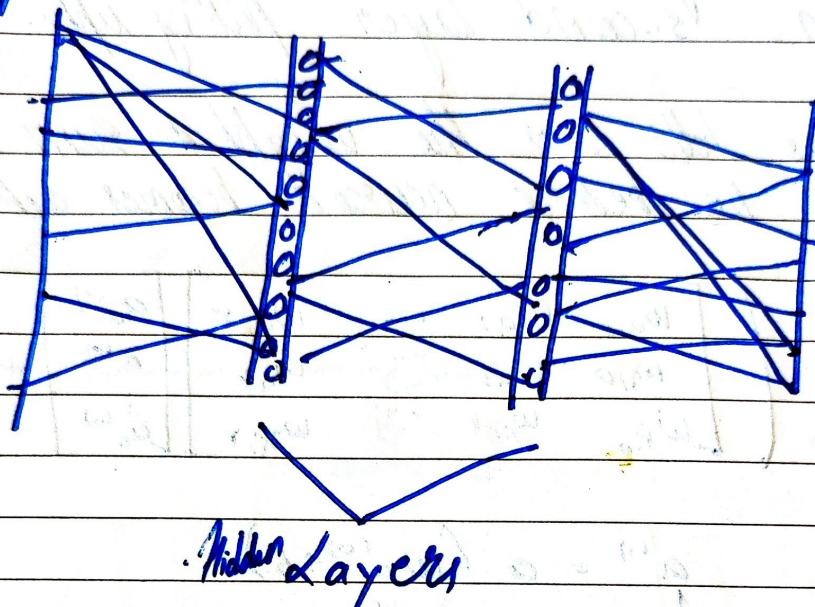


Neural networks

- Chap - 1 , Deep learning.

0 - 1
neuron

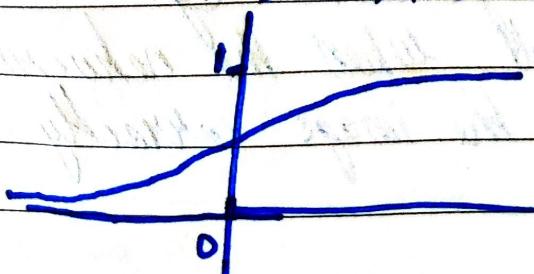
layers :-



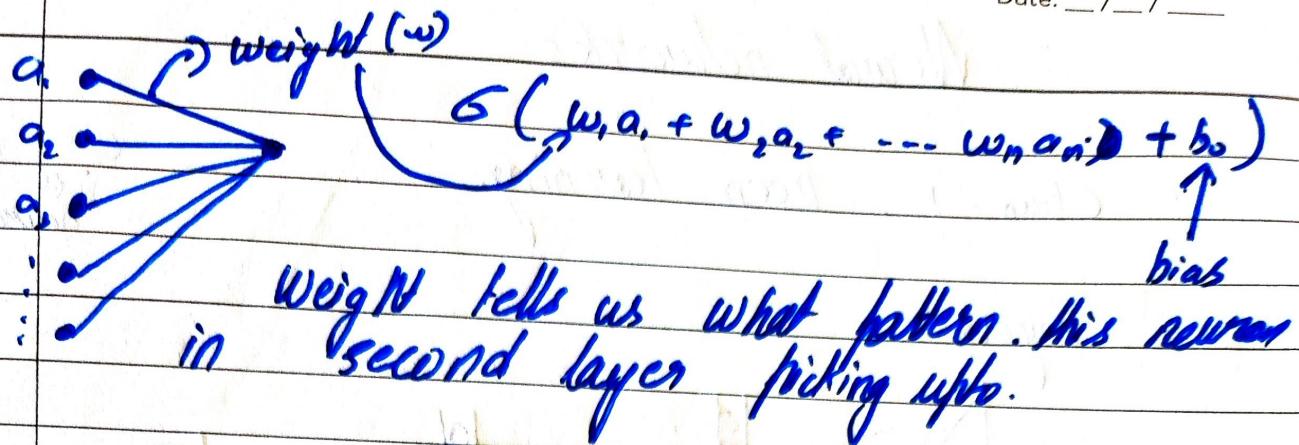
- The activation in 1 layer determines the activation in next layer
- iring of some neurons cause some other neurons to fire.
- ~~Any pattern is easily findable~~
- layers can do the job of dividing looking for components in an input pattern.

