

# The scaling factor in self-attention

William Chuang

October 6, 2023

## 1 Introduction

In the context of deep learning neural networks, attention algorithms are designed for determining connections by weights between two elements in a sequential input. Self-attention was proposed by Vaswani et. al.[2] in 2017. A recent study[1] shows that the paper written by Vaswani et. al.[2] was the fourth most cited paper in artificial intelligence upto 2020. In 2023, [2] became the second most cited paper in that original list. Self-attention mechanism has been implemented into almost all most recent LLM models.

Suppose a sequential input  $X \in \mathbb{R}^N \times \mathbb{R}^{d_x}$  is given. Denote  $Q^t := (XW_q)^t \in \mathbb{R}^{d_q} \times \mathbb{R}^N$ ,  $K^t := (XW_k)^t \in \mathbb{R}^{d_k} \times \mathbb{R}^N$ , and  $V^t := (XW_v)^t \in \mathbb{R}^{d_v} \times \mathbb{R}^N$ . The notations for denoting column vectors of each these three matrices are:  $q_i^t \in \mathbb{R}^{d_q}$ ,  $k_i^t \in \mathbb{R}^{d_k}$ , and  $v_i^t \in \mathbb{R}^{d_v}$ . The  $j$ -th component of each column vector is denoted as  $q_{ij}^t \in \mathbb{R}$ ,  $k_{ij}^t \in \mathbb{R}$ , and  $v_{ij}^t \in \mathbb{R}$ .

Now, in a standard implementation of a transformer-based model or any models that have self-attention mechanism implemented in it,  $W_k, W_v$ , and  $W_q$  are initiated in a way to satisfy the following four conditions:

- (i)  $W_q \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_q}$ ,  $W_k \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_k}$ , and  $W_v \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_v}$ ,
- (ii) the variance of the inner product of each column vector in  $X^t$  with each column vector in either  $W_q, W_k$ , or  $W_v$  to be unit (all these inner products results in entries matrices:  $Q^t$ ,  $K^t$ , and  $V^t$ ),

- (iii) each  $q_i^t = (q_{i1}, \dots, q_{id_q})^t$  and each  $k_i^t = (k_{i1}, \dots, k_{id_k})^t$  to be random variables with zero means,  
i.e.  $\mathbb{E}(q_i^t) = \mathbb{E}(\{q_{ij}\}_{j=1}^{d_q}) = 0 = \mathbb{E}(\{k_{ij}\}_{j=1}^{d_k}) = \mathbb{E}(k_i^t), \forall i \in \{1, \dots, N\}$ , and unit variances,  
i.e.  $\text{Var}(q_i^t) = \text{Var}(\{q_{ij}\}_{j=1}^{d_q}) = 1 = \text{Var}(\{k_{ij}\}_{j=1}^{d_k}) = \text{Var}(k_i^t), \forall i \in \{1, \dots, N\}$ , and
- (iv)  $q_{i_1}^t$  and  $k_{i_2}^t$  are all independent in the following ways:  $\text{Cov}(q_{i_1}^t, k_{i_2}^t) = 0, \forall i_1, i_2$ ,  $\text{Cov}(q_{i_1}^t, q_{i_2}^t) = 0, \forall i_1 \neq i_2$ ,  $\text{Cov}(k_{i_1}^t, k_{i_2}^t) = 0, \forall i_1 \neq i_2$ .

Furthermore, by convention, let  $d_k = d_q$ . Then, with the four above conditions satisfied, any scalar products of column vectors in  $Q^t$  and  $K$  have variance  $d_k$  (see Appendix).

**Definition 1.** *The Softmax function is defined using the Boltzmann distribution function: Given a sequential inputs (i.e. consider the finite sequence as a vector):  $\{z_i\}, z_j \in \mathbb{R}, \forall j \in \{1, \dots, n\}$ , then*

$$\text{Softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}.$$

Furthermore, this distribution can be scaled by using a parameter  $\beta$ , and  $\beta$  is the inverse of the temperature of the physical system described by the distribution:

$$\text{Softmax}(\beta z_i) := \frac{\exp(\beta z_i)}{\sum_j \exp(\beta z_j)}.$$

In order to take inputs that are tensors, for each component of a tensor  $a_{i_1 i_2 \dots i_n}$ , the softmax is operated on its last dimension  $i_n$ , and also output a tensor which has indices  $i_1 \dots i_n$ . For example, for a  $N \times d$  tensor (i.e. a matrix), the softmax is taking on its last dimension, i.e. which equals  $d$ . The output is also an  $N \times d$  matrix, say it is denoted by  $\text{softmax}(a_{i_1 i_2 \dots i_n})_{N \times d}$ . Then, each entry  $\text{softmax}(a_{i_1 i_2})_{ij}$  of the matrix  $\text{softmax}(a_{i_1 i_2})_{N \times d}$  is defined as

$$\text{softmax}(a_{i_1 i_2})_{ij} := \frac{\exp(\beta a_{ij})}{\sum_{i_1=1}^N \exp(\beta a_{i_1 j})}.$$

Operating on the last dimension means to fix the last index in the normalization. In this example that means to consider each column vector in  $a_{i_1 i_2}$  as a set of data or a physical system that Boltzmann distribution was applied to, and the result  $\text{softmax}(a_{i_1 i_2})_{ij}$  is the probability of the

configuration  $a_{ij}$ . The denominator is also called the normalization factor or the partition function that it sums all possible configurations using to given data  $a_{i_1 i_2}$  where  $i_2 = j$  to be the input of the exponential function  $\exp(\beta s)$ ,  $s \in \mathbb{R}$ .

