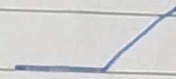
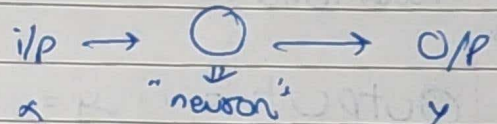


1) Neural Network

ReLU \rightarrow 
Rectified linear unit

Neuron \rightarrow basically takes input & gives out



Supervised \rightarrow Best

2) Supervised Learning

sequence \rightarrow RNN

image \rightarrow CNN

Structured data

Unstructured data

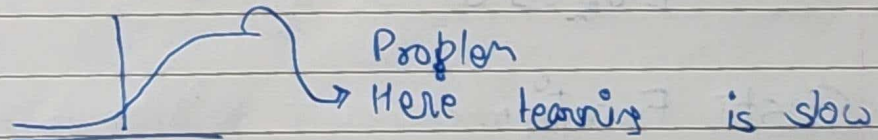
\rightarrow table

\rightarrow audio, img, ~~txt~~ text

easy for comp

easy for humans

3) Sigmoid



ReLU



WEEK 2

4) Binary Classification

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

$n \rightarrow$ no. of training ex.

$m = \{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)}) \}$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)} \end{bmatrix}$$

$$X \in \mathbb{R}^{n \times n}$$

$$X \text{ shape} = (n, n)$$

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

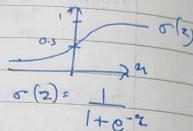
$$Y \in \mathbb{R}^{1 \times n}$$

$$Y \text{ shape} = (1, n)$$

5) Logistic Regression
 $\hat{y} = P(y=1|x)$ (Probability)
 Parameter: $w \in \mathbb{R}^n, b \in \mathbb{R}$

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

$b \rightarrow$ intercept



$$x_0 = 1, \quad x \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \begin{matrix} b \\ w \end{matrix}$$

5) Cost Function

want, $\hat{y}^{(i)} \approx y^{(i)}$

loss (error func):

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

if $y=1$, $L(\hat{y}, y) = -\log \hat{y}$ \leftarrow want \hat{y} large
 if $y=0$, $L(\hat{y}, y) = -\log(1-\hat{y})$ \leftarrow \hat{y} want min

$$\therefore P(y=1|x) = \hat{y} \theta (1-\hat{y})^{1-\theta}$$

$$\therefore \log p(y|x) = -L(\hat{y}, y)$$

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

isme hai problem hai.

3) Cost funcⁿ: $J(w, b) = \frac{1}{n} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})$

loss funcⁿ only for i^{th} steps

Cost function for n steps

We need to minimize Cost Funcⁿ

7) Gradient descent.

want to find w, b that minimize $J(w, b)$

In here we repeat \downarrow

$$w := w - \alpha \frac{dJ(w)}{dw}$$

α is learning rate
 w is updated b

$\alpha \rightarrow$ learning rate

further we use $dw = \frac{dJ(w)}{dw}$

$$\therefore w := w - \alpha dw$$

$\frac{dJ(w)}{dw} \rightarrow$ represents slope of parameter

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

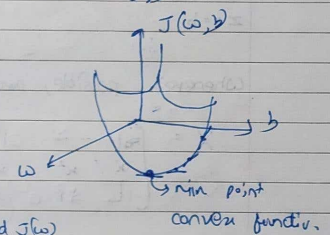
ACTUAL

$$Z = w x_1 + w_2 x_2 + b = w^T x + b$$

$$a = \sigma(z)$$

$$L = L(a, y)$$

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



Vectorisation helps reduction in for loops.

$$z = w^T x + b \quad w = \begin{bmatrix} : \\ \end{bmatrix} \quad x = \begin{bmatrix} i \\ \end{bmatrix}$$