- Sigmoid function :- (σ)

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Range (0 - 1)
used for activation values.

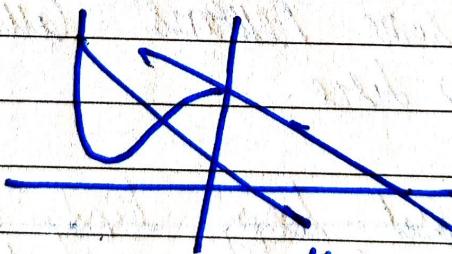


Bias tells how the weighted sum needs to be before neuron becomes active.

$$\sigma \left(\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,r} \\ w_{1,0} & w_{1,1} & \dots & w_{1,r} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n,0} & w_{n,1} & \dots & w_{n,r} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$

$$a^{(1)} = \sigma(w a^{(0)} + b)$$

↓ ↓
 1st layer initial layer



~~Cost functions :-~~ For that you add up the Sq. of the differences between each thresh outputs activations and the value that you want them to have and this is just the cost of a single training programme e.g.

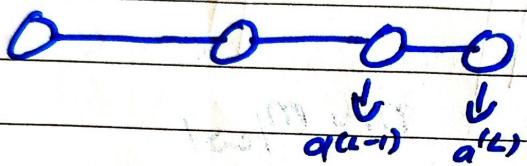
(The sum is small when the network confidently classifies the image correctly)

But it's large when the network doesn't really know what it's doing.

So how you calculate the avg. costs of all the training examples.

- The avg. cost is our measure for how bad our network is.

- Consider 4 networks



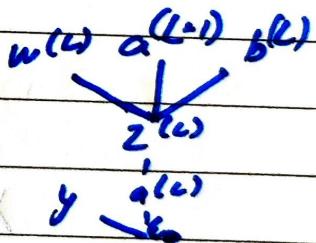
(we are considering only last 2)

$(y) \rightarrow$ desired output

$$\text{cost} \rightarrow C_0 = (a^{(L)} - y)^2$$

$$a^{(L)} = w^{(L)} \cdot a^{(L-1)} + b^{(L)} = z^{(L)}$$

$$a^{(L)} = g(z^{(L)})$$



- Computing gradient derivatives.

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y) \cdot a'(z^{(L)}) \cdot a^{(L-1)}$$

$$z^{(L)} = w^L a^{(L-1)} + b^{(L)}$$

$$\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$$

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$

o Some notation :-

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

\downarrow

n_x dimensional feature vector

$$m \text{ training examples : } \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^m, y^m)\}$$

$$M = M_{\text{train}}, \quad m = m_{\text{test}}$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \\ \downarrow \end{array}$$

$\leftarrow m \rightarrow$

$$X \in \mathbb{R}^{n_x \times m} \quad X.\text{Shape} = (n_x, m)$$

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

$$y \quad Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{Shape} = (1, m)$$

* Logistic Regression Cost function

$$\hat{y}^{(i)} = a(w^T x^{(i)} + b)$$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$.

Want $\hat{y}^{(i)} \approx y^{(i)}$ \rightarrow true label.
 \downarrow
 our output

- Loss (error) function: (for a single training ex.)

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$$\text{If } y=1 \Rightarrow L(\hat{y}, y) = -\log \hat{y}$$

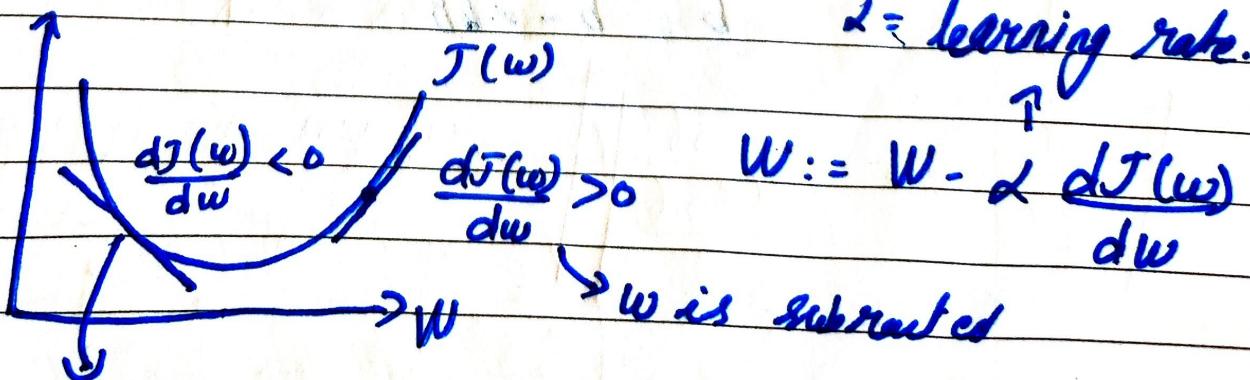
$$\text{If } y=0 \Rightarrow L(\hat{y}, y) = -\log (1-\hat{y})$$

