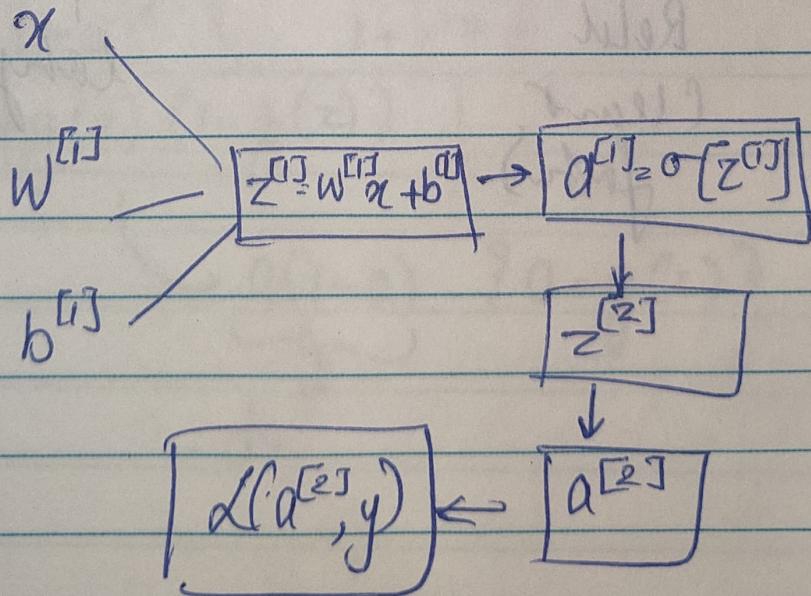
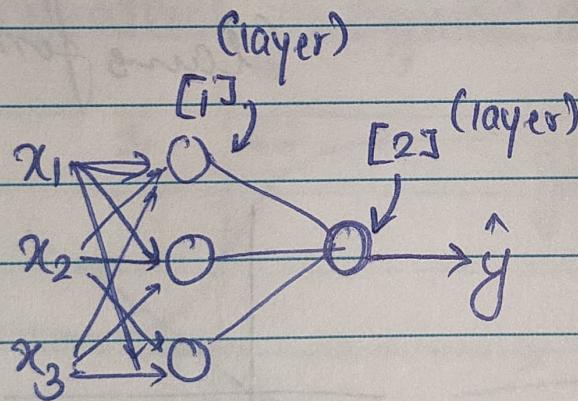
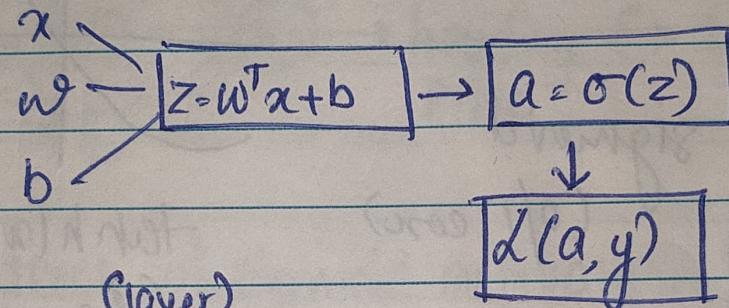
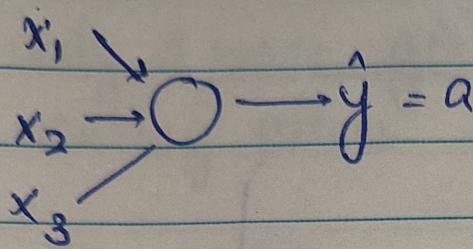


