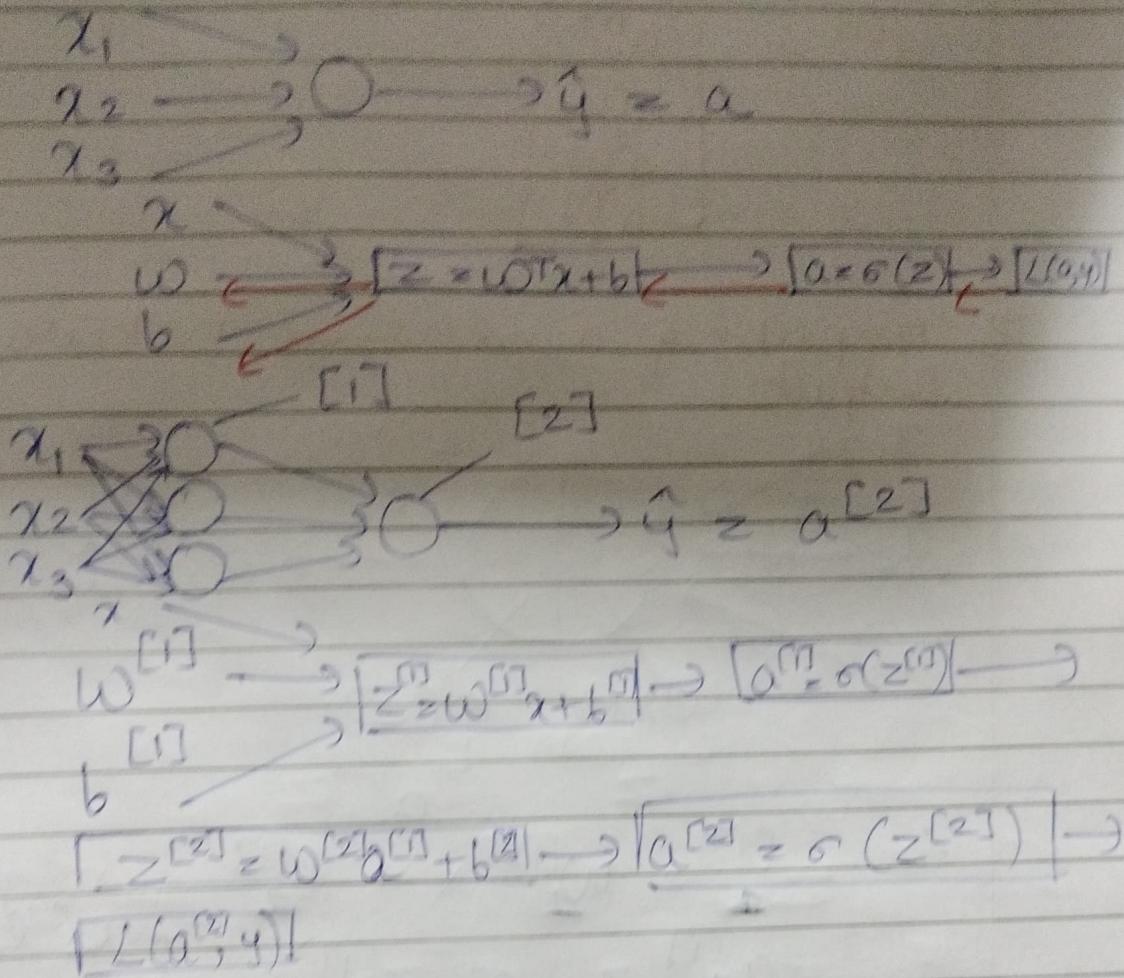
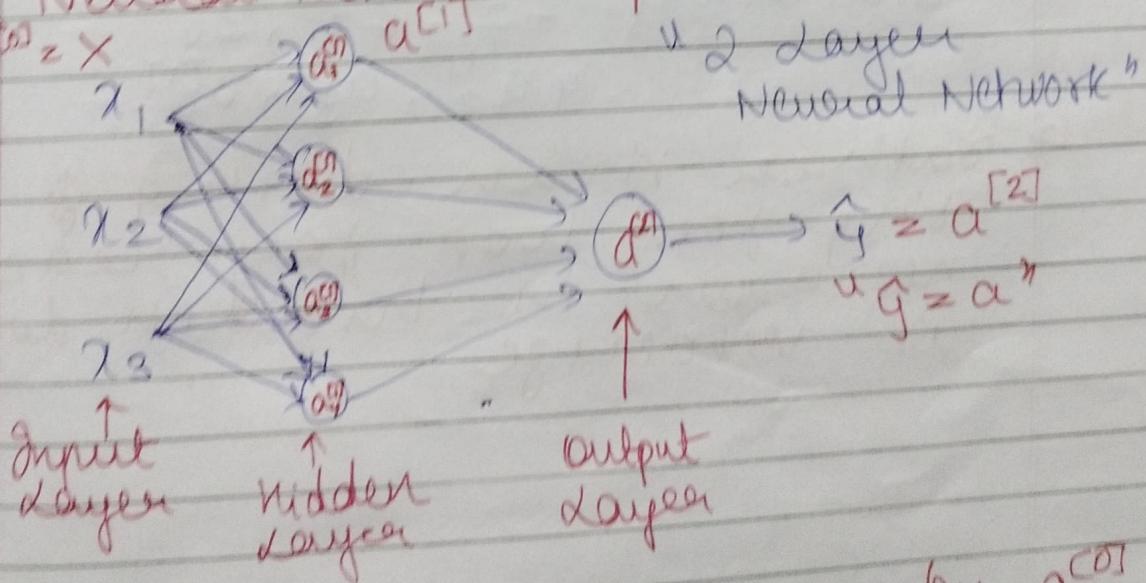