Non vectorized

```
z = 0
for i in range(n):
    z += w[i] * x[i]
z += b
```

vectorized

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b$$

Whenever possible, avoid explicit for loops

$$X = \begin{bmatrix} x^1 & x^2 & x^3 & \dots & x^n \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix} \quad R^{n \times m}$$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(n)}] = w^T X + [b \ b \ b \ \dots \ b]$$

$$Z = \text{np.dot}(w.T, X) + b \quad [\text{Broadcasting in python}]$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(n)}] = \sigma(Z)$$

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

$$dZ = [dz^{(1)} \ dz^{(2)} \ \dots] \quad dA = [d^{(1)} \ \dots \ d^{(n)}] \quad Y = [y^{(1)} \ \dots \ y^{(n)}]$$

$$dZ = A - Y$$

$$dw = \frac{1}{n} X dZ, \quad db = \frac{1}{n} \text{np.sum}(dZ)$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

8) Broadcasting

In vector & nonvector cols

$$\text{col} = A \cdot \text{sum}(\text{axis}=0)$$

↓

0 = sum vertical

1 = " horizontal

$$\text{percentage} = 100 * A / (\text{col_reshape}(1, 2))$$

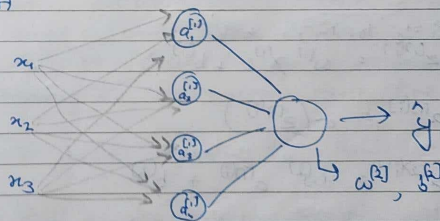
→ Don't use Rank 1 array

WEEK 3

9) NN Representation

$a^{(l)}$ ← layer
 $a_i^{(l)}$ → node in layer

$$X = A^{(0)}$$



$$z_1^{(1)} = w_1^{(1)} x + b_1^{(1)}, \quad d_1^{(1)} = \sigma'(z_1^{(1)})$$

$$z_2^{(1)} = w_2^{(1)} x + b_2^{(1)}, \quad d_2^{(1)} = \sigma'(z_2^{(1)})$$

$$\vdots$$

$$Z^{(1)} = \begin{bmatrix} \omega_{11}^{(1)} & - \\ -\omega_{21}^{(1)} & - \\ -\omega_{31}^{(1)} & - \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ \vdots \end{bmatrix} = \begin{bmatrix} \omega_{11}^{(1)} x_1 + b_1^{(1)} \\ \omega_{21}^{(1)} x_1 + b_2^{(1)} \\ \vdots \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ \vdots \end{bmatrix}$$

$$\begin{aligned} 4 \times 3 & \quad 3 \times 1 & 4 \times 1 \\ Z^{(1)} &= \omega^{(1)} x + b^{(1)} \\ a^{(1)} &= \sigma(Z^{(1)}) \\ Z^{(2)} &= \omega^{(2)} a^{(1)} + b^{(2)} \\ a^{(2)} &= \sigma(Z^{(2)}) \end{aligned}$$

$$x \rightarrow a^{(2)} = \hat{y}$$

$$x^{(1)} \rightarrow a^{(2)(1)} = \hat{y}^{(1)}$$

$$x^2 \rightarrow a^{(2)(2)} = \hat{y}^{(2)}$$

$$\vdots$$

$$x^n \rightarrow a^{(2)(n)} = \hat{y}^{(n)}$$

$a^{(2)(1)}$ example?
→ layer 2

for $i=1$ to n ,
 $z^{(1)(i)} = \omega^{(1)} x^{(i)} + b^{(1)}$

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

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

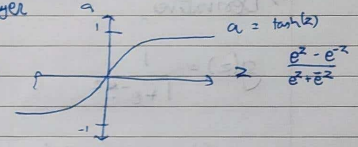
$$a^{(2)(i)} = \sigma(z^{(2)(i)})$$

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

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

horizontal = n , ~~not~~ vertically depends on no. of hidden layers.

