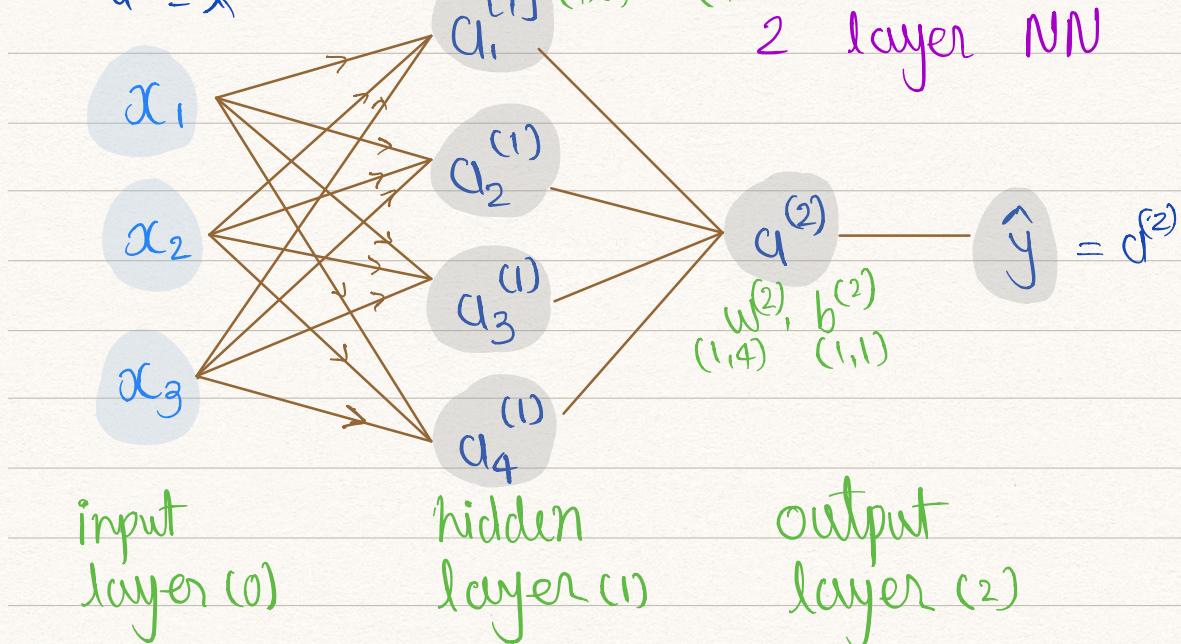


Neural Network Representation:

$$u^{(0)} = x$$

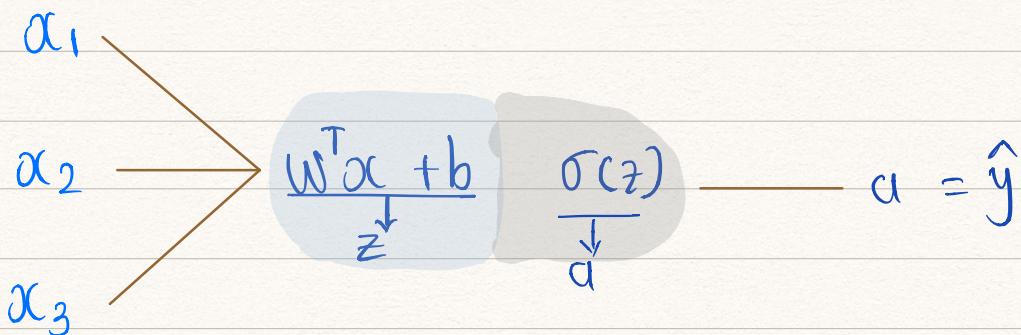


$$u^{(1)} = \begin{bmatrix} c_1^{(1)} \\ c_2^{(1)} \\ \vdots \\ c_4^{(1)} \end{bmatrix}$$

$\hat{y} = u^{(2)}$ because in logistic reg. we only get single o/p from N.N.