Shallow neural network



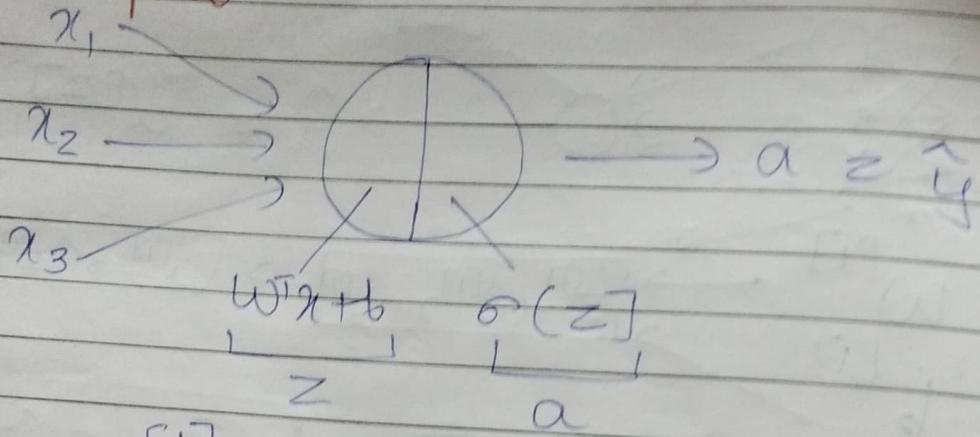
Neural Network representation



activations of input layer by $a^{[0]}$
of hidden layer $a^{[1]}$

$$a^{[1]} = \begin{bmatrix} a_1^{[0]} \\ a_2^{[1]} \\ a_3^{[2]} \\ a_4^{[3]} \end{bmatrix}$$

Computing a neural network's output



$$z_1^{[1]} = w_1^{[0]T} x + b_1^{[1]} \quad \left. \begin{array}{l} \text{Similarly} \\ \text{for } z_2^{[2]}, z_3^{[3]}, z_4^{[4]} \end{array} \right\}$$

$$a_1^{[1]} = \sigma(z_1^{[1]}) \quad \left. \begin{array}{l} \text{Similarly} \\ \text{for } a_2^{[2]}, a_3^{[3]}, a_4^{[4]} \end{array} \right\}$$

→ rather than doing \$ using for loop for each neuron in hidden layer we vectorize

ω^* (ω^{IR})

