

# DLV Assignment - 4

Vinta Reethu

ES18BTECH11028



We know that

$$h_t = \tanh [v z_t + w h_{t-1}]$$

$$\hat{y}_t = \text{softmax} [v h_t]$$

- let us assume xent loss is used then we can write the derivative of  $E_t$  wrt  $h_t$  as given below

$$\begin{aligned} \frac{\partial E_t}{\partial h_t} &= \frac{\partial E_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} \cdot \frac{\partial z_t}{\partial h_t} \\ &= (\hat{y}_t - y_t) v \end{aligned}$$

- Calculating derivative of  $h_t$ ,

$$\begin{aligned} \frac{\partial h_t}{\partial h_{t-1}} &= \frac{\partial}{\partial h_{t-1}} [w h_{t-1} + v z_t] \\ &= w (1 - \tanh^2 [v z_t + w h_{t-1}]) \\ &= w (1 - h_t^2) \end{aligned}$$

- Calculating  $\frac{\partial E_1}{\partial w}, \frac{\partial E_1}{\partial v}, \frac{\partial E_1}{\partial I_1}$  :

$$\begin{aligned} \frac{\partial E_1}{\partial w} &= \frac{\partial E_1}{\partial h_1} \cdot \frac{\partial h_1}{\partial w} \\ &= (\hat{y}_1 - y_1) v (1 - h_1^2) h_0 \end{aligned}$$

$$\begin{aligned} \frac{\partial E_1}{\partial v} &= \frac{\partial E_1}{\partial h_1} \cdot \frac{\partial h_1}{\partial v} \\ &= (\hat{y}_1 - y_1) v (1 - h_1^2) I_1 \end{aligned}$$

$$\begin{aligned} \frac{\partial E_1}{\partial I_1} &= \frac{\partial E_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial I_1} \\ &= (\hat{y}_1 - y_1) h_1 \end{aligned}$$

b) Calculating  $\frac{\partial E_2}{\partial \omega}$ ,  $\frac{\partial E_2}{\partial U}$ ,  $\frac{\partial E_2}{\partial V}$ :

$$\begin{aligned}
 \frac{\partial E_2}{\partial \omega} &= \frac{\partial E_2}{\partial h_2} \sum_{k=1}^2 \left( \frac{\partial h_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial \omega} \right) \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot \left( \frac{\partial h_2}{\partial \omega} + \left[ \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial \omega} \right] \right) \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot \left[ (1 - h_2^*) h_1 + [(1 - h_2^*) w (1 - h_1^*) h_0] \right] \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot (1 - h_2^*) \cdot [h_1 + w(1 - h_1^*) h_0] \\
 \frac{\partial E_2}{\partial U} &= \frac{\partial E_2}{\partial h_2} \sum_{k=1}^2 \left[ \frac{\partial h_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial U} \right] \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot \left( \frac{\partial h_2}{\partial U} + \left[ \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial U} \right] \right) \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot \left[ (1 - h_2^*) I_2 + (1 - h_2^*) w (1 - h_1^*) I_1 \right] \\
 &= (\hat{y}_2 - y_2) \cdot v \cdot (1 - h_2^*) \cdot [I_2 + w(1 - h_1^*) I_1]
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial E_2}{\partial V} &= \frac{\partial E_2}{\partial E_2} \cdot \frac{\partial E_2}{\partial V} \\
 &= (\hat{y}_2 - y_2) \cdot \underline{\underline{h_2}}
 \end{aligned}$$

c) Calculating  $\frac{\partial E_3}{\partial \omega}$ ,  $\frac{\partial E_3}{\partial U}$ ,  $\frac{\partial E_3}{\partial V}$ :

$$\begin{aligned}
 \frac{\partial E_3}{\partial \omega} &= \frac{\partial E_3}{\partial h_2} \sum_{k=1}^3 \left[ \frac{\partial h_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial \omega} \right] \\
 &= (\hat{y}_3 - y_3) \cdot v \cdot \left( \frac{\partial h_2}{\partial \omega} + \left[ \frac{\partial h_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial \omega} \right] + \left[ \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial \omega} \right] \right) \\
 &= (\hat{y}_3 - y_3) \cdot v \cdot \left( (1 - h_2^*) h_2 + (1 - h_2^*) w (1 - h_2^*) h_1 + (1 - h_2^*) w (1 - h_2^*) \frac{w}{(1 - h_1^*) h_0} \right) \\
 &= (\hat{y}_3 - y_3) \cdot v \cdot (1 - h_2^*) \cdot [h_2 + w(1 - h_2^*) (h_1 + w(1 - h_1^*) h_0)]
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial E_3}{\partial U} &= \frac{\partial E_3}{\partial h_2} \sum_{k=1}^3 \left[ \frac{\partial h_2}{\partial h_k} \cdot \frac{\partial h_k}{\partial U} \right] \\
 &= (\hat{y}_3 - y_3) \cdot v \cdot \left( \frac{\partial h_2}{\partial U} + \left[ \frac{\partial h_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial U} \right] + \left[ \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial U} \right] \right) \\
 &= (\hat{y}_3 - y_3) \cdot v \cdot \left( (1 - h_2^*) I_2 + (1 - h_2^*) w (1 - h_2^*) I_2 + (1 - h_2^*) w (1 - h_2^*) \frac{w}{(1 - h_1^*) I_1} \right)
 \end{aligned}$$