$w^{(1)} \rightarrow (4 \times 3) \rightarrow 4 \text{ output unit}$
 3 input unit

Computing NN's Output:



$a_i^{(l)}$ \rightarrow layer number
 a_i \rightarrow node in layer

$$\begin{aligned} z_1^{(1)} &= w_1^{(1)T} x + b_1^{(1)} \\ a_1^{(1)} &= \sigma(z_1^{(1)}) \end{aligned} \quad \left. \begin{array}{l} \text{layer 1, node 1} \\ \text{computation} \end{array} \right.$$

$$\begin{aligned} z_2^{(1)} &= w_2^{(1)T} x + b_2^{(1)} \\ a_2^{(1)} &= \sigma(z_2^{(1)}) \end{aligned} \quad \left. \begin{array}{l} \text{layer 1, node 2} \\ \text{computation} \end{array} \right.$$

Vectorize the whole process:

$$\left[\begin{array}{c} -w_1^{(1)T} \\ -w_2^{(1)T} \\ -w_3^{(1)T} \\ -w_4^{(1)T} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{array} \right] = \left[\begin{array}{c} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{array} \right] = \underline{z^{(1)}}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ \vdots \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$$

given input x :

$$z^{(1)}_{(4 \times 1)} = w^{(1)}_{(4 \times 3)} x_{(3 \times 1)} + b^{(1)}_{(4,1)}$$

$$a^{(1)}_{(4,1)} = \sigma(z^{(1)})_{(4,1)}$$

Algorithm:

for $i=1$ to m :

$$\begin{aligned} z^{(1)(i)} &= w^{(1)} \cdot x^{(i)} + b^{(1)} \\ a^{(1)(i)} &= \sigma(z^{(1)(i)}) \\ z^{(2)(i)} &= w^{(2)} a^{(1)(i)} + b^{(2)} \\ a^{(2)(i)} &= \sigma(z^{(2)(i)}) \end{aligned}$$

$\left. \begin{array}{l} \text{4 lines} \\ \text{of for} \\ \text{loop} \end{array} \right\}$

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{bmatrix}$$

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

training examples.

$$A^{(1)} = \underbrace{\left[\begin{array}{c|c|c|c} 1 & & & \\ \hline u^{(1)(1)} & u^{(1)(2)} & \cdots & u^{(1)(m)} \end{array} \right]}_{\text{hidden units}}$$

Explanation for vectorized implementation:

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

$$z^{(1)(2)} = w^{(1)} x^{(2)} + b^{(1)}$$

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

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

$$w^{(1)} \left[\begin{array}{c|c|c} 1 & x^{(1)} & x^{(2)} \\ \hline x^{(1)} & \vdots & \vdots \\ \hline \end{array} \right] = \left[\begin{array}{c|c|c} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{array} \right] = \left[\begin{array}{c|c|c} z^{(1)(1)} & z^{(1)(2)} & \dots \\ \hline +b^{(1)} & +b^{(1)} & \dots \\ \hline z^{(1)} & z^{(2)} & \dots \end{array} \right]$$

↓
Set of all training examples

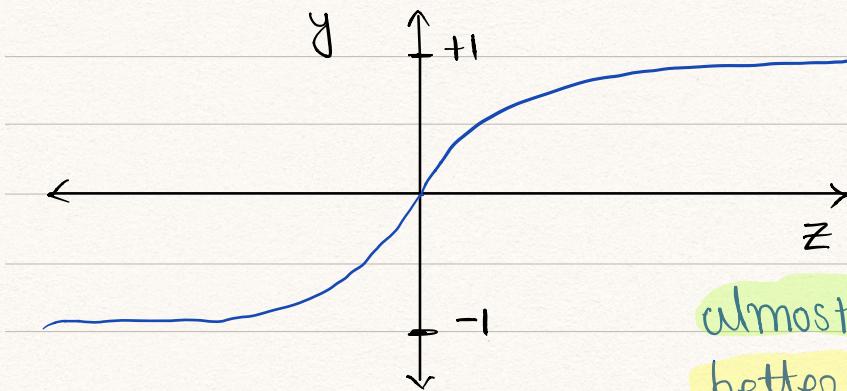
$$\left. \begin{array}{l} z^{(1)} = w^{(1)} X + b^{(1)} \\ A^{(1)} = \sigma(z^{(1)}) \\ z^{(2)} = w^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} = \sigma(z^{(2)}) \end{array} \right\} \begin{array}{l} \text{vectorization} \\ \text{of 4 lines} \\ \text{in for loop} \end{array}$$

