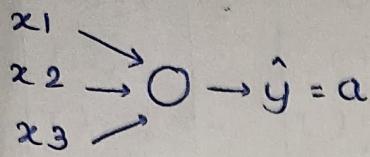


• NEURAL NETWORKS OVERVIEW



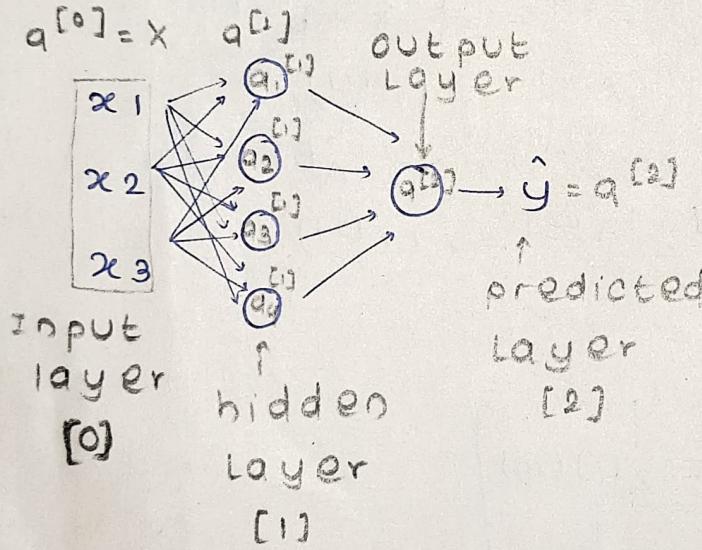
$$x \rightarrow w \rightarrow z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

- [i] → ith layer
- (i) → ith training example

$$\begin{array}{c} x \\ w^{[1]} \rightarrow z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \\ b^{[1]} \end{array} \quad \begin{array}{c} w^{[2]} \rightarrow z^{[2]} = w^{[2]} x + b^{[2]} \\ b^{[2]} \end{array}$$

$L(a^{[2]}, y) \leftarrow a^{[2]} = \sigma(z^{[2]}) \leftarrow$

• NEURAL NETWORK REPRESENTATION:



• Hidden Layer: values
not observed in training set

$$a^{[L]} = \begin{bmatrix} a_1^{[L]} \\ a_2^{[L]} \\ \vdots \\ a_n^{[L]} \end{bmatrix}$$

Hidden Layer: $w^{[1]}, b^{[1]}$
nodes - (4, 3); (4, 1)
↑ ↑
inputs input

output: $w^{[2]}, b^{[2]}$
(1, 4); (1, 1)

node: $\underbrace{w^T x + b}_z$ and $\sigma(z)$

a vector transpose

layer 1 → node 1: $z_1^{[1]} = \underbrace{w_1^{[1] T} x + b_1^{[1]}}_{\text{node } 1} \& a_1^{[1]} = \sigma(z_1^{[1]})$
node 2: $z_2^{[1]} = \underbrace{w_2^{[1] T} x + b_2^{[1]}}_{\text{node } 2} \& a_2^{[1]} = \sigma(z_2^{[1]})$

$$\underbrace{\begin{bmatrix} -\omega_1^{[1]T} \\ -\omega_2^{[1]T} \\ -\omega_3^{[1]T} \\ -\omega_4^{[1]T} \end{bmatrix}}_{\omega^{[0]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} \omega_1^{[1]T} + b_1^{[1]} \\ \omega_2^{[1]T} + b_2^{[1]} \\ \omega_3^{[1]T} + b_3^{[1]} \\ \omega_4^{[1]T} + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

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

Given input x : $x = a^{[0]}$

$$(4,1) \rightarrow z^{[1]} = \underbrace{\omega^{[1]} x}_{(4,4)} + b^{[1]} \quad (4,1)$$