- Cost function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

used for entire training set.

$$= \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

- Gradient descent



Here w is added

• Logistic regression derivatives. (arrow)

$$Z = w_1x_1 + w_2x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

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

$$\frac{da}{dz} = \frac{dL(a, y)}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$$

$$\frac{dz}{dL} = \frac{dL}{da} \cdot \frac{da}{dz}$$

$$= \frac{dL}{da} a(1-a)$$

$$\frac{dz}{dL} = \left[\frac{-y}{a} + \frac{1-y}{1-a} \right]_{a(1-a)} = \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right)$$

$$= \left[-y + ay + a - ay \right] = \frac{e^{-z}}{(1+e^{-z})^2}$$

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

$$w_1 = w_1 - \alpha d w_1$$

$$w_2 = w_2 - \alpha d w_2$$

$$b = b - \alpha d b$$

→ * Logistic regression on m examples.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_j} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_j} l(a^{(i)}, y^{(i)})}_{d w_j^{(i)}}$$

* $J = 0; dw_1 = 0; dw_2 = 0; db = 0$

For $i=1$ to m

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

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

$$J \leftarrow -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \uparrow n=2 \text{ (features)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J/ = m$$

$$dw_1 / = m \quad dw_1 = \frac{\partial J}{\partial w_1}$$

$$db / = m$$

$$dw_2 / = m$$

Q3

- Vectorization [To avoid for loops.]

- Helps to come in faster computation, in very large data sets.

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

Non-vectorized

$$z = 0$$

For i in range(n_m):

$$Z += w^T[x_i] \times x[i]$$

vectorized

$$Z = np.\text{dot}(w, x) + b$$

$$Z += b$$

* Eg:-

$$V = \begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$$

If you want to calculate exponential of each element of V

$$U = \begin{bmatrix} e^{V_1} \\ \vdots \\ e^{V_n} \end{bmatrix}$$

- non-vectorized $\Rightarrow V = np.\text{zeros}((n, 1))$

For i in range(n):

$$V[i] = \text{math.exp}(V[i])$$

- vectorized

import numpy as np

$$U = np.\text{exp}(V)$$

$$J=0, \boxed{dw_1 = 0; dw_2 = 0}, db = 0$$

$$dw = np \cdot \text{zeros}((n_n, 1))$$

for $i=1$ to M :

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

$$d^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log d^{(i)} + (1-y^{(i)}) \log (1-d^{(i)})]$$

$$dz^{(i)} = d^{(i)} [1-d^{(i)}]$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$dw += x^{(i)} d_2^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/M, dw_1 = dw_1/M, dw_2/M, db = db/M$$

$$dw/M$$

• vectorizing logistic regression

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

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

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

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

$$z^{(3)} = w^T x^{(3)} + b$$

$$d^{(3)} = \sigma(z^{(3)})$$

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

$$[z^{(1)} \ z^{(2)} \ z^{(3)} \dots \ z^{(m)}] = w^T \cdot x + [b \ b \dots \ b]_{1 \times m}$$

$$= [w^{(1)} \text{ (row vector)}] \begin{bmatrix} x^{(1)} & x^{(2)} \dots & x^{(m)} \end{bmatrix} + [b \ b \dots]$$

$$Z = [z^{(1)} \ z^{(2)} \ z^{(3)} \dots \ z^{(m)}] = [w^T \cdot x^{(1)} + b \\ w^T \cdot x^{(2)} + b \dots \\ w^T \cdot x^{(m)} + b]$$

$Z = \text{np. dot}(w^T, x) + b$
 (broadcasting in python)

~~$A = f(a)$~~

$$A = [d^{(1)} a^{(2)} a^{(3)} \dots \ a^{(m)}] = \sigma(z)$$

* Vectorizing gradient computation.