Instead of sigmoid $\sigma(z)$, using $g(z)$ i.e. $\tanh(z)$ works better, for hidden layer

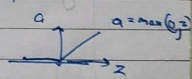


Now instead of avg at 0.5, we get average closer to 0.

For output layer it is better to use sigmoid as $0 \leq \hat{y} \leq 1$ and of -1 to 1 .

However problem with both $g(z)$ is at very large and very small values of z the gradient descent/slope becomes closer to 0, ~~hence~~ so consumes more time.

So one better option is ReLU



defn ReLU is also helpful, ~~not~~ not used much, check if necessary $a = \max(0, z)$

⇒ why do we need non-linear activation funcⁿ.

$g(z) = z \Rightarrow$ linear.
if we use linear activation funcⁿ the NN behaves as just a net if there is no hidden layer, i.e. hidden layers are useless.

if $y \in \mathbb{R}$, you may use linear activation. use it only in output layer if necessary.

→ Derivatives

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

$$g'(z) = \frac{dg(z)}{dz} = g(z)(1 - g(z)) = a(1 - a)$$

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

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

$$= 1 - (\tanh(z))^2$$

$$g(z) = 1 - a^2$$

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

$$g'(z) = \begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

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

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

→ Back Propagation

$$dz^{(k)} = A^{(k)} - y$$

$$dw^{(k)} = \frac{1}{n} dz^{(k)} A^{(k)T}$$

$$db^{(k)} = \frac{1}{n} \text{np.sum}(dz^{(k)}, \text{axis}=1, \text{keepdims}=\text{True})$$

to prevent Rank 1 array

$$dz^{(l)} = \underbrace{w^{(l+1)T}}_{n^{(l+1) \times n^{(l)}}} \underbrace{dz^{(l+1)}}_{\text{means element wise product}} \underbrace{g^{(l+1)'}}_{n^{(l+1)} \times n^{(l)}}(z^{(l+1)})$$

$$dw^{(l)} = \frac{1}{n} dz^{(l)} X^T, \quad db^{(l)} = \frac{1}{n} \text{np.sum}(dz^{(l)}, \text{axis}=1, \text{keepdims}=\text{True})$$

→ We cannot initialize w as 0 for DL.
as then, $a_1^{(1)} = a_2^{(1)}$ also, $dz_1^{(1)} = dz_2^{(1)}$

Hence, ini. it randomly.

$w^{(1)} = \text{np.random.rand}(2,2) * 0.01$
 $b^{(1)} = \text{np.zeros}(2,1)$
↳ small because
if $w \uparrow$ then
 $g(z)$ will have slope
zero.

WEEK-4

Dimensional Analysis.

$n^{(i)}$ = no of ~~nodes~~ hidden units in layer i

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

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

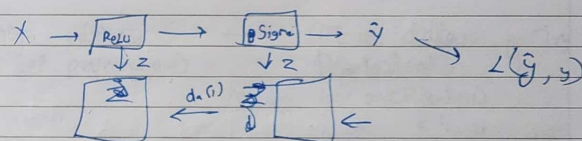
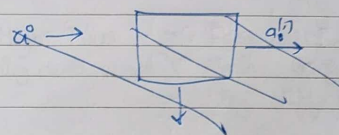
$$(n^{(1)} \times 1) \quad (n^{(1)} \times n^{(1)}) \quad (n^{(1)} \times 1) \quad (n^{(1)} \times 1)$$

dw & db have same dimensions as w & b resp.

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

$$(n^{(1)} \times m) \quad (n^{(1)} \times n^{(0)}) \quad (n^{(0)} \times m) \quad (n^{(1)} \times 1)$$

$$dz^{(1)}, dA^{(1)}, z^{(2)}, A^{(2)} : (n^{(2)} \times m)$$



Parameters : w & b

Hyperparameters : α (learning rate)
iterations

hidden layers l
hidden units $n^{(1)}, n^{(2)}, \dots$
choice of activation funcⁿ

data : momentum, minibatch size ...
we'll see