Week 3

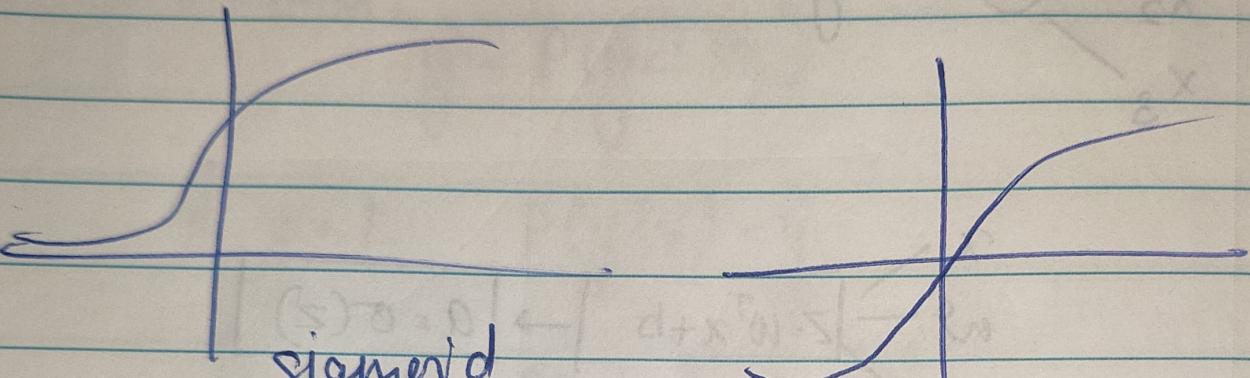
Lec-1 Neural Network



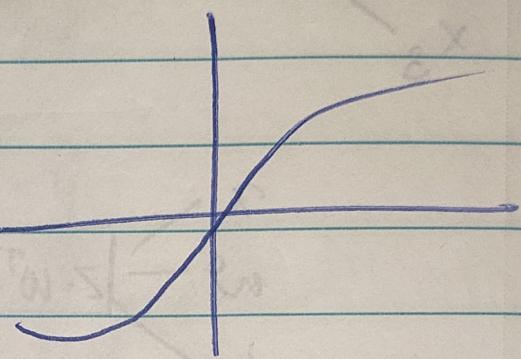
L-2

representation

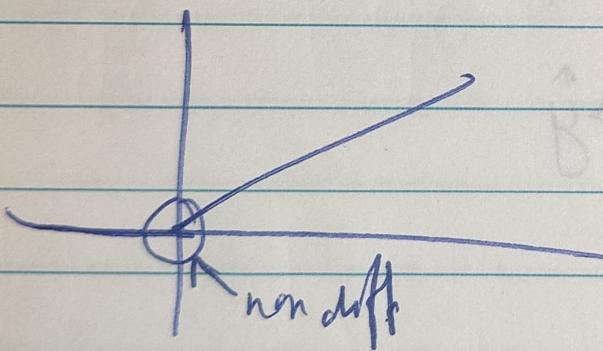
Activation fn.



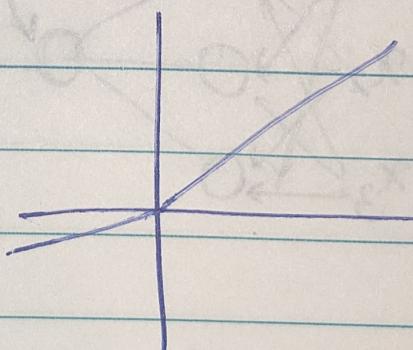
sigmoid
(0/1 case)



$\tanh(x)$
(langs festw.)



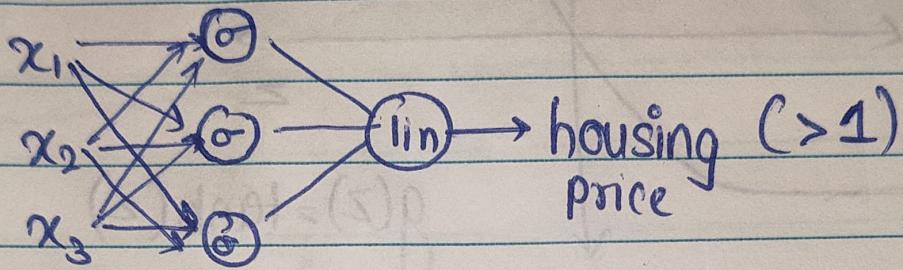
ReLU
(langs
festw.)



