

Theory of Deep Learning  
Homework 4  
Part II

We have input vectors  $[101211]$ ,  $[120111]$ .

1. Output vectors of first attention need will be given by

$$\text{softmax} \left( \frac{Q_i K_i^T}{\sqrt{d_k}} \right) V_i \text{ where } Q_i, K_i, V_i$$

$$\Rightarrow Q \times W_i^Q$$

$$K \times W_i^K$$

$$V \times W_i^V$$

$$\text{So, } Q_i = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 4 & 3 \\ 3 & 0 & 2 \end{bmatrix}_{(2 \times 3)}$$

$$K_i = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \\ 2 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 7 & 8 & 6 \\ 6 & 8 & 5 \end{bmatrix}$$

$$V_i = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 3 & 0 & 2 & 1 \\ 1 & 2 & 3 & 1 \\ 1 & 2 & 1 & 1 \\ 2 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 \\ 3 & 1 & 3 & 1 \end{bmatrix} = \begin{bmatrix} 12 & 8 & 10 \\ 11 & 9 & 13 \end{bmatrix}$$

Now,  $\Omega, K^T$  values for the two vectors are:

$$\cancel{\begin{bmatrix} 4 & 4 & 3 \end{bmatrix}}_{(1 \times 3)} \begin{bmatrix} 7 \\ 8 \\ 6 \end{bmatrix}_{(3 \times 1)} = \cancel{\begin{bmatrix} 78 \end{bmatrix}} \text{ & } \cancel{\begin{bmatrix} 3 & 5 & 2 \end{bmatrix}} \begin{bmatrix} 6 \\ 8 \\ 5 \end{bmatrix}_{(3 \times 1)} = \cancel{\begin{bmatrix} 68 \end{bmatrix}}$$

$$\Omega, K^T = \begin{bmatrix} 4 & 4 & 3 \\ 3 & 5 & 2 \end{bmatrix}_{(2 \times 3)} \begin{bmatrix} 7 & 6 \\ 8 & 8 \\ 6 & 5 \end{bmatrix}_{(3 \times 2)} = \begin{bmatrix} 78 & 71 \\ 20 & 68 \\ 23 & 71 \end{bmatrix}$$

(1) After Dividing each element by  $\sqrt{6k} \Rightarrow \sqrt{3}$  and running softmax we get output vector as multiplied by V

$$\begin{bmatrix} 0.983 & 0.017 \\ 0.947 & 0.053 \end{bmatrix} \begin{bmatrix} 0.359 & 0.987 \\ 0.640 & 0.017 \end{bmatrix}_{(2 \times 2)} \begin{bmatrix} 12, 8, 10 \\ 11, 9, 13 \end{bmatrix}_{(2 \times 3)}$$

$$= \begin{bmatrix} 11.983 & 8.017 & 10.051 \\ 11.947 & 8.053 & 10.159 \end{bmatrix} \approx \begin{bmatrix} 15.165 & 11.755 & 16.421 \\ 7.867 & 5.273 & 6.621 \end{bmatrix}_{(2 \times 3)}$$

Now, for Head 2 :

$$\Omega_2 = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 2 \\ 2 & 1 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 3 & 6 & 4 \\ 5 & 5 & 7 \end{bmatrix}$$

$$K_2 = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 3 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 3 & 3 \\ 2 & 3 & 2 \end{bmatrix}$$

ribution of masses  
would

and the blanks

D.

parts plan table

if you  
them)

$$V_2 = \begin{bmatrix} 1 & 0 & 1 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 3 & 2 & 4 \\ 4 & 2 & 3 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 3 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 9 & 6 & 9 \\ 16 & 9 & 14 \end{bmatrix}$$

Now, similarly to need 1 we have:

$$\alpha_2 k_2^T = \begin{bmatrix} 3 & 6 & 4 \\ 5 & 5 & 7 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 3 & 3 \\ 3 & 2 \end{bmatrix} = \begin{bmatrix} 36 & 32 \\ 46 & 39 \end{bmatrix}$$

(2) After dividing matrix by  $\sqrt{3}/8$  applying softmax then multiplying by values vector matrix we get:

$$\begin{bmatrix} 0.909 & 0.091 \\ 0.983 & 0.017 \end{bmatrix} \begin{bmatrix} 0.003 & 0.017 \\ 0.996 & 0.182 \end{bmatrix} \begin{bmatrix} 9 & 6 & 9 \\ 16 & 9 & 14 \end{bmatrix} = \begin{bmatrix} 0.299 & 0.171 & 0.265 \\ 24.676 & 14.814 & 22.712 \end{bmatrix}$$

(3) Now, to get final output, we concatenate the two outputs from the 2 needs & multiply with W<sub>o</sub> matrix giving us:

$$\begin{bmatrix} 15.165 & 11.755 & 16.421 & 0.299 & 0.171 & 0.265 \\ 27.867 & 5.273 & 6.621 & 24.676 & 14.814 & 22.712 \end{bmatrix}$$

$$\begin{bmatrix} 9.637 & 6.273 & 9.455 \\ 9.119 & 6.051 & 9.095 \end{bmatrix}$$

$$\begin{bmatrix} 9.637 & 6.273 & 9.455 \\ 9.119 \end{bmatrix} \begin{bmatrix} 11.983 & 9.017 & 10.051 & 9.637 & 6.273 & 9.455 \\ 11.947 & 8.053 & 10.159 & 9.119 & 6.051 & 9.095 \end{bmatrix} .1$$

$$\Rightarrow \begin{bmatrix} 53.854 & 43.615 & 29.325 & 37.944 & 57.723 & 20.102 \\ 52.446 & 42.501 & 28.397 & 37.276 & 56.571 & 20.318 \end{bmatrix}$$

To explain the complexity expression given above, we'll go step by step to see if we can arrive at the same expression ourselves.

Firstly, we let  $X$  be the complete input to the self-attention layer where  $X$  would have a dimensionality of  $(N', D)$  with  $N'$  being the number of input embedding vectors &  $D$  the overall dimension of each.

Now, we need to perform a linear transformation to each of the vectors in  $X$  with the weight matrices to compute the  $Q, K, V$  matrices. This would result in  $3 N'D^2$  operations / complexity.

Next, to compute  $Q \times K^T$  & then  $(Q \times K^T) \times V$  we need another  $2 N'^2 D$  operations to be fed.

After which, we have the final matrix multiplication of this result with the weight matrix  $W^O$  which would again be  $N'D^2$  complexity.

Ultimately, we have the overall complexity of the MSA as:

$$\Omega(3N'D^2 + 2N'D^2 + N'D^2) \Rightarrow \Omega(4N'D^2 + 2N'D^2)$$

$x$	$y$	$x$	$y$
1	2	3	4 5 6 7

- d3 The cost of transporting the distribution of messes  
 $P_x$  to  $P_y$  as shown below would  
be
- $$8 \cdot 14 - 11 + 16 - 21 \Rightarrow 71$$

This is the optimal way to transport the blocks  
of messes here this is the TMD.

We can also confirm this using the transport plan table

		$y(TB)$	
		$y$	
		$y(TD)$	
$x$ (From)	$x$		
	$y$		
$x$ (From)	$y$	4 6	
1		1 0	$\Rightarrow 71$ . (If you
2		0 1	move them)