1962年
1月1日

$$z = \begin{bmatrix} w_1^{(0)\top} x + b_1^{(0)} \\ w_2^{(0)\top} x + b_2^{(0)} \\ w_3^{(0)\top} x + b_3^{(0)} \\ w_4^{(0)\top} x + b_4^{(0)} \end{bmatrix} = \begin{bmatrix} z_1^{(0)} \\ z_2^{(0)} \\ z_3^{(0)} \\ z_4^{(0)} \end{bmatrix}$$

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

Given input x:

$$\begin{aligned} \rightarrow & Z^{(4,1)} = W^{(4,3)} \chi^{(3,1)} + b^{(4,1)} \\ \rightarrow & a^{(1)} = g(Z^{(1)}) \\ \rightarrow & Z^{(2,1)} = W^{(2,2)} a^{(1,1)} + b^{(2,1)} \\ \rightarrow & a^{(2)} = g(Z^{(2)}) \end{aligned}$$

Vectorizing across multiple examples

$$\begin{array}{l} x \longrightarrow a^{[2]} = \hat{y} \\ x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)} \\ x^{(2)} \longrightarrow a^{2} = \hat{y}^{(2)} \\ \vdots \\ x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)} \end{array}$$

(2)(i) — training eg. i
a layer

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

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

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

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

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

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

as seen
earlier

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

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

no of
neurons
in hidden
layer

~~no~~ m
(training egs)

Justification for vectorized implementation

$$z^{[1]}(1) = w^{[1]}x^{(1)} + b^{[1]} \rightarrow 0$$

$$z^{[1]}(2) = w^{[1]}x^{(2)} + b^{[1]} \rightarrow 0$$

$$z^{[1]}(3) = w^{[1]}x^{(3)} + b^{[1]} \rightarrow 0$$

$$w^{[1]} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix}$$

$$w^{[1]}x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{[1]}x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

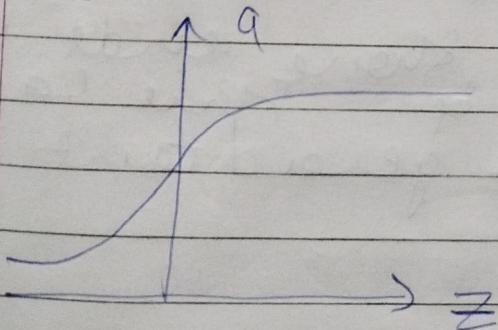
$$w^{[1]}x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{[1]} \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$z = \begin{bmatrix} z^{[0]}(1) & z^{[0]}(2) & z^{[0]}(3) & \dots \end{bmatrix} = \underbrace{\mathbf{z}}_{\mathbf{Z}}$$

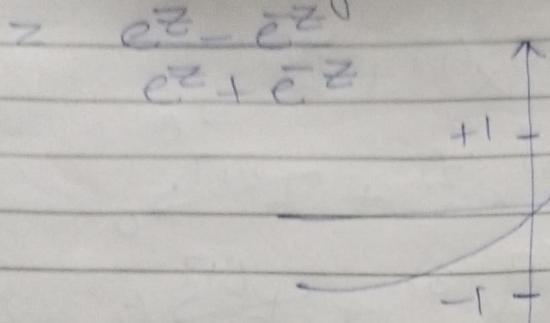
Activation function

→ So far we have used sigmoid function



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

→ A function that works better than sigmoid is $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$



This works better for hidden units than $\sigma(z)$ as values b/w +1 & -1 are the mean of activations that come out of hidden layer are closer to 0 mean. This makes learning for next layer a little bit easier

→ But $\tanh(z)$ not so used for o/p layer instead sigmoid used as o/p must be b/w 0 & 1

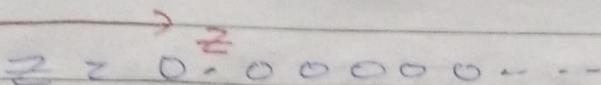
→ Sometimes activations are diff for diff layers eq $g^{[1]}(z), g^{[2]}(z)$

→ downside of both sigmoid & \tanh func is that if z is very large or small gradient of derivative of slope becomes very small & slope ends up close to 0 & this can slow down gradient descent.

