

10/9/22

## Homework (3)

Page No.:

Date:

Q1. ① We know that,

$$w_n = w_{n-1} - \alpha \left. \frac{dL}{dw} \right|_{w_{n-1}} \rightarrow (1)$$

where  $w_n, w_{n-1}$  are consecutive weights &  $L$  is the training loss. with  $\alpha \rightarrow$  learning rate.

\* A very high learning rate causes large error gradients to accumulate which results in very large updates to the weights during training. This is called the exploding gradient problem.

\* This causes model training to diverge & the loss increases rapidly as a result.

② \* Learning rate scheduler updates the learning rate during the training process after a certain number of epochs, according to some rules.

\* In the figure, the learning rate decays [decreases] after 15 epochs each for some part of the training process.

\* Such a schedule with decreasing learning rate after some epochs during the training process can lead to improved accuracy, faster convergence, and reduce overfitting of the model.

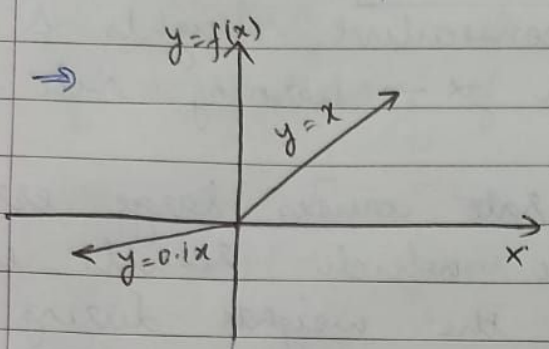
e.g., Reducing <sup>(a very high)</sup> learning rate uniformly during the training will cause model to converge faster towards the beginning, and moving reasonably ~~fast~~ <sup>slow</sup> in later epochs.

\* Also, changing learning rate makes the model training independent of it.



Q2. (1)  $f(x) = \max(x, 0) + \min(x, 0) \times 0.1$

ans.  $f(x) = \begin{cases} x, & x \geq 0 \\ 0.1x, & x < 0 \end{cases} \rightarrow \textcircled{1}$



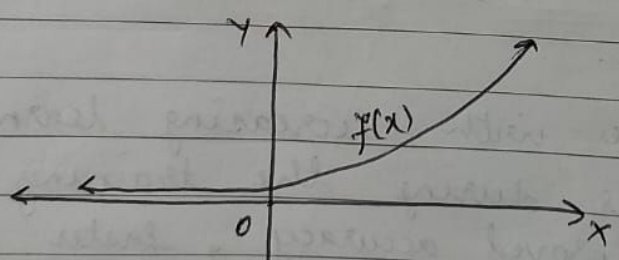
$\therefore$  Since  $f(x)$  is non-linear, it may be a suitable activation  $ft^n$ .

$f'(x) = \begin{cases} 1, & x \geq 0 \\ 0.1, & x < 0 \end{cases} \rightarrow \textcircled{2}$

$\therefore$  Gradient of  $f(x)$  is never going to be 0.

(2)  $f(x) = \ln(e^{3x} + 1)$

ans.  $f'(x) = \frac{e^{3x}}{e^{3x} + 1} \cdot 3 = \frac{3e^x}{e^{3x} + 1} \rightarrow \textcircled{1}$   
 $= \frac{3}{1 + e^{-3x}}$



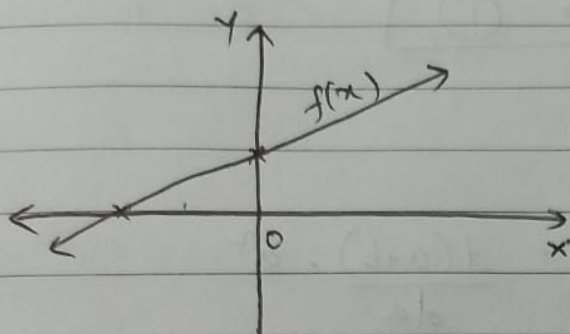
$\therefore f(x)$  is a non-linear activation  $ft^n$ .

Using  $\textcircled{1}$ ,  
 $f'(x) = \begin{cases} 3, & x \rightarrow \infty \\ 0, & x \rightarrow -\infty \end{cases}$

$\Rightarrow f(x)$  is prone to vanishing gradient for high negative values of  $x$ .