Activation Functions:

Sigmoid Function:

$$\sigma = \frac{1}{1 + e^{-z}}$$

tanh Function :



$$y = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

almost always works
better than sigmoid

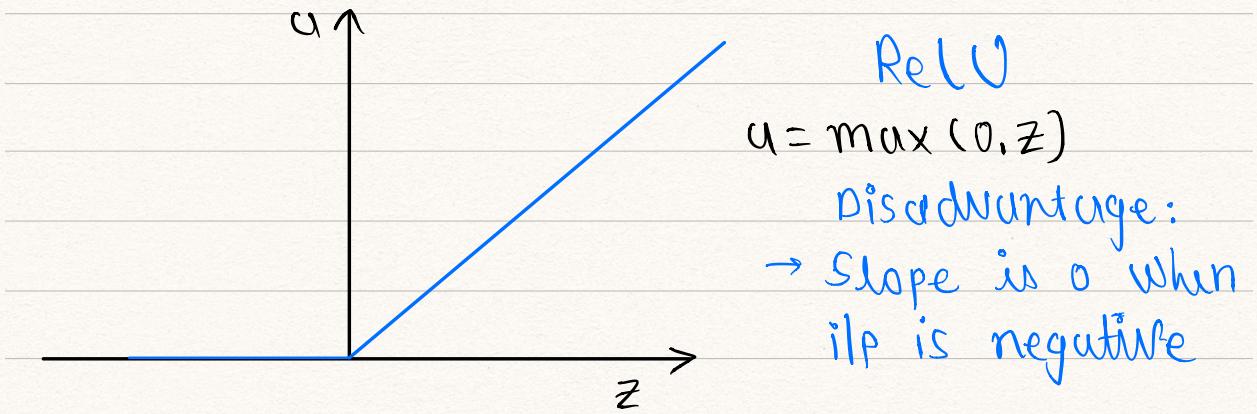
- have close to zero mean
- centering the data
- easier for further calculation

→ when you have binary classification output layer should be sigmoid.
keep hidden layers as tanh.

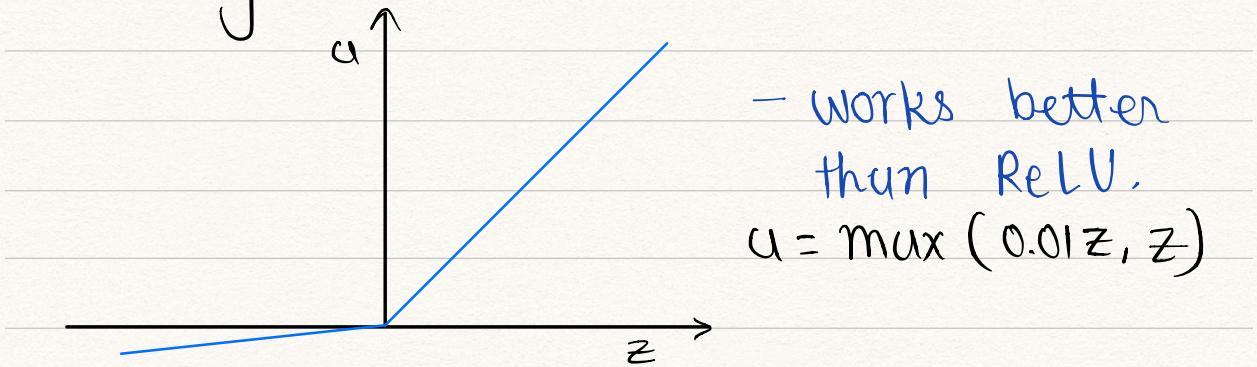
→ You can have different activation function for different layers.

→ When you have very large or very small i/p, slope of function is close to zero, so it can slow down gradient descent.