→ ∵ A more popular choice is rectified linear Unit (ReLU)

 $\uparrow a$

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

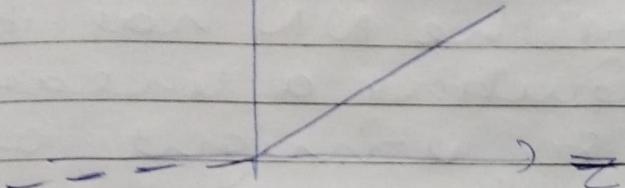


→ one disadvantage is that when $z < 0$ derivative is ~~zero~~

→ ∵ another version called leaky ReLU is used

 $\uparrow a$

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



Sigmoid — Almost never use
only use in O/P layer for binary
classification

$\tanh z$ — used in hidden layers
in pref to σ

RELU — when not sure what
to use

Why do we need activation
functions

→ If no activation functions

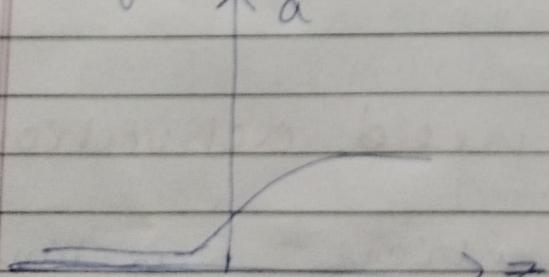
$$a^{[0]} = z^{[0]} = w^{[0]T}x + b^{[0]}$$

$$a^{[1]} = z^{[1]} = w^{[1]T}a^{[0]} + b^{[1]} \\ = (w^{[1]}w^{[0]})x + (w^{[1]b^{[0]}} + b^{[1]}) \\ = w'$$

- if we use this linear activation function the network is just outputting a linear func of input
- if u use a linear activation func / don't have an activation func no matter how many layers the NN has it's all just computing a linear activation function so u might as well have no hidden layers
- Model is no more expressive than standard logistic regression without any hidden layer

Derivatives of activation functions

1) Sigmoid activation function



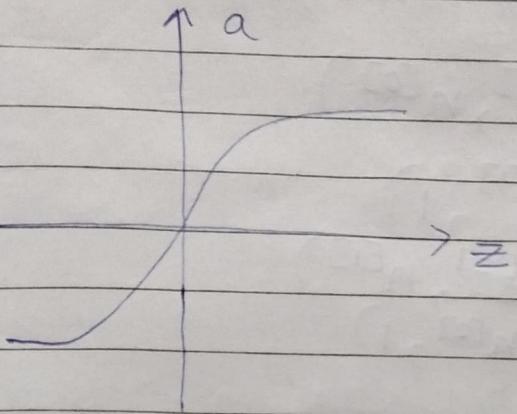
$$g(z) = \frac{1}{1+e^{-z}}$$

$$\frac{d}{dz}(g(z)) = \text{slope of } g(z) \text{ at } z$$

$$z \frac{e^{-z}}{(1+e^{-z})^2} = g(z)(1-g(z))$$

$$= a(1-a)$$

2) Tanh activation function



$$g(z) = \tanh(z)$$

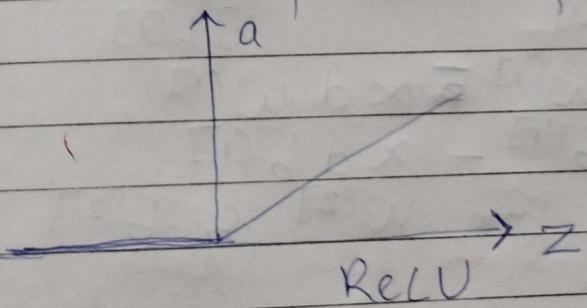
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

$$\therefore a = g(z)$$

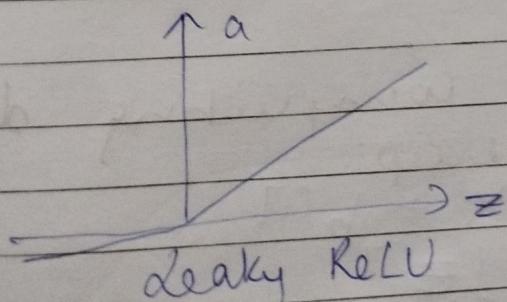
$$\therefore g'(z) = \underline{\underline{1 - a^2}}$$