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

$$A = [a^{(1)} \dots \ a^{(m)}]$$

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

$$dz = A - y = [d^1 y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$

~~dot = 0~~

$$dw = 0$$

we got
rid of this
for loop

$$\begin{aligned} dw^f &= x^{(1)} dz^{(1)} \\ dw^f &= x^{(2)} dz^{(2)} \\ &\vdots \\ dw^f &= m \end{aligned}$$

$$\begin{aligned} db &= 0 \\ db^f &= dz^{(1)} \\ db^f &= dz^{(2)} \\ &\vdots \\ db^f &= dz^{(m)} \\ db^f &= m \end{aligned}$$

For this,

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$db = \frac{1}{m} np \cdot \text{sum}(dz)$$

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

$$= \frac{1}{m} \left[\begin{array}{c|c} x^{(1)} & \cdots & x^{(m)} \\ \hline 1 & \cdots & 1 \end{array} \right] \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \cdots + x^{(m)} dz^{(m)}]$$

Modified logistic regression.

Fals "itera" in range (1000):

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

$$Z = \phi D.P.\cdot \text{dot}(w.T, x) + b$$

$$A = a(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} \times dZ^T$$

$$db = \frac{1}{m} n.p. \sum (dZ)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

Broadcasting example

Calories from Carbs, proteins, Fats in 100g of different foods:-

| | Apples | Beet | Eggs | Potatoes | |
|---------|--------|-------|------|----------|----------------------|
| Carb | 56.0 | 0.0 | 4.4 | 68.0 | |
| protein | 1.2 | 104.0 | 52.0 | 8.0 | |
| Fat | 1.8 | 135.0 | 99.0 | 0.9 | = A _(3,4) |

Suppose you want to calculate % of calories from crab, []

- means it will sum ↓
- ① $\text{col} = A \cdot \text{Sum}(\text{axis}=0)$; i.e "row vector []"
 - ② percentage = $100 \cdot A / (\text{col})$; if $\text{axis}=0$ it will sum → ie [] vector

* Broad Broadcasting example.

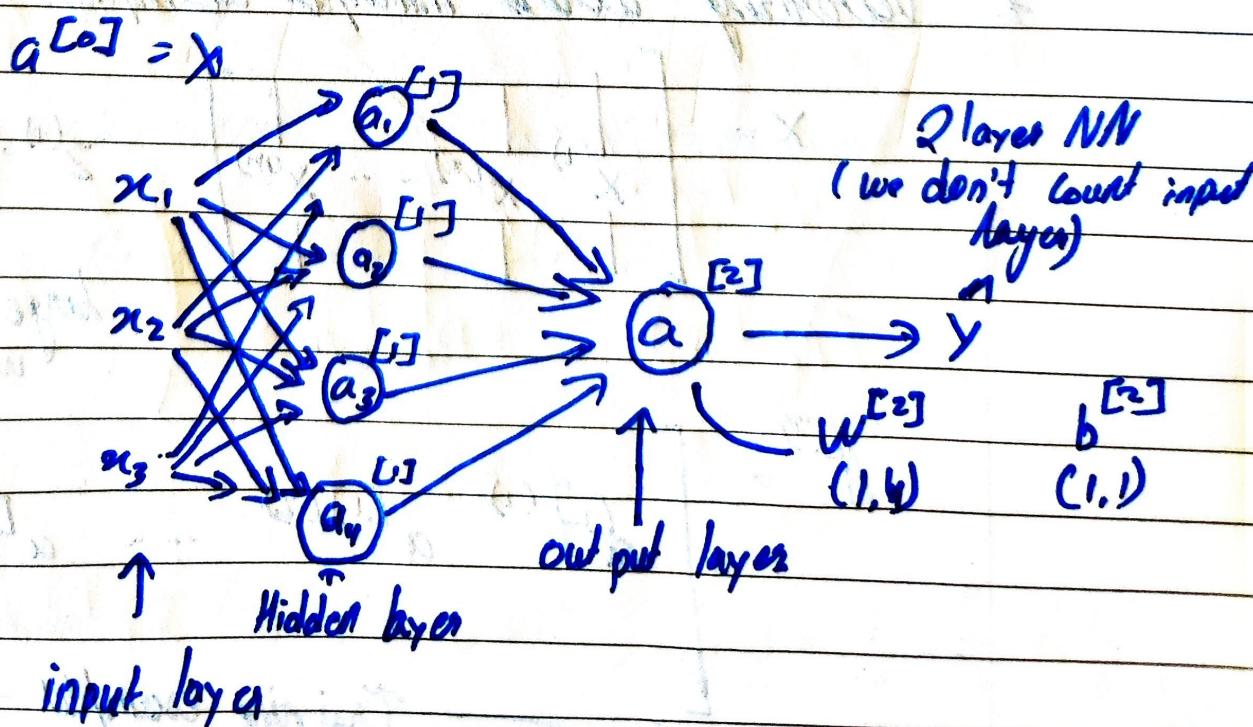
$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \xrightarrow{100} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 5 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