Rectified Linear Unit :



Leaky ReLU:



Why do we need non linear Activation Functions?

linear activation function:

$$g(z) = z$$

↳ identity function.

$$\begin{aligned}a^{(1)} &= z^{(1)} = w^{(1)}x + b^{(1)} \\a^{(2)} &= z^{(2)} = w^{(2)}a^{(1)} + b^{(2)} \\a^{(2)} &= w^{(2)}(w^{(1)}x + b^{(1)}) + b^{(2)} \\&= \frac{(w^{(2)}w^{(1)})x}{w^{(1)}} + \frac{(w^{(2)}b^{(1)} + b^{(2)})}{w^{(1)}}\end{aligned}$$
$$a^{(2)} = w'x + b'$$

No matter how many layer NN has
its just doing linear activation.
No benefit of doing this.

You can use linear activation on output layer, not recommended at hidden layer.

Derivatives of Activation Functions:

When we implement back prop. on NN
We need to compute slope or the
derivative of activation function.

Sigmoid Activation Function:

$$g(z) = \frac{1}{1+e^{-z}} = a$$

$$\begin{aligned} g'(z) &= \frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ &= g(z) \cdot (1 - g(z)) \\ &= a(1-a) \end{aligned}$$

tanh activation:

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

$$\begin{aligned} g'(z) &= \frac{d}{dz} (g(z)) = 1 - (\tanh(z))^2 \\ &= 1 - (g(z))^2 \\ &= \underline{\underline{1 - a^2}} \end{aligned}$$

ReLU & Leaky ReLU:

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

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

~~undefined " $z=0$ "~~

↳ due to discontinuity.

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

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

~~undefined $z=0$~~

Gradient Descent for NNs:

Parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$
 $(n^{(1)}, n^{(0)}) \quad (n^{(1)}, 1) \quad (n^{(2)}, n^{(1)}) \quad (n^{(2)}, 1)$
 $n_x = n^{(0)}, n^{(1)}, n^{(2)} = 1$

Cost Function: $J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)})$

$$= \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}, y)$$

$\hat{y} \leftarrow a^{(2)}$

Repeat : { compute predictions:

$$\partial w^{(1)} = (\hat{y}^{(i)}, i=1, 2, \dots, m) \quad \partial b^{(1)} - \partial J$$

$$\text{update: } \begin{aligned} \hat{\frac{\partial w^{(1)}}{\partial w^{(1)}}} &= w^{(1)} - \frac{\partial b^{(1)}}{\partial b^{(1)}} \\ b^{(1)} &= b^{(1)} - \alpha \frac{\partial b^{(1)}}{\partial b^{(1)}} \end{aligned}$$

Formulas for Computing derivatives:

Forward Propagation:

$$\begin{aligned} z^{(1)} &= w^{(1)}X + b^{(1)} \\ a^{(1)} &= g^{(1)}(z^{(1)}) \\ z^{(2)} &= w^{(2)}a^{(1)} + b^{(2)} \\ a^{(2)} &= g^{(2)}(z^{(2)}) = \sigma(z^{(2)}) \end{aligned}$$

Back Propagation:

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

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

$$dw^{(2)} = \frac{1}{m} dz^{(2)} A^{(1)T}$$

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

$$(n^{(2)}, 1)$$

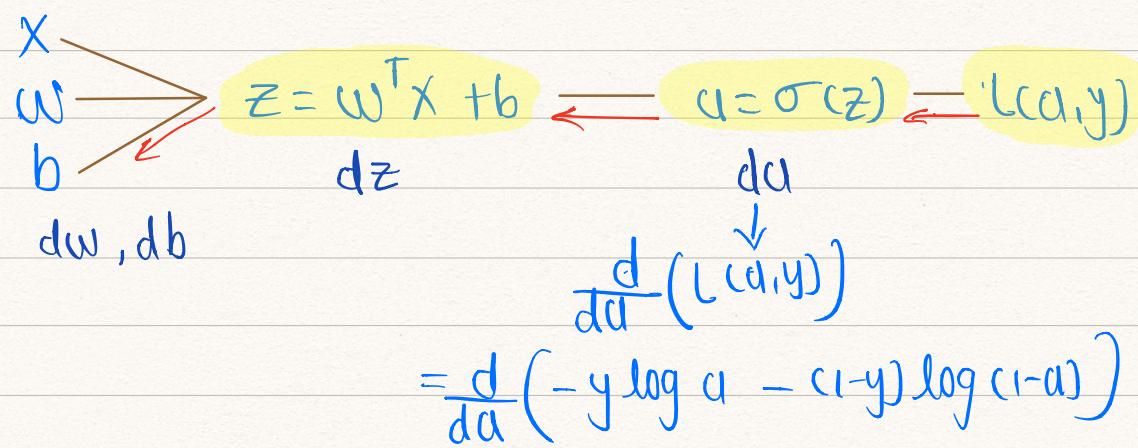
$$dz^{(1)} = \underbrace{w^{(2)T} dz^{(2)}}_{(n^{(1)}, m)} * \underbrace{g^{(1)'}(z^{(1)})}_{\text{elementwise product}} (n^{(1)}, m)$$

