

Puneet - Mangla  
CS17BTECH11029.

Q.1

$$a) \frac{\partial E}{\partial W} = \frac{\partial E_1}{\partial W} + \frac{\partial E_2}{\partial W} + \frac{\partial E_3}{\partial W}$$

$$\frac{\partial E}{\partial U} = \frac{\partial E_1}{\partial U} + \frac{\partial E_2}{\partial U} + \frac{\partial E_3}{\partial U}$$

$$\frac{\partial E}{\partial V} = \frac{\partial E_1}{\partial V} + \frac{\partial E_2}{\partial V} + \frac{\partial E_3}{\partial V}$$

$$b) \frac{\partial E_2}{\partial W} = \sum_{k=0}^2 \frac{\partial E_2}{\partial O_2} \frac{\partial O_2}{\partial H_2} \left( \prod_{j=k+1}^2 \frac{\partial H_j}{\partial H_k} \right) \frac{\partial H_k}{\partial W}$$

$$\frac{\partial E_2}{\partial U} = \sum_{k=0}^2 \frac{\partial E_2}{\partial O_2} \frac{\partial O_2}{\partial H_2} \left( \prod_{j=k+1}^{j-1} \frac{\partial H_j}{\partial H_{j-1}} \right) \frac{\partial H_k}{\partial U}$$

$$\frac{\partial E_2}{\partial V} = \frac{\partial E_2}{\partial O_2} \frac{\partial O_2}{\partial V} = \frac{\partial E_2}{\partial O_2} \frac{\partial O_2}{\partial Z_2} \frac{\partial Z_2}{\partial V}$$

$$\left\{ Z_2 = V H_2 \right\}$$

$$c) \frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial O_3} \frac{\partial O_3}{\partial H_3} \left( \prod_{j=k+1}^3 \frac{\partial H_j}{\partial H_{j-1}} \right) \frac{\partial H_k}{\partial W}$$

$$\frac{\partial E_3}{\partial U} = \sum_{k=0}^3 \frac{\partial E_3}{\partial O_3} \frac{\partial O_3}{\partial H_3} \left( \prod_{j=k+1}^3 \frac{\partial H_j}{\partial H_{j-1}} \right) \frac{\partial H_k}{\partial U}$$

$$\frac{\partial E_3}{\partial V} = \frac{\partial E_3}{\partial O_3} \frac{\partial O_3}{\partial Z_3} \frac{\partial Z_3}{\partial V} \quad \left\{ Z_3 = V H_3 \right\}$$

Q.3

Rank	0	1	2	3	4	5	6	7	8	9	10
precision		1	1	0.6	0.5	0.4	0.5	0.57	0.5	0.44	0.5
Recall	0.2	0.4	0.4	0.4	0.4	0.4	0.6	0.8	0.8	0.8	1

$\lambda$	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
$P_{interp}(\lambda)$	1	1	1	1	1	0.57	0.57	0.57	0.57	0.5	0.5

$$AP_{interp} = \frac{1}{11} \sum_{\lambda} P_{interp}(\lambda)$$

$$= \frac{1}{11} (8.28) = 0.7527$$

Q.2 a)  $\frac{\partial E_T}{\partial w} = \sum_{k=0}^{3^T} \frac{\partial E_T}{\partial y_T} \cdot \frac{\partial y_T}{\partial h_T} \left( \prod_{j=k+1}^T \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_k}{\partial w}$

