

Prof. Ryan Cotterell

## Danilo Dordevic: Assignment 6

ddordevic@student.ethz.ch, 21-954-888.

15/01/2023 - 22:19h

## Question 1

## a.1)

Each head  $h \in [N_h]$  in the multi-head attention unit is defined by the matrices  $\mathbf{W}_q^{(h)} \in \mathbb{R}^{D_i \times D_{hid}}$ ,  $\mathbf{W}_k^{(h)} \in \mathbb{R}^{D_i \times D_{hid}}$ ,  $\mathbf{W}_v^{(h)} \in \mathbb{R}^{D_i \times D_o}$ . These matrices are learnable parameters, so each head contains  $D_i \cdot D_{hid} + D_i \cdot D_{hid} + D_i \cdot D_o = \boxed{D_i \cdot (2D_{hid} + D_o)}$  parameters.

Since the multi-head self-attention module contains  $N_h$  independent heads, the total number of parameters of all heads is  $\boxed{N_h \cdot D_i \cdot (2D_{hid} + D_o)}$ .

Additionally, the multi-head self-attention module contains an output matrix  $\mathbf{W}_{out} \in \mathbb{R}^{N_h D_o \times D_{out}}$  with  $\boxed{N_h \cdot D_o \cdot D_{out}}$  parameters, as well as the bias vector  $b_{out} \in \mathbb{R}^{D_{out}}$  with  $\boxed{D_{out}}$  parameters.

Adding the values in dashed boxed, we get the total number of parameters in the multi-head self-attention module:  $\boxed{N_h \cdot D_i (2D_{hid} + D_o) + N_h \cdot D_o \cdot D_{out} + D_{out}}$ .

## a.2

Let images  $\mathbf{Z}_{W \times H \times D_i}$  be indexed by a single vector value  $\mathbf{p} = (i, j)$ , which can be interpreted as a linear index, with the following notation:  $\mathbf{Z}_{\mathbf{p},:} = \mathbf{Z}_{i,j,:}$ .

Self-attention for index  $\mathbf{q}$  can be written as the following:

$$\begin{aligned}
 \mathbf{O}_{\mathbf{q},:} &= \text{Self-Attention}(\mathbf{Z})_{\mathbf{q},:} \\
 &= [\text{softmax}(\mathbf{Z}\mathbf{W}_q\mathbf{W}_k^T\mathbf{Z}^T) \cdot \mathbf{Z}\mathbf{W}_v]_{\mathbf{q},:} \\
 &= \text{softmax}(\mathbf{Z}\mathbf{W}_q\mathbf{W}_k^T\mathbf{Z}^T)_{\mathbf{q},:} \cdot \mathbf{Z}\mathbf{W}_v \\
 &= \boxed{\text{softmax}(\mathbf{Z}_{\mathbf{q},:}\mathbf{W}_q\mathbf{W}_k^T\mathbf{Z}^T) \cdot \mathbf{Z}\mathbf{W}_v} \\
 &= \left( \sum_{\mathbf{k}} \text{softmax}(\mathbf{Z}_{\mathbf{q},:}\mathbf{W}_q\mathbf{W}_k^T\mathbf{Z}^T)_{\mathbf{k}} \mathbf{Z}_{\mathbf{k},:} \right) \mathbf{W}_v
 \end{aligned}$$

Same can be applied to the case of MultiHead-SelfAttention:

$$\begin{aligned}
\text{MultiHeadSelf-Attention}(\mathbf{Z})_{q,:} &= \text{concat}_{h \in [N_h]} [O_{q,:}] \cdot \mathbf{W}_{out} + \mathbf{b}_{out} \\
&= \boxed{\text{concat}_{h \in [N_h]} [\text{softmax}(\mathbf{Z}_{q,:} \mathbf{W}_q \mathbf{W}_k^T \mathbf{Z}^T) \cdot \mathbf{Z} \mathbf{W}_v] \cdot \mathbf{W}_{out} + \mathbf{b}_{out}} \\
&= \text{concat}_{h \in [N_h]} \left[ \left( \sum_k \text{softmax}(\mathbf{Z}_{q,:} \mathbf{W}_q \mathbf{W}_k^T \mathbf{Z}^T)_k \mathbf{Z}_{k,:} \right) \mathbf{W}_v \right] \cdot \mathbf{W}_{out} + \mathbf{b}_{out}
\end{aligned}$$