③  $f(x) = \ln(e^{3x+1})$

ans:  $f(x) = (3x+1) \rightarrow ①$



$\therefore f(x)$  is a linear activation  $f_t^n$ .

$f'(x) = 3 \rightarrow ②$

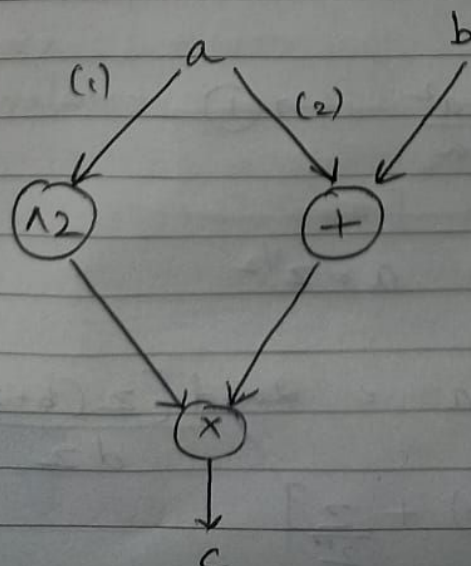
$\Rightarrow f'(x)$  is not prone to gradient vanishing.

$\therefore$  Combining ①, ②, ③ parts,

(i)  $f(x) = \ln(e^{3x+1})$  is not a suitable activation  $f_t^n$  since it's linear & won't be able to fit a variety of data.

(ii)  $f(x) = \ln(e^{3x+1})$  is prone to gradient vanishing.

Q3.



$a = -1, b = 4$



(1.)  $dc/da = ?$

ans:

$$C = (a^2) \times (a+b)$$

$$c = a^2 \cdot (a+b) \quad \text{---} \quad \textcircled{1}$$

$$\frac{dc}{da} = \frac{d[a^2 \cdot (a+b)]}{da}$$

$$= \left( \frac{d(a^2)}{da} \right) \cdot (a+b) + \frac{d(a+b)}{da} \cdot a^2$$

$$= 2a(a+b) + a^2 \quad \text{---} \quad \textcircled{2}$$

★ Putting values in eq. (2),

$$\left( \frac{dc}{da} \right) = 2(-1)[-1+4] + (-1)^2$$

$$= -6 + 1 = \underline{\underline{-5}}$$

(2.)  $\left( \frac{dc}{da} \right)$  component at (1) & (2) ?

ans:  $\left( \frac{dc}{da} \right) = \left[ \frac{\partial c}{\partial v_1} \cdot \frac{\partial v_1}{\partial a} \right] + \left[ \frac{\partial c}{\partial v_2} \cdot \frac{\partial v_2}{\partial a} \right] \quad \text{---} \quad \textcircled{1}$

where  $v_1 = a^2$  &  $v_2 = (a+b)$

&  $c = v_1 + v_2$

( $v_1 \equiv \text{location (1)}$ )  
( $v_2 \equiv \text{location (2)}$ )

At (1):  $\rightarrow$

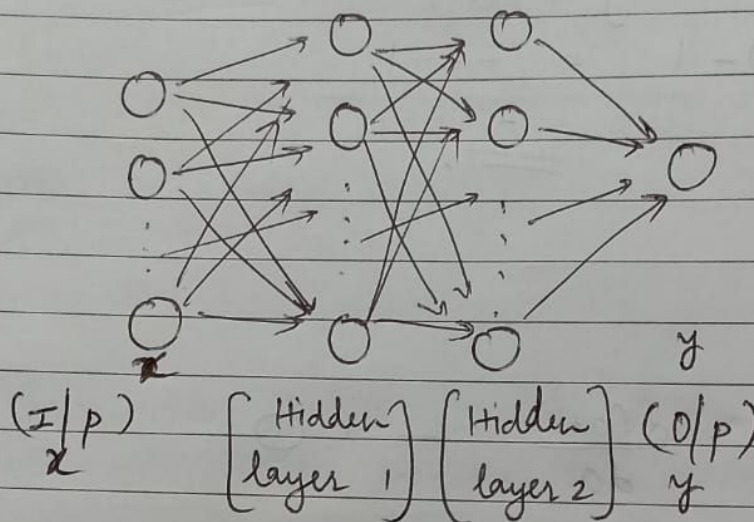
$$\frac{\partial c}{\partial v_1} \cdot \frac{\partial v_1}{\partial a} = (a+b) \cdot (1) \cdot (2a) = (-1+4) \cdot 2(-1)$$

$$= \underline{\underline{-6}}$$

At (2):  $\rightarrow$

$$\frac{\partial c}{\partial v_2} \cdot \frac{\partial v_2}{\partial a} = a^2 \cdot (1) \cdot (1) = (-1)^2 = \underline{\underline{1}}$$

Q4.