↓

$$\begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

* How Neural network representation



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

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

$$z_2^{[l]} = w^{[l]T} x + b_2^{[l]}$$

$$a_2^{[l]} = \sigma(z_2^{[l]})$$

$[l]$ — layer
 a :

node in layer

$$\begin{bmatrix} - & w_1^{[1]T} \\ - & w_2^{[1]T} \\ - & w_3^{[1]T} \\ - & w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \phi$$

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

* Vectorizing across multiple examples.

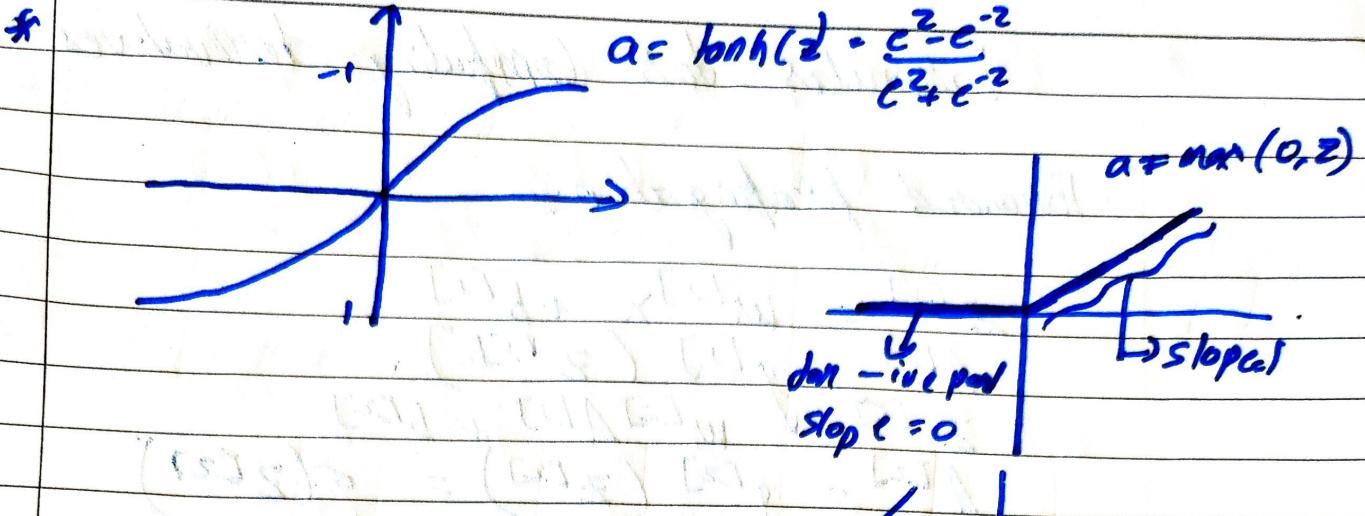
$$x = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

$$z^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ \vdots \\ z_n^{(1)} \end{bmatrix}$$