$$db^{(1)} = \frac{1}{m} \cdot \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims=T})$$

$$d\omega^{(1)} = \frac{1}{m} dz^{(1)} x^T$$

Backpropagation Intuition:

Computing Gradient:



$$du = -y/a + \frac{1-y}{1-a}$$

$$dz = da \cdot g'(z) \xrightarrow{g(z) \rightarrow \text{sigmoid}} \sigma(z)$$

$$a = \sigma(z)$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial u} \cdot \frac{\partial u}{\partial z} \uparrow \frac{d}{dz}(O(z))$$

$$dw = dz \cdot x$$

$$dh = dz$$

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

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

$$(a^{(2)}, y)$$

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

$$z^{(2)} = w^{(2)T} a^{(1)} + b^{(2)}$$

$$da^{(2)}$$

$$dz^{(2)}$$

$$dw^{(2)} = dz^{(2)} \cdot a^{(1)T}$$

$$db^{(2)}$$

$$a^{(2)} - y$$

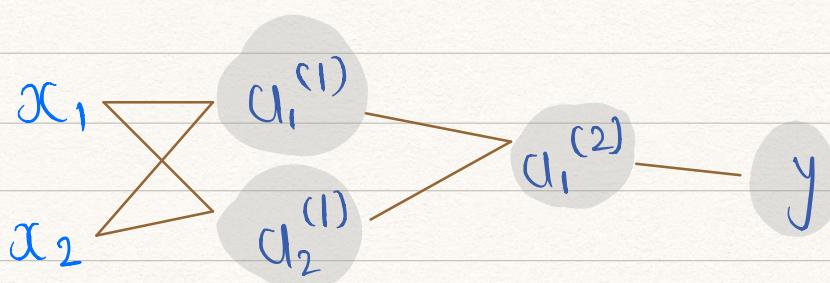
$$w^{(2)}, b^{(2)}$$

$$dZ^{(1)} = \frac{(n^{(1)}, n^{(2)})}{(n^{(2)}, 1)} \cdot \underbrace{\frac{dZ^{(1)}}{(n^{(2)}, 1)} * g^{(1)}(z^{(1)})}_{\text{element wise product}}$$

$$\frac{w^{(2)}}{(n^{(2)}, n^{(1)})}$$

Random Initialization:

What happens if you initialize weights to zero?



$$w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(1)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$u_1^{(1)} = u_2^{(1)}$$

$$dz^{(1)} = dz^{(2)}$$

$$dw = \begin{pmatrix} u & v \\ u & v \end{pmatrix}$$

$$W^{(1)} = W^{(1)} - 2dw$$

$$W^{(1)} = \begin{bmatrix} - & - & - & - \\ - & - & - & - \end{bmatrix}$$

identical rows.

If you initialize with zero, hidden layer will be symmetric and we will keep on computing same function, regardless of no. of layers.

So we do...

$$W^{(1)} = \text{np.random.randn}(2, 2) * 0.01$$

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

multiply with small no.
we do not have symmetry breaking problem

$$W^{(2)} = \text{random}$$
$$b^{(2)} = 0$$

→ If w values are too large , z will be very big or very small. in which case we will not get good gradient decent because it will have dw close to zero.

(end up at flat part of activation function)