Binary Classification Problem

①  $H.L-1 = 100$  units

$H.L-2 = 20$  units

No. of parameters =  $(100 \times 100) + (100 \times 20) + (20 \times 1)$   
 $= 10,000 + 2,000 + 20$   
 $= \underline{\underline{12,020}}$

②  $H.L-1 = 20$  units

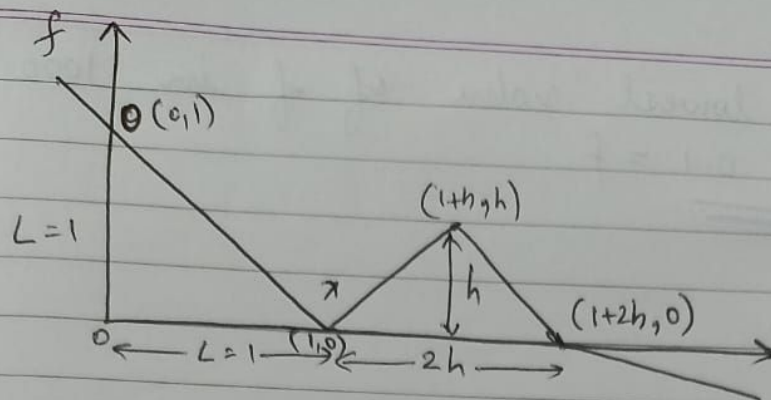
$H.L-2 = 100$  units

No. of parameters =  $(100 \times 20) + (20 \times 100) + (100 \times 1)$   
 $= 2,000 + 2,000 + 100$   
 $= \underline{\underline{4,100}}$

P.T.O.



Q.5.



① Lowest value of  $f$  in 1000 steps?

ans.  $f(x) = \begin{cases} (1-x) & , x \leq 1 \\ (x-1) & , 1 < x \leq (h+1) \\ (1-x) & , (h+1) < x \leq (1+2h) \\ \beta(1-x) & , (1+2h) < x \quad \forall \beta > 0 \end{cases}$

$$f'(x) = \begin{cases} -1 & , x \leq 1 \\ 1 & , 1 < x \leq (h+1) \\ -1 & , (h+1) < x \leq (1+2h) \\ -\beta & , x > (1+2h) \end{cases}$$

\* At point 0,  $x_0 = 0$   
 $\alpha = 0.3$

$$\Rightarrow x_1 = x_0 - \alpha f'(x_0) = 0 - 0.3(-1) = 0.3$$

$$x_2 = 0.3 - 0.3(-1) = 0.6$$

$$x_3 = 0.6 - 0.3(-1) = 0.9$$

$$x_4 = 0.9 - 0.3(-1) = 1.2$$

$$x_5 = 1.2 - 0.3(1) = 0.9$$

$$x_6 = 0.9 - (0.3)(-1) = 1.2$$

$\Rightarrow$  We see that after  $x_3 = 0.9$ ,  $x$  starts to oscillate b/w 0.9 & 1.2.

$\hookrightarrow$  We are oscillating around a local minima.

$$\Rightarrow y_3 = (1-0.9) = 0.1 \quad \& \quad y_4 = (1.2-1) = 0.2$$

Date: \_\_\_\_\_

$\therefore$  The lowest value of  $f$  in 1000 steps  
is  $0.1 = f$

(2)  $\alpha = 0.001$

$\beta_1 = 0.9$

$\beta_2 = 0.999$

$\epsilon = 0$

ans.  $\star$  From the graph, we know that Adam optimizer will escape the local minima at  $x$

if  $\boxed{x > (1+h)}$   $\rightarrow$  (1)

$\Rightarrow$  ~~(x-1)~~  $(x-1) > h$

$\Rightarrow$   $\boxed{h_{\max} = (x_{\max} - 1)}$   $\rightarrow$  (2)

$\star \therefore$  We need to find the max. value of  $x$  such that it ~~is greater than~~ ~~(x-1)~~ doesn't start descending towards local minima again.

NOTE: Please find attached code for finding  $h_{\max}$ .

$\boxed{h_{\max} = 0.002434}$