3) ReLU & Leaky ReLU



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

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



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

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

Gradient descent for neural network with 1 hidden layer

→ Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 $n_2 = n^{[0]} \rightarrow n^{[1]} \rightarrow n^{[2]} = \frac{n^{[2]}}{2} = 2$

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

$$b^{[1]} = (n^{[1]}, 1)$$

$$w^{[2]} = (n^{[2]}, n^{[1]})$$

$$b^{[2]} = (n^{[2]}, 1)$$

→ cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$
 $= \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$
 $\uparrow a^{[2]}$

→ gradient descent:
 compute predictions

repeat $\left(\begin{array}{l} \hat{y}_i^{(i)} \text{ for } i = 1 \text{ to } m \\ \frac{dJ}{dw^{[1]}} \frac{dJ}{db^{[1]}} \end{array} \right) \frac{dJ}{db^{[1]}}$
 until parameters look like they are converging

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

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

similarly for $w^{[2]}$ & $b^{[2]}$

formula for computing derivatives

- forward pass:

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

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

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

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

back propagation

$$d\hat{z}^{[2]} = A^{[2]} - y$$

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

$$dW^{[2]} = \frac{1}{m} d\hat{z}^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np \cdot \text{sum}(d\hat{z}^{[2]}, \text{axis}=1)$$

↑
keeps dimension
 $(h^{[2]}, 1)$

Prevents rank 1 matrix

$$d\hat{z}^{[1]} = \underbrace{w^{[2]T} d\hat{z}^{[2]}}_{n_3 \times n_4 (h^{[1]}, m)} * \underbrace{g^{[1]'}(\hat{z}^{[1]})}_{\text{element wise product}} (n^{[1]}, m)$$

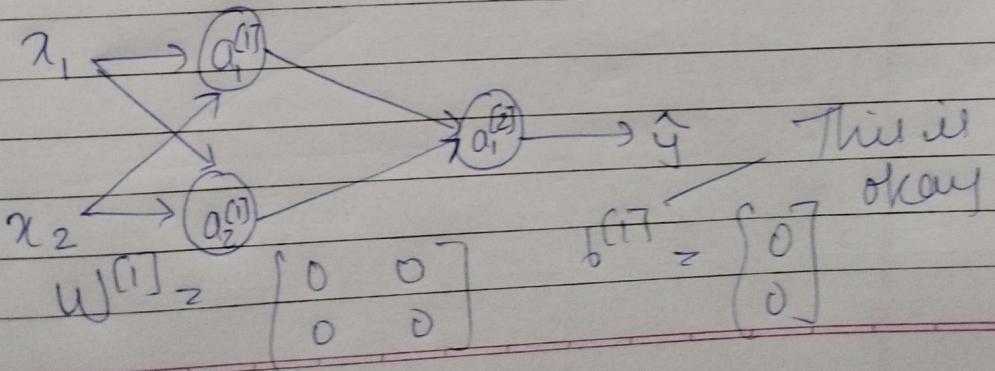
$$dW^{[1]} = \frac{1}{m} d\hat{z}^{[1]} x^T$$

$$db^{[1]} = \frac{1}{m} np \cdot \text{sum}(d\hat{z}^{[1]}, \text{axis}=1)$$

Keep dimension
 $(n^{[1]}, 1)$

$$\begin{aligned} d\hat{z}^{[2]} &= (n^{[2]}, m) \\ d\hat{z}^{[1]} &= (n^{[1]}, m) \end{aligned}$$

Random Initialization



→ problem with initializing
With this mat
 $a_1^{[1]}$ & $a_1^{[2]}$ will be equal

Eq in back prop
 $dz_1^{[1]} = dz_2^{[1]}$

→ We want diff layers to compute
diff functions : we initialize
randomly

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

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

$w^{[2]} = np.random.rand((1, 2)) * 0.01$ we want
 $b^{[2]} = 0$ to initialize
wt to very
small val

→ if w is very big, z will be
very big :-

~~if~~ → more likely to end
up at the flat part
∴ slope is very small