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

$$(1,1) z^{[2]} = \underbrace{\omega^{[2]} a^{[1]}}_{(1,4)} + b^{[2]} \quad (1,1)$$

$$a^{[2]} = \sigma(z^{[2]}) \quad (1,1)$$

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

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

$$\vdots^{(n)} \longrightarrow a^{[2](n)} = \hat{y}^{(n)}$$

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad n \times m$$

for $i = 1 \text{ to } m$:

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

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

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

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

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

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

$$z^{[2]} = \omega^{[2]} A^{[1]} + b^{[2]}$$

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

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

$$A^{[1]} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ a^{1} & a^{[1](2)} & a^{[1](3)} & \dots & a^{[1](m)} \end{bmatrix}$$

↑ input features
↓ training examples

• VECTORIZED IMPLEMENTATION.

$$z^{1} = \omega^{[1]} x^{(1)} + b^{[1]} ; z^{[1](2)} = \omega^{[1]} x^{(2)} + b^{[1]}, z^{[1](3)} = \omega^{[1]} x^{(3)} + b^{[1]}$$

$$\omega^{[1]} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \omega^{[1]} x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \omega^{[1]} x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \omega^{[1]} x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$x = \begin{bmatrix} 1 & 1 & 1 & \dots \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ 1 & 1 & 1 & \dots \end{bmatrix} \Rightarrow \omega^{[1]} x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

* $\omega^{[1]} x^{(i)} = z^{[1](i)}$

$\Rightarrow z^{[1]} = \omega^{[1]} x + b^{[1]}$

* x stacked in column vector $\Rightarrow z$ stacked in col-vec
(input) (output)

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

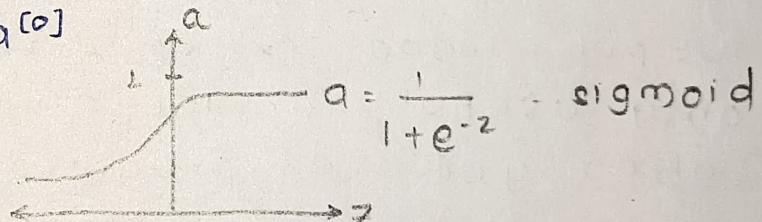
• ACTIVATION FUNCTION: (sigmoid is one choice)

$$z^{[1]} = \omega^{[1]} x + b^{[1]} \dots a = a^{[0]}$$

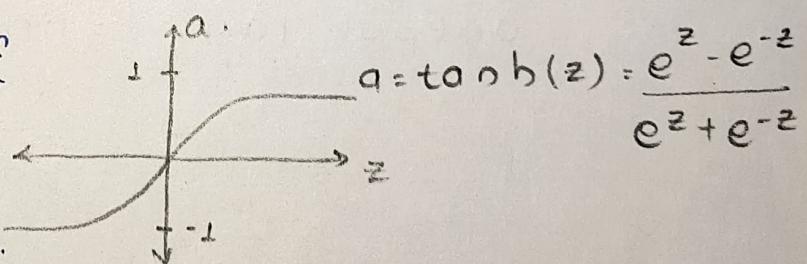
$$a^{[1]} = \sigma(z^{[1]})$$

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

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



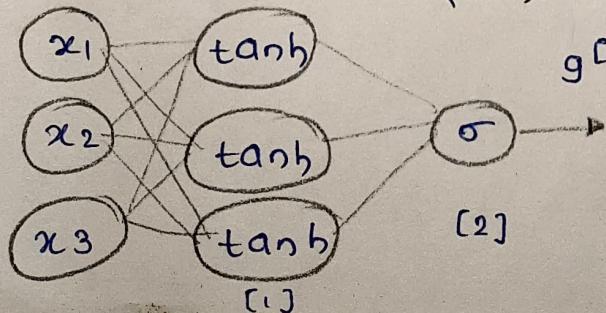
NOTE: For hidden layers $g(z^{[0]}) = \tanh(z^{[0]})$ works better than $\sigma(z^{[0]})$ as mean is closer to zero.