$$= (\hat{y}_3 - y_3) \cdot v \cdot (1 - h_3^v) \cdot \left( x_3 + w(1 - h_2^v) (x_2 + w(1 - h_1^v) x_1) \right)$$

$$\begin{aligned} \frac{\partial E_3}{\partial v} &= \frac{\partial E_3}{\partial x_3} \cdot \frac{\partial x_3}{\partial v} \\ &= (\hat{y}_3 - y_3) \cdot \underline{h_3} \end{aligned}$$

a)

$$\frac{\partial E}{\partial w} = \sum_{k=1}^3 \frac{\partial E_k}{\partial w}$$

$$\begin{aligned} &= (\hat{y}_1 - y_1) v (1 - h_1^v) h_0 + \\ &\quad (\hat{y}_2 - y_2) v (1 - h_2^v) [h_1 + w(1 - h_1^v) h_0] + \\ &\quad (\hat{y}_3 - y_3) v (1 - h_3^v) [h_2 + w(1 - h_2^v) (h_1 + w(1 - h_1^v) h_0)] \end{aligned}$$

$$\frac{\partial E}{\partial v} = \sum_{k=1}^3 \frac{\partial E_k}{\partial v}$$

$$\begin{aligned} &= (\hat{y}_1 - y_1) v (1 - h_1^v) x_1 + \\ &\quad (\hat{y}_2 - y_2) v (1 - h_2^v) [x_2 + w(1 - h_1^v) x_1] + \\ &\quad (\hat{y}_3 - y_3) v (1 - h_3^v) [x_3 + w(1 - h_2^v) (x_2 + w(1 - h_1^v) x_1)] \end{aligned}$$

$$\frac{\partial E}{\partial v} = \sum_{k=1}^3 \frac{\partial E_k}{\partial v}$$

$$\begin{aligned} &= (\hat{y}_1 - y_1) h_1 + \\ &\quad (\hat{y}_2 - y_2) h_2 + \\ &\quad (\hat{y}_3 - y_3) \underline{h_3} \end{aligned}$$

a) If we look at the gradient term of  $E$  wrt  $W$ ,

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \cdot \frac{\partial \hat{y}_3}{\partial h_3} \left( \prod_{j=k+1}^3 \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_k}{\partial W}$$

For sigmoid activation functions, this gradient is upper bounded by 1. The gradients that carry information become smaller and smaller and will vanish over time. Hence long term dependencies will not contribute for learning.

We can solve this vanishing gradient problem in RNN by

- Using ReLU as activation function.
- Using Regularization.
- Better initialization of weights.
- Using short time sequences.
- Using architectures such as LSTMs and GRUs.

b) i) The above given dataset has repetitive words. If there are words that are repeated in almost entire dataset, it has to capture the long term dependencies in order to predict accurately. So to capture the long term dependencies, the network has to back propagate through time over a long interval. As discussed above this would lead to vanishing gradient problem.

ii) We can use ReLU as an activation function or we can use regularization. We could use LSTM's instead of RNN's as it proven to avoid vanishing gradient problem (as it has all states connected by highway).

3. Let us look at some of the formulas that are needed to compute interpolated AP

• Precision =  $\frac{cTP}{cTP + cFP}$  c - cumulative

• Recall =  $\frac{cTP}{cTP + cFN} = \frac{cTP}{\text{Bounding boxes in total}}$

CTP	CTP	Rank	Precision	Recall
0	0	0	0	0
1	0	1	1	$\frac{1}{5}$
2	0	2	1	$\frac{2}{5}$
2	1	3	$\frac{2}{3}$	$\frac{2}{5}$
2	2	4	$\frac{1}{2}$	$\frac{2}{5}$
2	3	5	$\frac{1}{2}$	$\frac{3}{5}$
3	3	6	$\frac{2}{5}$	$\frac{3}{5}$
4	3	7	$\frac{1}{2}$	$\frac{4}{5}$
4	4	8	$\frac{4}{7}$	$\frac{4}{5}$
4	5	9	$\frac{1}{2}$	$\frac{4}{5}$
5	5	10	$\frac{4}{9}$	$\frac{4}{5}$
			$\frac{1}{2}$	1