$\text{if } \left| \frac{\partial h_j}{\partial h_{j-1}} \right| < 1$  } might be because of use of sigmoid activations

then  $\left| \prod_{j=k+1}^T \frac{\partial h_j}{\partial h_{j-1}} \right|$  will vanish, resulting in

no or very minor gradient of  $E_T$  w.r.t to

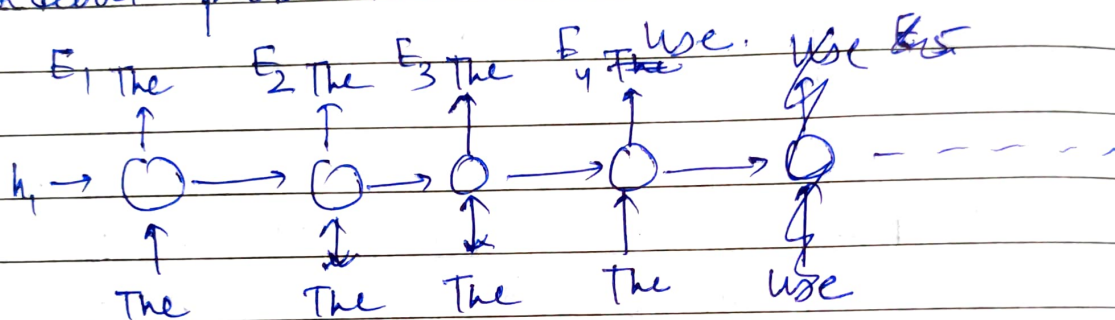
i.e.  $E_T$  will not have any contribution in updating weight  $w$ .

Solution: - Use of gated networks or activations that are not bounded between -1, 1



b)(i) The data above can be confusing for a normal recurrent network as, many same words are repeated many times before a new word comes.

A model can get confused if a 'the' should follow 'the' or 'use' should follow a 'the'. The model can become biased to former. The reason of this is because of vanishing gradient problem.



$E_i$  represents error. now while calculating gradient of total error w.r.t to  $w$ , we get following

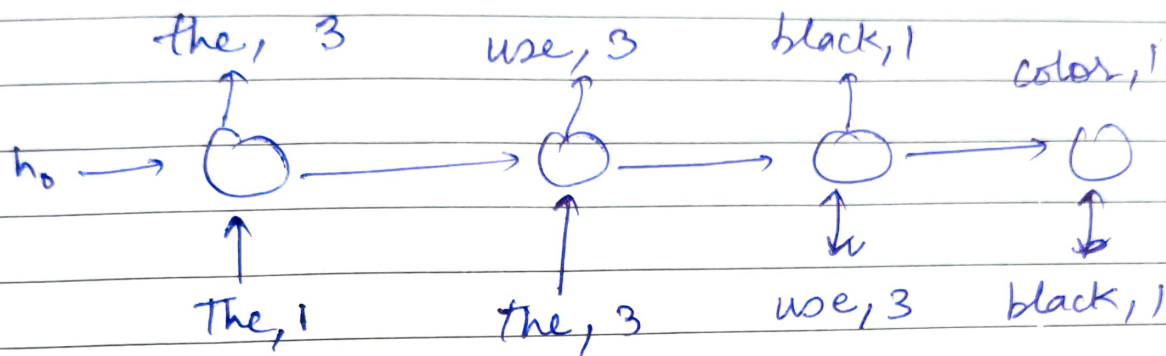
$$E = E_1 + E_2 + E_3 + E_4 + \dots$$

$$\frac{\partial E}{\partial w} = \frac{\partial E_1}{\partial w} + \frac{\partial E_2}{\partial w} + \frac{\partial E_3}{\partial w} + \frac{\partial E_4}{\partial w} + \dots$$

Errors  $E_1, E_2, E_3$  focus on predicting 'the' after seeing 'the', whereas  $E_4$  will focus on predicting 'use' after seeing 'the'.

Hence the contribution of  $E_4$  in updating  $w$  will be 4 times less than  $E_1, E_2, E_3$  which is equivalent to vanishing of gradient.

- (ii) The suggested modification is :  
Along with predicting the final next occurrence, predict the no. of times it will appear consecutively before a new word.



Now, if we use gated networks, the vanishing gradient problem will be resolved.



Q.4  $FL = -(1-p)^{\gamma} \log p$

When  $\gamma=0$   $FL = -\log p$  which is same as standard cross Entropy loss.

For correctly classified examples, Focal loss is relatively less than cross entropy loss thus preventing the model to make over-confident predictions which can negatively affect its overall performance.