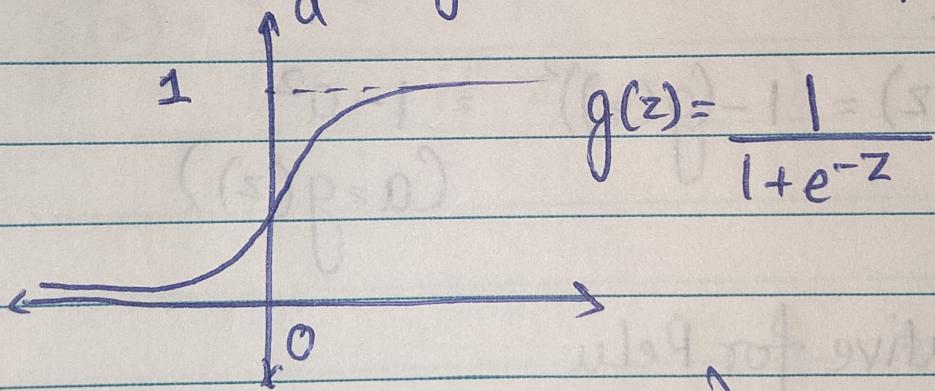
leaky relu

Why non-linear activation functions?

only 1 place $g(z) = z$ used
housing prices

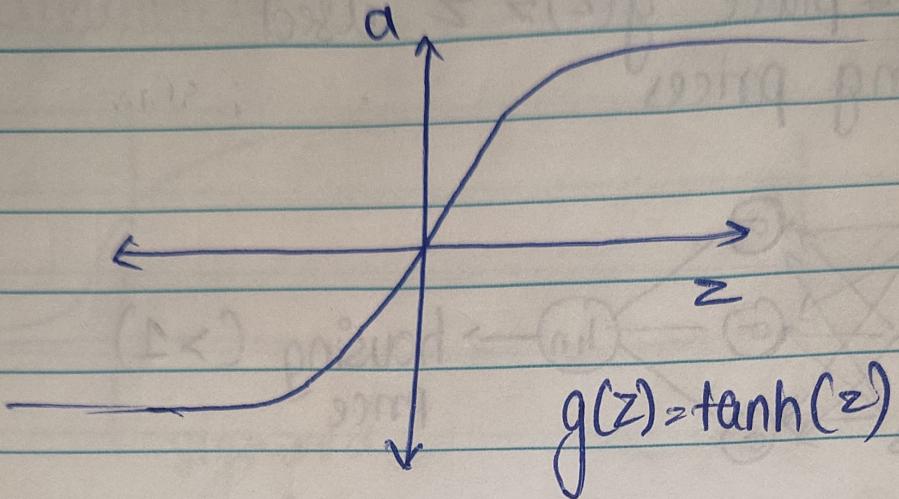


Derivatives of sigmoid activation fns.



$$\frac{dg(z)}{dz} = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right)$$
$$= \frac{g(z)(1-g(z))}{g(z)(1-g(z))}$$
$$= \frac{a(1-a)}{da} \quad \{a=g(z)\}$$