- The formula for calculating interpolated AP is

$$AP = \frac{1}{11} \sum_{r \in (0, 0.1, \dots, 1)} P_{\text{interp}}(r)$$

where

$$P_{\text{interp}}(r) = \max_{\bar{r}: \bar{r} \geq r} p(\bar{r})$$

- Therefore  $P_{\text{interp}}(r)$  can be as below

Recal level	$P_{\text{interp}}(r)$
0	1
0.1	1
0.2	1
0.3	1
0.4	1
0.5	$\frac{4}{7}$
0.6	$\frac{5}{7}$
0.7	$\frac{4}{7}$
0.8	$\frac{4}{7}$
0.9	$\frac{1}{2}$
1	$\frac{1}{2}$

$$\text{Interpolated AP} = \frac{1}{11} [1+1+1+1+1+\frac{4}{3}+\frac{4}{3}+\frac{4}{3}+\frac{4}{3}+\frac{1}{2}+\frac{1}{2}]$$

$$= \frac{58}{77} = \underline{\underline{0.753}}$$

- 4.
- Focal loss function is given by  $-(1-p)^{\gamma} \log p$ .
  - The standard cross entropy loss is given by  $-\log p$ .
  - When  $\gamma=0$ , focal loss is equal to standard cross entropy loss.
  - The way focal loss work is, it increases (or penalizes) loss for incorrectly classified data and decreases loss for correctly classified data.

5. We can prove the claim as follows.

- Let  $(x_1, y_1, x_2, y_2)$  be a tuple denoting a bounding box.
- Let's find the distance between one of the two corners say  $P$ . Consider a circle of radius  $r$  around the other corner of the ground truth bounding box.
- For any predicted bounding box with corresponding opposite corner lying on this circle will have the same  $L_2$  norm of  $\sqrt{r^2 + r^2}$ . Despite of having same  $L_2$  norm these boxes can have different IOUs.
- This can happen because  $L_2$  norm is a distance measure, whereas IoU measures the extent of overlap between the bounding boxes.

6. a) Given that,

$$\begin{aligned} \text{input-size} &= 3 \times 3 \\ \text{filter} &= 7 \times 7 \\ \text{stride} &= 1 \\ \text{padding} &= 0 \end{aligned}$$

We know that the output shape of transpose convolution is

$$\text{input-size} \times \text{stride} - \text{stride} + \text{filter} - 2 \times \text{padding}$$

$$\begin{aligned} \text{Therefore, output-size} &= 3 \times 1 - 1 + 7 - 2(0) \\ &= 3 - 1 + 7 \\ &= 9 \end{aligned}$$

∴ output  $\rightarrow 9 \times 9$  matrix

b) Input =  $2 \times 2$   
 kernel =  $2 \times 2$

for example let us consider,

$$\text{Input} = \begin{bmatrix} 25 & 31 \\ 43 & 49 \end{bmatrix}$$

$$\text{kernel} = \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 3 & 2 & 1 & 0 \\ 0 & 3 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 3 & 2 \\ 0 & 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 25 \\ 31 \\ 43 \\ 49 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \cdot 25 \\ 1 \cdot 31 \\ 2 \cdot 25 \\ 3 \cdot 25 + 2 \cdot 31 + 1 \cdot 43 \\ 8 \cdot 31 + 1 \cdot 49 \\ 2 \cdot 43 \\ 3 \cdot 43 + 2 \cdot 49 \\ 3 \cdot 49 \end{bmatrix}$$

↓  
 Transposed  
 weight matrix

$$\begin{bmatrix} 0 \\ 25 \\ 31 \\ 50 \\ 180 \\ 142 \\ 86 \\ 227 \\ 147 \end{bmatrix} \xrightarrow{\text{reshape}} \begin{bmatrix} 0 & 25 & 31 \\ 50 & 180 & 142 \\ 86 & 227 & 147 \end{bmatrix}$$

- In order to achieve the transposed weight matrix, we can use the following code snippet.

```
def matrix_transposed(Filter):
    k, W = torch.zeros(5), torch.zeros((4, 9))
    k[:, 2], k[:, 3:5] = filter[0, :], k[1, :]
    W[0, :5], W[1, 1:6], W[2, 3:8], W[3, 4:] = k, k, k, k
    return W.T
```

- Once you get the above transposed matrix we just have to multiply and reshape it, which can be done as follows

```
def transpose_convolution(x, filter):
    W = matrix_transposed(filter)
    y = torch.mv(W, x.reshape(-1))
    return y.reshape((3, 3))
```