For incorrectly classified examples, the <sup>focal</sup> loss is not affected much - just decreased slightly as  $(1-p)^{\gamma}$  is more near 1 and in mis classification  $p < 0.5$

Q.5 A bounding box can be represented as 4 dimensional vector -  $(x, y, h, w)$

$x, y \rightarrow$  coordinates of center

$h, w \rightarrow$  height and width of bounding box.

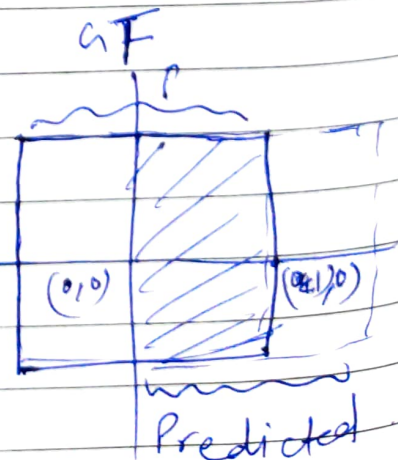
Consider following 2 cases

(1) GT bbox =  $(0, 0, 2, 2)$

Predicted bbox =  $(0, 0, 2, 2)$

$L_2$  norm is  $= \sqrt{1^2 + 0^2 + 0^2 + 0^2}$   
 $= 1$

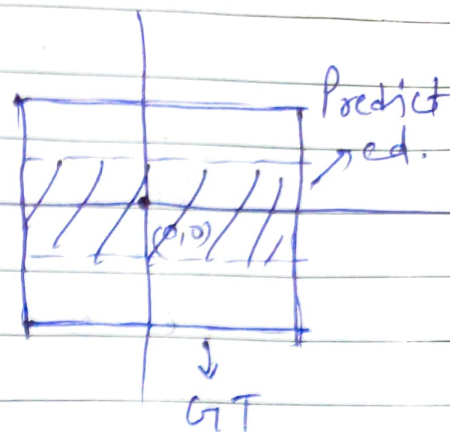
$IOU = \frac{2}{6} = \frac{1}{3}$



(2) GT bbox = (0, 0, 2, 2)  
 Predicted bbox = (0, 0, 1, 2)

$$L_2 \text{ norm} = \sqrt{0^2 + 0^2 + 1^2 + 0^2} = 1$$

$$IOU = \frac{2}{4} = \frac{1}{2}$$



The reason why it happens is because the  $L_2$  norm of vector representation of a bounding box -  $\|(x, y, h, w)\|$  doesn't take into the/capture the relative shift/positioning of bboxes. It considers each of the elements -  $x, y, h, w$  as separate quantity while calculating the norm.

IOU on other hand, takes into account the relative positioning of two bboxes.

Q.6 a)  $H_o = (H_I - 1) * \text{stride} + H_k - 2 * \text{padding}$   
 $W_o = (W_I - 1) * \text{stride} + W_k - 2 * \text{padding}$   
 $H_o, W_o$  - height, width of output  
 $H_I, W_I$  - " " of input  
 $H_k, W_k$  - " " of kernel

by applying above equations

$$H_o = 9$$

$$W_o = 9$$

So size of result is  $9 \times 9$ .



b)

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$\begin{bmatrix} e & f \\ g & h \end{bmatrix}$$

Input  
↓

$$\begin{bmatrix} e & 0 & 0 & 0 \\ f & e & 0 & 0 \\ 0 & f & 0 & 0 \\ g & 0 & e & 0 \\ h & g & f & e \\ 0 & h & 0 & f \\ 0 & 0 & g & 0 \\ 0 & 0 & h & 0 \\ 0 & 0 & 0 & h \end{bmatrix}$$

output  
↘

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

$$\begin{bmatrix} ae \\ af + be \\ fb \\ ag + ec \\ ah + gb + fc + de \\ hb + fd \\ gc \\ hc + dg \\ hd \end{bmatrix}$$

↓

resize into 3x3.  
output.