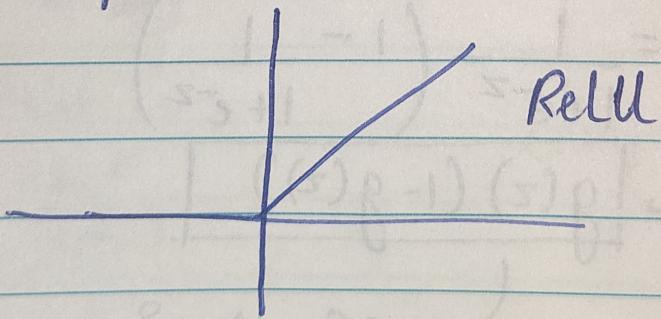
derivative for $\tanh(x)$ function



$$\frac{d g(z)}{dz} = 1 - (\tan(z))^2$$

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

derivative for ReLu

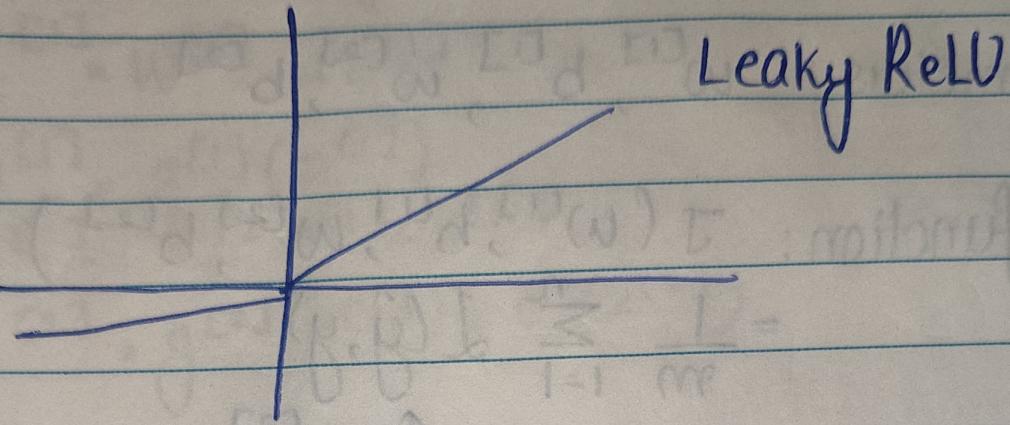


$$g(z) = \max(0, z)$$

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

undef. if $z = 0$

derivative for Leaky ReLU



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

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

gradient descent for neural networks

parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

Cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$

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

\hat{y} ← $a^{[2]}$

gradient descent:

Repeat {

 Compute $\hat{y}^{[1]}(g^{[1]}, i-1, \dots, m)$
 $dw^{[1]} = \frac{dJ}{dw^{[1]}}, db^{[1]} = \frac{dJ}{db^{[1]}}$

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

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

⋮

}

Forward propagation

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

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

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

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

Back propagation

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

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

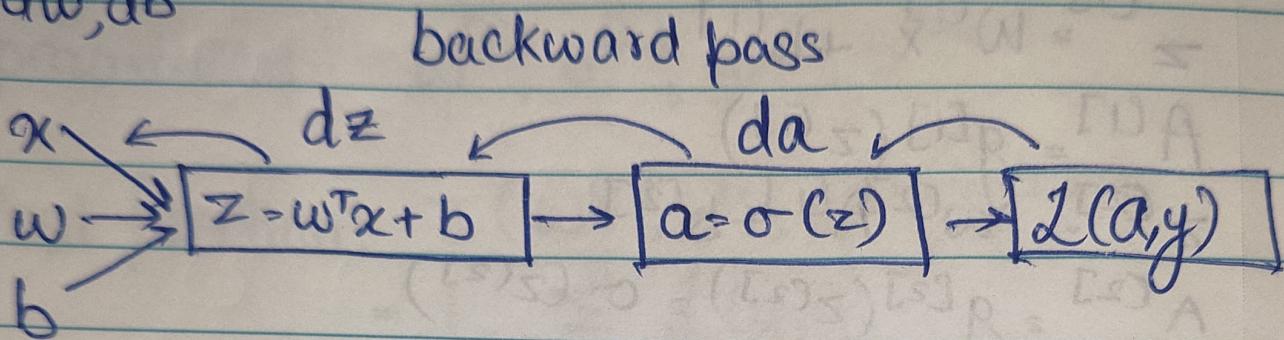
$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdms} = \text{True})$$

$$dz^{[1]} = \underbrace{w^{[2] \top} dz^{[2]}}_{(m^{[1]}, m)} \times \underbrace{g^{[2]'}(z^{[2]})}_{\text{elementwise product}}$$

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

Backpropagation Intuition

$d\mathbf{w}, d\mathbf{b}$



Forward Pass

$$da = \frac{d(L(a, y))}{da} = -y \log a - (1-y) \log(1-a)$$

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

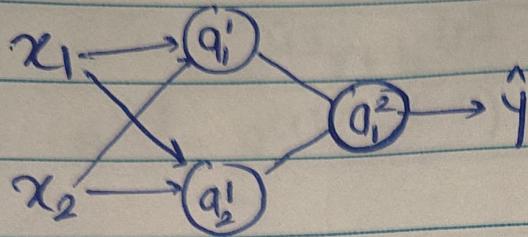
$$dz = da \cdot g'(z) = \left(\frac{-y}{a} + \frac{1-y}{1-a} \right) (a(1-a))$$