Exception: σ activation function works better for output layer in binary classification $[0 \leq y \leq 1]$

$$g^{[1]}(z^{[0]}) = \tanh(z^{[0]})$$

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



• Limitations:

If $z \gg 0$ or $z \ll -\infty$ then slope ≈ 0 slows gradient descent.

better choice = ReLU (Rectified Linear Unit)

RULES :

(i) FOR binary classification

output layer ... $\hat{y} \in \{0,1\}$

use : SIGMOID FUNCTION $0 \leq \hat{y} \leq 1$

(ii) For all other cases ... activation fx? = ReLU
when $z = -ve$, slope = 0
* Leaky ReLU better than ReLU

NOTE: **ReLU** faster than sigmoid / tanh

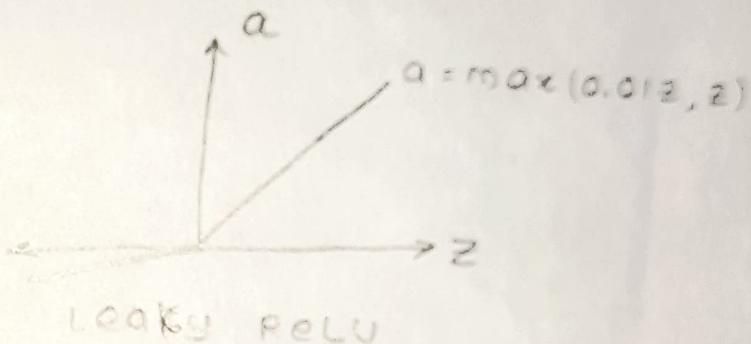
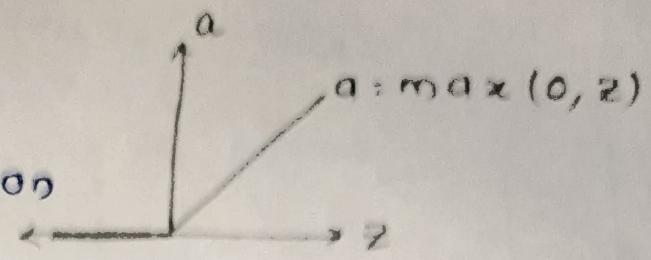
σ : 0 to 1

tanh : -1 to 1

ReLU \Rightarrow introduces

non-linearity

$\hat{y} \geq 0$



* Linear / identity activation function : $g(z) = z$

output = linear fx? of input

- only useful for ML on a regression problem
(when $y \in \mathbb{R}$... ex. prediction of housing price)

Here : hidden layers = ReLU or tanh

Output layer = linear activation function

3

• DERIVATIVES OF ACTIVATION FUNCTIONS:

sigmoid:

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

$$\therefore \frac{d}{dz} g(z) = \frac{1}{(1+e^{-z})} \left(1 - \frac{1}{1+e^{-z}} \right)$$

slope = $g(z)(1-g(z)) = g'(z)$
of $g(z)$ at z

$$g'(z) = a(1-a) \dots g(z)=a$$

If: (i) $z \gg +ve$ no ... $g'(z)=0$

(ii) $z \ll -ve$ no ... $g'(z)=0$

$$(iii) z=0 \dots g'(z) = \frac{1}{4}$$

tanh:

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

$$\therefore g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ = 1 - [\tanh(z)]^2$$

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

If: (i) $z \gg +ve$ no ... $g'(z)=0$

(ii) $z \ll -ve$ no ... $g'(z)=0$

$$(iii) z=0 \dots g'(z)=1$$

ReLU:

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

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

* technically undefined for $\boxed{z=0}$

leaky ReLU: $g(z) = \max(0.01z, z)$

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

• COMPUTING GRADIENTS:
→ Logistic Regression

$$a = \hat{y} \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

x

$$\omega \rightarrow z = \omega^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