## b.1)

Plugging  $\mathbf{W}_q = \mathbf{W}_k = \mathbf{0}$  into the expression for  $\mathbf{A}_{q,k}^{relative}$  removes the first three terms, as they all contain at least one of the matrices set to  $\mathbf{0}$ . This leaves the new expression:

$$\begin{aligned}
\mathbf{A}_{q,k}^{relative} &= \mathbf{v}^T \widetilde{\mathbf{W}_k} \mathbf{r}_\delta \\
&= \mathbf{v}^T \mathbf{I} \mathbf{r}_\delta \\
&= \mathbf{v}^T \mathbf{r}_\delta \\
&= -\alpha^{(h)} \begin{pmatrix} 1 \\ -2\Delta_1^{(h)} \\ -2\Delta_1^{(h)} \end{pmatrix}^T \begin{pmatrix} \|\delta\|^2 \\ \delta_1 \\ \delta_2 \end{pmatrix} \\
&= \boxed{-\alpha^{(h)} (\|\delta\|^2 - 2\delta_1 \Delta_1^{(h)} - 2\delta_2 \Delta_2^{(h)})}
\end{aligned}$$

## b.2

### Absolute encoding

Looking at equation 5), we can see that the expression for  $\mathbf{A}_{q,k}^{absolute}$  consists of 4 terms, each of which is produced by the same order of operations over vectors and matrices of the same dimensions, so analyzing the computational complexity of only one of them is enough to determine the computational complexity of  $\mathbf{A}_{q,k}^{absolute}$ . Without the loss of generality, we analyze the first term:  $\mathbf{Z}_{q,:} \mathbf{W}_q \mathbf{W}_k^T \mathbf{Z}_{k,:}^T$ . First, multiplying  $\mathbf{Z}_{q,:} \in \mathbb{R}^{1 \times D_x}$  and  $\mathbf{W}_q \in \mathbb{R}^{D_x \times D_x}$  requires  $D_x \cdot (2D_x - 1)$  operations (there are  $D_x$  entries in the result of this vector-matrix multiplication, and each one requires  $D_x$  multiplications and  $D_x - 1$  additions). This is repeated for  $\mathbf{W}_k^T \mathbf{Z}_{k,:}^T$ , which also requires  $D_x \cdot (2D_x - 1)$  operations. Result of these two multiplications are two  $D_x$ -dimensional vectors, which are then multiplied. This requires  $D_x$  operations. The total cost of computing one term is thus  $2 \cdot D_x(2D_x - 1) + D_x - 1$ . As there are 4 terms, the cost of computing  $\mathbf{A}_{q,k}^{absolute}$  is  $4 \cdot (2 \cdot D_x(2D_x - 1) + D_x - 1) = O(D_x^3)$ . The attention matrix  $\mathbf{A}^{absolute}$  contains  $N^2$  entries, so the cost of computing it is  $O(N^2 D_x^3)$ . To get the self-attention matrix  $O^{absolute}$ , we need to compute the softmax of the matrix  $\mathbf{A}^{absolute}$ , which takes  $O(N^2)$  time. After that, the resulting  $N \times N$  matrix is multiplied by  $\mathbf{Z}$  which is  $N \times D_x$ , requiring  $O(N^2 D_x)$  operations. Finally the result is

multiplied by a  $D_x \times D_x$  matrix  $\mathbf{W}_v$  to yield  $O^{absolute}$ , which takes  $\boxed{O(ND_x^2)}$ . Summing up all of the expressions in dashed boxes, the final computational complexity of computing self-attention in the case of absolute positional encoding is  $\boxed{O(N^2D_x^2 + N^2D_x + ND_x^2)}$

## Relative encoding

It is evident from Equation 6, which defines the expression for the entries to the  $\mathbf{A}^{relative}$  matrix, that the computational costs of computing  $\mathbf{A}^{relative}$  is the same as computing  $\mathbf{A}^{absolute}$ , as there is the same number of terms with the same operations over matrices and vectors of same dimensions. However, if the relative encoding is gaussian, then the first three terms from equation 6 can be ignored, and computing  $\mathbf{A}_{q,k}$  boils down to vector inner product,  $\mathbf{A}_{q,k} = \mathbf{v}^T \mathbf{r}_\delta$ , which has complexity  $O(D_p)$ . Computing the whole  $\mathbf{A}^{relative}$  matrix requires doing the same computation  $N^2$  times, so its computational complexity is  $O(N^2D_p)$ . Now that the  $\mathbf{A}^{relative}$  is found, to get  $\mathbf{O}^{relative} = \text{softmax}(\mathbf{A}^{relative})\mathbf{Z}\mathbf{W}_v$ , we need to perform the softmax operation, and two more matrix multiplications, that are exactly the same as in the absolute encoding part, which take  $O(N^2)$ ,  $O(N^2D_x)$  and  $O(ND_x^2)$ , respectively. the total computation complexity is then  $O(N^2D_p + N^2D_x + ND_x^2)$ , which, if we assume that  $D_p \ll D_x$ , boils down to  $\boxed{O(N^2D_x + ND_x^2)}$ .

