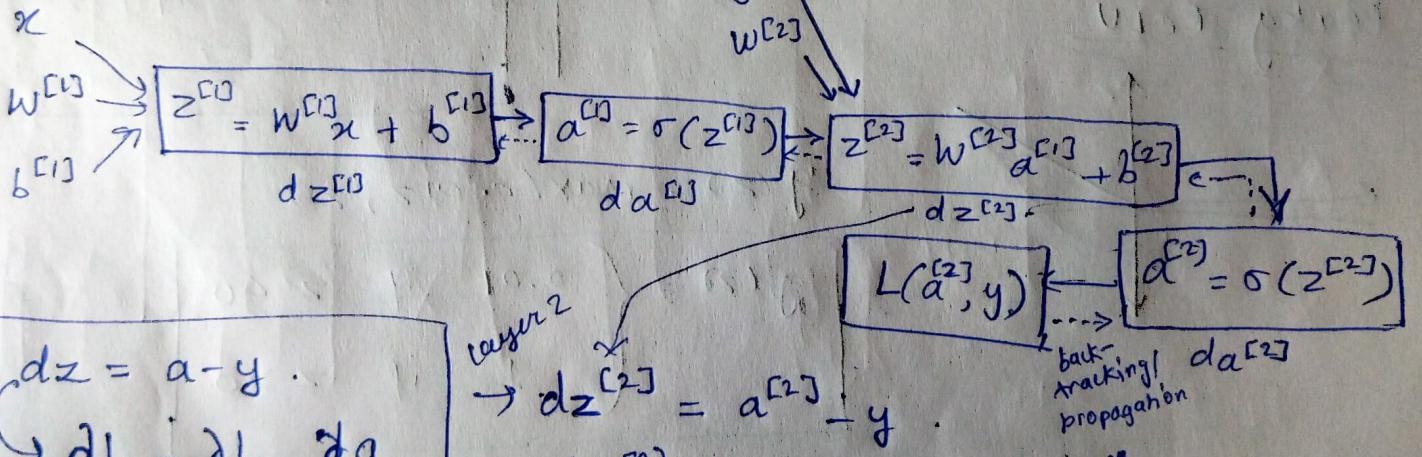
$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

$$w^{[2]} = w^{[2]} - \alpha d w^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$$

}

* Neural Network Gradients (Backpropagation derivation)



$$d z = a - y$$

$$\frac{d L}{d z} = \frac{d L}{d a} \cdot \frac{d a}{d z}$$

$$d z = d a \cdot \frac{d g(z)}{d z}$$

$$d z = d a \cdot g'(z)$$

Dimensions.

$$x_1 \quad 0$$

$$x_2 \quad 0$$

$$x_3 \quad 0$$

$$x = n^{[2]} \quad \text{(in binary classification)}$$

$$n^{[1]} = 1$$

$$\begin{aligned} \rightarrow d z^{[1]} &= d z^{[2]} \cdot a^{[1]T} \cdot g'(z^{[1]}) \quad (\text{similar to layer 2 values}) \\ \rightarrow d w^{[1]} &= d z^{[1]} \cdot x^T \\ \rightarrow d b^{[1]} &= d z^{[1]} \end{aligned}$$

$$\left. \begin{aligned} W^{[2]} &= (n^{[2]} \times n^{[1]}) \\ Z^{[2]} \cdot d z^{[2]} &= (n^{[2]} \times 1) \quad \{ \text{in this case } (1, 1) \} \\ Z^{[1]} \cdot d z^{[1]} &= (n^{[1]} \times 1) \end{aligned} \right\}$$

Summary of Gradient Descent:

$$dz^{[2]} = a^{[2]} - y$$

$$dw^{[2]} = dz^{[2]} \cdot a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dw^{[1]} = dz^{[1]} \cdot x^{[0]T}$$

$$db^{[1]} = dz^{[1]}$$

Vectorized implementation (for 'm' no. of examples).

$$dz^{[2]} = A^{[2]} - Y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$$

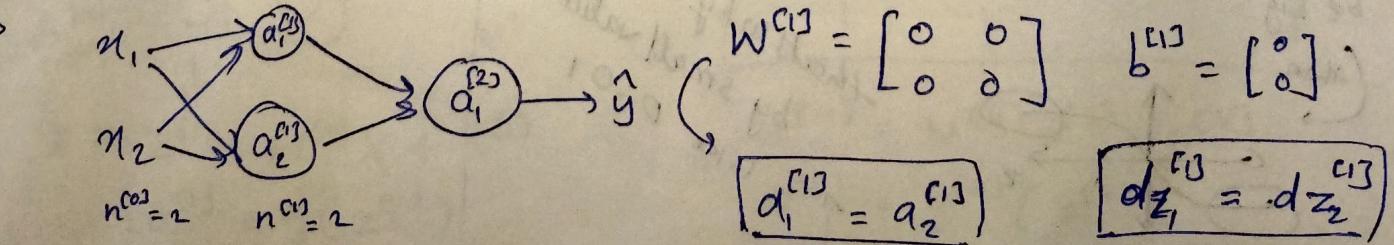
avoids something like $(m, 1)$ and instead writes $(m, 1)$

$$(n^{[1]}, m) \quad dz^{[1]} = \underbrace{w^{[2]T} dz^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)}$$

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$$

* Why We should randomly initialize parameters instead of initializing with zero?



both nodes will have same values ultimately because they have the same function

So, both the nodes are exactly symmetrical, they are computing exactly the same thing.

due to this :- $d\mathbf{w} = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \rightarrow$ all rows end up exactly the same.

after every

\rightarrow gradient descent : $\mathbf{W}^{[1]} = \mathbf{W}^{[1]} - \alpha d\mathbf{w}$

$$\mathbf{W}^{[1]} = \begin{bmatrix} -11 & -11 \\ -11 & -11 \end{bmatrix}$$

So, even after training it many times they will still compute the same thing.

all rows are exactly same.

This is why we do random initialization.

Hence,

$$\left\{ \begin{array}{l} \mathbf{W}^{[1]} = \text{np.random.randn}(2, 2) * 0.01 \\ \mathbf{b}^{[1]} = \text{np.zeros}(2, 1) \end{array} \right.$$

we prefer very small values why?

$$\left\{ \begin{array}{l} \mathbf{W}^{[2]} = \dots \\ \mathbf{b}^{[2]} = 0 \end{array} \right.$$

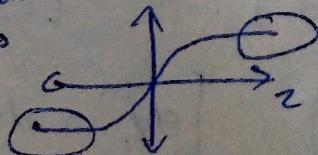
because as long as $\mathbf{W}^{[1]}$ is randomized, $d\mathbf{z}^{[1]}$ will be unique.

0.01 \rightarrow we prefer small values.

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

this will be big if this is big.

then



at big z values,

slope ≈ 0 , no/very slow gradient descent.

that's why we multiply by small values like 0.01