$\frac{d\omega}{dz} = d\omega \cdot x$

$\frac{db}{dz} = dz$

$\frac{da}{dz} = a - y$

$L(a, y) = -[y \log a + (1-y) \log(1-a)]$

$$\frac{d\omega}{dz} = d\omega \cdot x$$

$$\frac{db}{dz} = dz$$

$$dz = da \cdot g'(z)$$

$$g(z) = \sigma(z)$$

$$a = \sigma(z)$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{da}{dz}$$

$$\frac{da}{dz} = \frac{d}{da} \left[-\frac{y}{a} + \frac{(1-y)}{1-a} \right]$$

$$\frac{d}{da} g(z) = g'(z)$$

$$\frac{d}{dz} g(z) = g'(z)$$

* Neural network gradients

x

$$\omega^{[1]} \rightarrow z^{[1]} = \omega^{[1]} x + b^{[1]} \rightarrow \sigma(z^{[1]}) = a^{[1]} \rightarrow z^{[2]} = \omega^{[2]} a^{[1]} + b^{[2]}$$

$$d\omega^{[1]} = d\omega^{[1]} x^T$$

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

$$d\omega^{[2]} = d\omega^{[2]} a^{[1]}^T$$

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

$$L(a^{[2]}, y) \leftarrow a^{[2]} = \sigma(z^{[2]})$$

$$dz^{[2]} = \omega^{[2]} T \cdot d\omega^{[2]} * g^{[2]}'(z^{[1]})$$

NOTE: dimensions $\rightarrow \omega^{[2]} = (\eta^{[2]}, \eta^{[1]})$

$$z^{[2]}, d\omega^{[2]} = (\eta^{[2]}, 1) \quad ; \quad z^{[1]}, dz^{[1]} = (\eta^{[1]}, 1)$$

$$dz^{[1]} = \omega^{[2]} T \cdot d\omega^{[2]} * g^{[2]}'(z^{[1]})$$

$$(\eta^{[1]}, 1) = (\eta^{[1]}, \eta^{[2]}), (\eta^{[2]}, 1) * (\eta^{[1]}, 1)$$

[IMP] * \Rightarrow element wise product

np.dot(a, b) \Rightarrow matrix multiplication.

Summary of Gradient Descent

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

$$d_w^{[2]} = d_z^{[2]} a^{[2]T}$$

$$d_b^{[2]} = d_z^{[2]}$$

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

$$d_w^{[1]} = d_z^{[1]} x^T$$

$$d_b^{[1]} = d_z^{[1]}$$

axis=0
Horizontal

4

vectorized

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

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

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

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

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

for back prop

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

$$d_w^{[2]} = \frac{1}{m} d_z^{[2]} A^{[1]T}$$

$$J = \frac{1}{m} \sum_1^m \mathcal{L}(\hat{y}, y) \dots \text{cost}$$

$$d_b^{[2]} = \frac{1}{m} np.\text{sum}(d_z^{[2]}, \underbrace{\text{axis}=1, \text{keepdims=True}}_{\text{vertical summing}})$$

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

$$d_w^{[1]} = \frac{1}{m} d_z^{[1]} x^T \quad \text{horizontal}$$

$$d_b^{[1]} = \frac{1}{m} np.\text{sum}(d_z^{[1]}, \underbrace{\text{axis}=1, \text{keepdims=True}}_{\text{vertical}})$$

- If all weights are initialized to zero, all hidden units compute the same function

* Large z slows down learning

RANDOM INITIALIZATION

$$w^{[1]} = np.random.rand(2, 2) * 0.01$$

$$b^{[1]} = np.zeros((2, 1))$$

$$w^{[2]} = np.random.rand(1, 2) * 0.01$$

$$b^{[2]} = 0$$

NOTE: $w^{[k]}$ dimensions \Rightarrow (no. of neurons in k^{th} layer, no. of inputs of layer)

$b^{[k]}$ dimensions \Rightarrow (no. of neurons in k^{th} layer, 1)