### c.1)

The self-attention matrix can be separated "horizontally" into blocks that correspond to each head of the attention unit. Likewise, the  $\mathbf{W}_{out}$  matrix can be "vertically" separated into the corresponding blocks of the appropriate height, such that block  $\mathbf{W}_{out}[(h-1)D_o+1:hD_o,:]$  is multiplied by head  $\mathbf{O}^{(h)}$ , and the results of this block multiplication are summed.

$$\begin{aligned}
\text{MultiHead-SelfAttention}(\mathbf{Z}) &= \text{concat}_{h \in [N_h]} [\text{SelfAttention}_h(\mathbf{Z})] \mathbf{W}_{out} + \mathbf{b}_{out} \\
&= \text{concat}_{h \in [N_h]} [\text{softmax}(\mathbf{A}^{(h)}) \mathbf{Z} \mathbf{W}_v^{(h)}] \mathbf{W}_{out} + \mathbf{b}_{out} \\
&= \text{concat}_{h \in [N_h]} [\mathbf{O}^{(h)}] \mathbf{W}_{out} + \mathbf{b}_{out} \\
&= \sum_{h \in [N_h]} \text{softmax}(\mathbf{A}^{(h)}) \mathbf{Z} \underbrace{\mathbf{W}_v^{(h)} \mathbf{W}_{out}[(h-1)D_o+1:hD_o,:]}_{\mathbf{W}^{(h)}} + \mathbf{b}_{out} \\
&= \sum_{h \in [N_h]} \text{softmax}(\mathbf{A}^{(h)}) \mathbf{Z} \mathbf{W}^{(h)} + \mathbf{b}_{out} \\
\mathbf{W}^{(h)} &= \boxed{\mathbf{W}_v^{(h)} \mathbf{W}_{out}[(h-1)D_o+1:hD_o,:]}
\end{aligned}$$

Selecting  $q$ -th entries requires just selecting the corresponding entries in the softmax at-

tention:

$$\text{MultiHead-SelfAttention}(\mathbf{Z})_{\mathbf{q},:} = \left( \text{concat}_{h \in [N_h]} [\mathbf{O}^{(h)}] \right)_{\mathbf{q},:} \mathbf{W}_{out} + \mathbf{b}_{out} \quad (1)$$

$$= (\text{concat}_{h \in [N_h]} [\mathbf{O}_{\mathbf{q},:}^{(h)}] \mathbf{W}_{out} + \mathbf{b}_{out} \quad (2)$$

$$= \sum_{h \in [N_h]} \text{softmax}(\mathbf{A}^{(h)})_{\mathbf{q},:} \mathbf{Z} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (3)$$

$$= \sum_{h \in [N_h]} \text{softmax}(\mathbf{A}_{\mathbf{q},:}^{(h)}) \mathbf{Z} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (4)$$

$$= \boxed{\sum_{h \in [N_h]} \left( \sum_{\mathbf{k}} \text{softmax}(\mathbf{A}_{\mathbf{q},:}^{(h)})_{\mathbf{k}} \mathbf{Z}_{\mathbf{k},:} \right) \mathbf{W}^{(h)} + \mathbf{b}_{out}} \quad (5)$$

c.2)

We can rewrite the defining equation for convolution, so that it uses the previously introduced vector indexing with  $\mathbf{q}$  instead of scalars  $(i, j)$ . Additionally, the pixel shift  $(\delta_1, \delta_2)$  is rewritten using the vector notation  $\boldsymbol{\delta} = (\delta_1, \delta_2)$ :

$$\begin{aligned} \text{Convolution}(\mathbf{Z})_{i,j,:} &= \sum_{(\delta_1, \delta_2) \in \Delta_K} \mathbf{Z}_{i+\delta_1, j+\delta_2,:} \mathbf{W}_{\delta_1, \delta_2, :, :}^{\text{conv}} + \mathbf{b}_{out} \\ \text{Convolution}(\mathbf{Z})_{\mathbf{q},:} &= \boxed{\sum_{\Delta \in \Delta_K} \mathbf{Z}_{\mathbf{q}+\Delta, :} \mathbf{W}_{\Delta, :, :}^{\text{conv}} + \mathbf{b}_{out}} \end{aligned}$$