$$\left(\frac{df}{da} \frac{da}{dz} \right)$$

$$d\mathbf{w} = x dz$$

$$d\mathbf{b} = dz$$

Random initialization



if $w, b \rightarrow 0$ matrices

$$a_1^{[i]} = a_2^{[i]} \rightarrow dz_1^{[i]} = dz_2^{[i]}$$

Both hidden units have same influence
on next layer - computing same function

\therefore we init. randomly

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

$$b^{[1]} = \text{np. zeros}((3, 1))$$

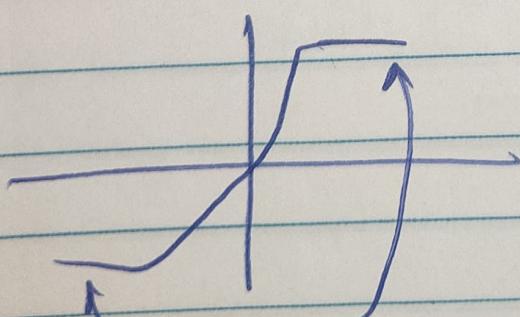
$$w^{[2]} = \dots$$

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

prefer to
init to smaller
vals.

because if

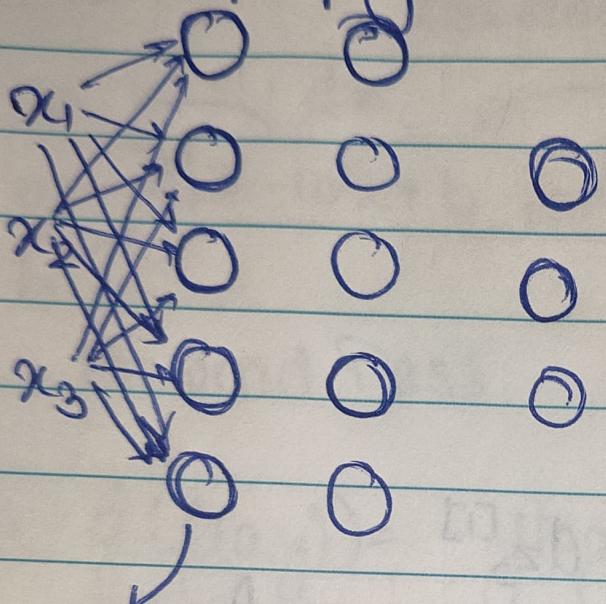
w is big,
 z is big



slow slope
 \therefore slow training

a may be
very large or
v. small

Deep L-Layer Neural Network forward propagation



$$x: z^{[0]} = w^{[1][0]} a + b^{[0]}$$

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

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

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

⋮

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

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

vectorized version

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

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

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

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

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

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

explicit for loop required for l to l

$n^{[1]}$ (no of hidden units in layer 1)

matrix dimensions

$l=5$ (5 layer NN)

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]} \rightarrow \text{bias}$$

$$n^{[1]} = 3, n^{[2]} = 5, n^{[3]} = 4, n^{[4]} = 2, n^{[5]} = 1$$

activation vector for $[1]$

\therefore dimension: $(3, 1)$ ($n^{[1]}, 1$)

$x \rightarrow$ 2 input features

$\therefore (2, 1) \text{ dim } [n^0, 1]$

$$\therefore [:] = w^T [:]$$

$(3 \times 2) \text{ dim}$

$$w^{[1]} = n^{[1]}, n^{[0]}$$

$$w^{[2]} = n^{[2]}, n^{[1]} = (5, 3) \text{ dim.}$$

$$\underline{\underline{w^{[l]} = (n^{[l]}, n^{[l-1]}) \text{ dim}}}$$

$$\underline{\underline{b^{[l]} = (n^{[l]}, 1) \text{ dim}}}$$

$$\underline{dw^{[l]} = (n^{[l]}, n^{[l-1]})} \quad db^{[l]} = (n^{[l]}, 1)$$

same as
w

same as
b

Vectorized Implementation

$$z^{[l]} = w^{[l]} \cdot x + b^{[l]}$$

$[n', n^o]$ $[n^o, m]$
(stacked)

$[n^{[l]}, m]$
(stacked
as $\begin{bmatrix} z^{[1]} & z^{[2]} & \dots \end{bmatrix}$)

$$z^{[l]}, A^{[l]}: [n^{[l]}, m] \text{ dim}$$

$$dz^{[l]}, dA^{[l]}:$$