layer no

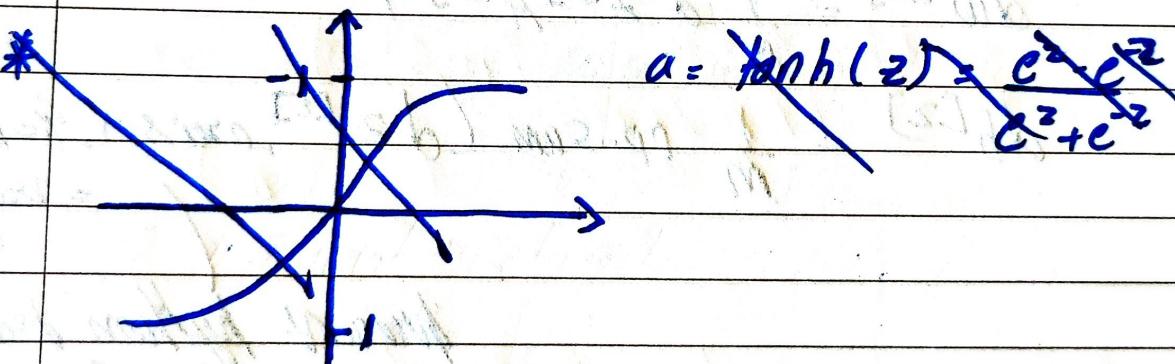
$$A^{(1)} = \begin{bmatrix} a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(m)} \end{bmatrix}$$

Training examples →



Mostly used RELU

- Sigmoid f^n isn't a good choice in maximum
- (and n 's)
- RELU is used in most



* Gradient descent does neural networks

Parameters: $w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}$
 $(n^{[0]}, n^{[0]})$ $(n^{[1]}, n^{[1]})$
 $(n^{[2]}, n^{[2]})$ $(n^{[2]}, 1)$

cost function: $J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}_i, y_i)$

Formulas for Computing derivatives

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:-

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

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

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

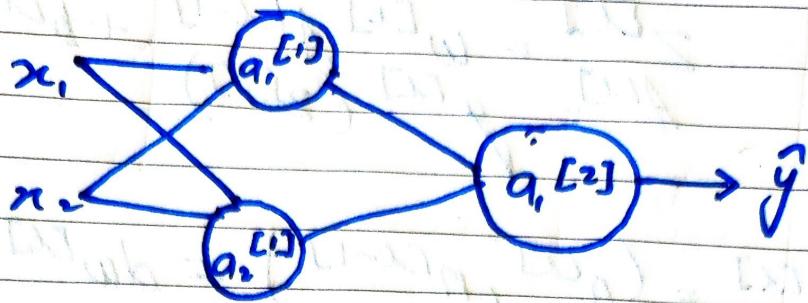
prevent python tensor
Sharing ($n^{[2]}$)

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

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

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

* ~~Topic~~ Random initialization



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

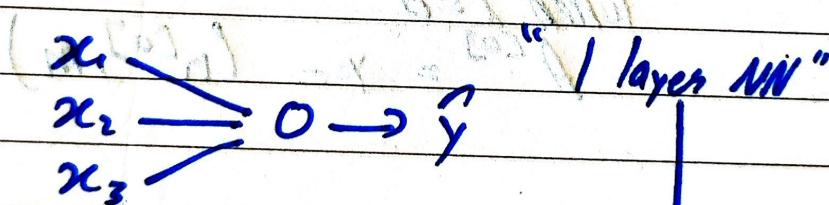
$$b^{[1]} = \text{np.zeros}(2, 1)$$

Similarly,

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

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

* What is a deep neural network?



$$\begin{matrix} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \\ x_3 & 0 & 0 & 0 \end{matrix} \rightarrow \hat{y} = a$$

$$L = 4$$

$n^{[0]}$ no. of units or

nodes in layer 1

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

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

$$n^{[0]} = n_x = 3$$

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

$$w^{[l]} = \text{weight of } z^{[l]}$$

* Forward propagation in a deep network.

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

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

* $w^{[l]} = (n^{[l]}, n^{[l-1]}) = dw^{[l]}$

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

* Vectorized implementation.

$$z^{[l]}, a^{[l]} = (n^{[l]}, 1)$$

$$z^{[l]}, A^{[l]} = (n^{[l]}, m) = dz^{[l]}, dA^{[l]}$$

when $l=0$

$$A^{[0]} = x = (n^{[0]}, m)$$

* Forward and backward functions.
layer l.