To prove that the multi-head self-attention for a chosen pixel  $\mathbf{q}$  is equal to the convolution, we start from equation 5, and utilizing the softmax assumption from Theorem 1, we get:

$$\text{MultiHead-SelfAttention}(\mathbf{Z})_{\mathbf{q},:} = \sum_{h \in [N_h]} \underbrace{\left( \sum_{\mathbf{k}} \text{softmax}(\mathbf{A}_{\mathbf{q},:}^{(h)})_{\mathbf{k}} \mathbf{Z}_{\mathbf{k},:} \right)}_{=1, \text{ for } \mathbf{k} \text{ such that } f(h)=\mathbf{q}-\mathbf{k}} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (6)$$

$$= \sum_{h \in [N_h]} \underbrace{\text{softmax}(\mathbf{A}_{\mathbf{q},:}^{(h)})_{\mathbf{q}-f(h)}}_1 \mathbf{Z}_{\mathbf{q}-f(h),:} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (7)$$

$$= \sum_{h \in [N_h]} \mathbf{Z}_{\mathbf{q}-f(h),:} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (8)$$

$$= \sum_{h \in [N_h]} \mathbf{Z}_{\mathbf{q}+f(h),:} \mathbf{W}^{(h)} + \mathbf{b}_{out} \quad (9)$$

$$= \sum_{\Delta \in \Delta_K} \mathbf{Z}_{\mathbf{q}+\Delta, :} \mathbf{W}^{(f^{-1}(\Delta))} + \mathbf{b}_{out} \quad (10)$$

$$= \sum_{\Delta \in \Delta_K} \mathbf{Z}_{\mathbf{q}+\Delta, :} \underbrace{\mathbf{W}_{f(h), :, :}^{\text{conv}}}_{\Delta} + \mathbf{b}_{out} \quad (11)$$

$$= \boxed{\sum_{\Delta \in \Delta_K} \mathbf{Z}_{\mathbf{q}+\Delta, :} \mathbf{W}_{\Delta, :, :}^{\text{conv}} + \mathbf{b}_{out}} \quad (12)$$

The transition from equation 9 to equation 10 is possible thanks to the symmetry of the shifts in the set  $\Delta_K$ . The equality between equation 10 and equation 11 hold because of the bijective map  $f : [N_h] \rightarrow \Delta_K$ , where we utilized the function to switch the sum from going over the head indices  $h$  to going over the shifts  $\Delta = f(h)$ . Moreover, here it can be seen that a mapping can be made between attention head weights and the convolutional filters:  $\mathbf{W}^{(h)} = \mathbf{W}^{f^{-1}(h)} = \mathbf{W}_{f(h),:, :}^{\text{conv}}$ .

### d.1)

Following are the expressions that I will use later in the proof.

$$\begin{aligned}\mathbf{A}_{\mathbf{q}, \mathbf{k}} &= -\alpha (\|\delta - \Delta\|^2 + c), \delta = \mathbf{q} - \mathbf{k} \\ \mathbf{A}_{\mathbf{q}, \mathbf{j}} &= -\alpha (\|\mathbf{q} - \mathbf{j} - \Delta\|^2 + c), \mathbf{j} \neq \mathbf{k}\end{aligned}$$

The following is the proof that when  $\alpha$  is large, the attention of the model for a given query pixel will focus only on one pixel:

$$\begin{aligned}\lim_{\alpha \rightarrow \infty} \text{softmax}(\mathbf{A}_{\mathbf{q}, :})_{\mathbf{k}} &= \lim_{\alpha \rightarrow \infty} \frac{\exp \mathbf{A}_{\mathbf{q}, \mathbf{k}}}{\exp(\mathbf{A}_{\mathbf{q}, \mathbf{k}}) + \sum_{\mathbf{j} \neq \mathbf{k}} \exp(\mathbf{A}_{\mathbf{q}, \mathbf{j}})} \\ &= \lim_{\alpha \rightarrow \infty} \frac{\exp(-\alpha c)}{\exp(-\alpha c) + \sum_{\mathbf{j} \neq \mathbf{k}} \exp(\mathbf{A}_{\mathbf{q}, \mathbf{j}})} \\ &= \lim_{\alpha \rightarrow \infty} \frac{\exp(-\alpha c)}{\exp(-\alpha c) + \sum_{\mathbf{j} \neq \mathbf{k}} \exp(-\alpha \|\mathbf{k} - \mathbf{j}\|^2 - \alpha c)} \\ &= \lim_{\alpha \rightarrow \infty} \frac{1}{1 + \sum_{\mathbf{j} \neq \mathbf{k}} \exp(-\alpha \|\mathbf{k} - \mathbf{j}\|^2)} \\ &= \frac{1}{1 + \sum_{\mathbf{j} \neq \mathbf{k}} \lim_{\alpha \rightarrow \infty} \exp(-\alpha \|\mathbf{k} - \mathbf{j}\|^2)} \\ &= \frac{1}{1 + 0} \\ &= \boxed{1}\end{aligned}$$

Since it holds that  $\sum_{\mathbf{k}'} \text{softmax}(\mathbf{A}_{\mathbf{q}, :})_{\mathbf{k}'} = 1$ , all other values  $\lim_{\alpha \rightarrow \infty} \text{softmax}(\mathbf{A}_{\mathbf{q}, :})_{\mathbf{j}}, \mathbf{j} \neq \mathbf{k}$  necessarily must be 0.

### d.2)

This is proven by just expanding the expression for  $\mathbf{A}_{q,k}^{relative}$  and applying the property of the 2-norm.

$$\begin{aligned}
\mathbf{A}_{q,k}^{relative} &= \mathbf{v}^T \mathbf{r}_\delta \\
&= -\alpha^{(h)} \begin{pmatrix} 1 \\ -2\Delta_1^{(h)} \\ -2\Delta_2^{(h)} \end{pmatrix}^T \begin{pmatrix} \|\delta\|^2 \\ \delta_1 \\ \delta_2 \end{pmatrix} \\
&= -\alpha^{(h)} (\|\delta\|^2 - 2\delta_1 \Delta_1^{(h)} - 2\delta_2 \Delta_2^{(h)}) \\
&= -\alpha^{(h)} (\|\delta\|^2 - 2\delta^T \Delta) \\
&= -\alpha^{(h)} (\|\delta\|^2 + \|\Delta\|^2 - 2\delta^T \Delta - \|\Delta\|^2) \\
&= -\alpha^{(h)} (\|\delta - \Delta\|^2 - \|\Delta\|^2) \\
&= -\alpha^{(h)} (\|\delta - \Delta\|^2 + c), \quad \boxed{c = -\|\Delta\|^2}
\end{aligned}$$

### e.1)

In order to keep the dimensions of the resulting images the same, as well as the pixel values equal, the image needs to be padded with  $\lfloor K/2 \rfloor$  zeros on each side. Here it is assumed that  $K$  is the kernel size of the convolutional filter.

### e.2)

Yes, a multi-head self-attention layer can express an arbitrarily-dilated convolution, because each head is not limited in terms of the pixels (values at pixel shifts) it attends to. So there exist configurations in which attention heads attend to pixels that form a pattern equivalent to dilated convolution.

### e.3)

Yes, a multi-head self-attention layer can learn to simulate a strided convolution, provided that a pooling layer is applied to the result of the multi-head self-attention layer's output, in order to reduce dimensions to the dimensions that correspond to the actual strided convolution output dimensions and to ignore values that would be skipped in the strided convolution due to the stride. The reasoning is similar to the one in section , as heads can attend to values at different pixel shifts.

## Question 2

At the beginning of training, all of the heads are mostly locally focused, not spread out and positioned randomly. As time goes on, some of the heads from the lower layers tend to spread their attention to a larger area, while keeping the centers near the query pixel.

On the higher layers, it happens that some heads form a circle of attention that is more narrow, and their centers are spread out farther from the query pixel. However, in all layers, there is a present tendency of the attention heads to cover a larger area with attention than at initialization time. As the behavior to support the Main Theorem, I would expect the centers of attention of all heads to form a shape of a kernel that a convolution operation is performed with, by the end of training.

<https://colab.research.google.com/drive/1JkRZJ31R5kLTw0VOWsc9SzKlw-B1RBq0?usp=sharing>