Furthermore, for a fixed  $j$ , the probability  $\text{softmax}(a_{i_1 i_2})_{ij}$  can be consider as the interaction strength between the  $i$ -th and  $j$ -th particles.

In the following example, the self-attention mechanism, the  $i$ -th particle would be the column vector  $q_i^t$  in the matrix  $Q^t$ , and the  $j$ -th particle would be the column vector  $k_j^t$  in the matrix  $K^t$ .

**Definition 2** (Self-attention). *With the notations  $Q, K$  and  $V$ , defined as above, the self-attention mechanism is the following function*

$$\text{Attention}(Q, K, V) := \text{softmax}(\beta (QK^T)) V,$$

where  $\beta := \frac{1}{\sqrt{d_k}} \in \mathbb{R}$ .

In self-attention, this  $\beta$  is  $\frac{1}{\sqrt{d_k}}$ , and  $T$  is the standard deviation of the scalar product  $\sqrt{d_k}$ .

By definition of softmax, apply softmax to the matrix  $\frac{QK^T}{\sqrt{d_k}}$  means to apply Softmax function on each column of the matrix  $\left(\frac{QK^T}{\sqrt{d_k}}\right)_{N \times N}$ , and for defining the self-attention function, take a matrix multiplication between  $\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)_{N \times N}$  and the value matrix  $V_{N \times d_v}$ .

**The question** that asked and answered by this paper is the following.

Because of the random initiation of each weight matrices, if this self-attention, i.e. weights assigning process were run in parallel, then the model can learn from more than one different perspectives. In [2], this notion is defined as multi-head attention. Each self-attention function  $\text{Attention}(Q, K, V)_i$  is called  $\text{head}_i$ , where  $i$  is from 1 to  $H$ . The output of all  $H$  heads are concatenated to a tensor denoted by  $\text{Concat}(\text{head}_1, \dots, \text{head}_H)$ , then defined  $\text{Multihead}(Q, K, V) := \text{Concat}(\text{head}_1, \dots, \text{head}_H) W^O$ , where  $W^O \in \mathbb{R}^{H \cdot d_v \times \mathbb{R}^{d_{\text{out}}}}$ .

Let  $\epsilon > 0$  be given. Because of the mathematical form of the Boltzmann distribution, if the input is  $z_i + \epsilon$ , then  $\text{Softmax}(\beta z_i + \epsilon) = \text{Softmax}(\beta z_i)$ . However, a naive but natural question to ask is what if the scaling factor  $\beta$  in Boltzmann distribution, or the self-attention function, is modified into the form  $\beta^\epsilon$  or  $\beta \cdot \epsilon$ ? It is necessary to ask this question, since we want to be able to justify why and when it makes sense to use the scaling factor  $\beta = \frac{1}{\sqrt{d_k}}$  in the self-attention function. Is

it sufficient to ask this question? That is the question this paper is going to answer. The next section provides a solution to this question.

However, a mathematical and statistical analysis on the question raised and answered in this paper was not found in the literature, and the issue pointed out by this paper can still be found in the most recent publications that implemented self-attention in their models.

## 2 Algorithm

To be able to make sense to use the scaling factor  $\beta = \frac{1}{\sqrt{d_k}}$  in the self-attention function, the four assumptions are necessary.

Within the four assumptions, the first assumption is usually fixed when self-attention is considered to be applied.

Since the way query and key matrices are defined, i.e.  $Q = XW_q$  and  $K = XW_k$ , assumption (iii) and (iv) are based on assumption (ii).

Suppose the assumptions (ii) is relaxed, the variance of the scalar product  $q_i k_i$  is approximated to  $\left(d_x \sigma_X^2 \sigma_{W_q}^2\right) \cdot \left(d_x \sigma_X^2 \sigma_{W_k}^2\right) \cdot d_k$ . Since the training set is given, i.e. it is fixed, and so are  $d_x$  and  $d_k$ ,  $\sigma_{W_q}$  and  $\sigma_{W_k}$  are the actual moving factors.

What can cause  $\sigma_{W_q}$  and  $\sigma_{W_k}$  to be changed? In the training phase, the gradient descent algorithm or the backpropagation is used and the algorithm is going to shift the initial probability distribution of each entry in the random matrices  $W_q$  and  $W_k$  to the final probability distribution at a local minimum in the moduli space (or on the loss surface).

Hence, the scaling factor  $\beta = \frac{1}{\sqrt{d_k}}$  is only valid at the initial point given the four assumptions at the initial point in the moduli space, and it can hold true if the thermodynamic system described by the Boltzmann distribution is in statics, i.e. its temperature is fixed. However, since the initial point is random and unlikely to be the final point in the moduli space throughout the training mode, the assumption (ii) is not valid right after the training of the model that used self-attention function is initiated. Without assumption (ii), assumption (iii) and (iv) are invalid throughout the training and predicting modes.

To solve this problem, the solution could be written as an algorithm that it measures  $\sigma_{W_q}$  and

$\sigma_{W_k}$  and updates the scaling factor  $\beta$  in after each iteration or batch during the training and the best scaling factor  $\beta$  that is going to be used in the predicting mode could be determined by the best model chosen from the ensemble of all trained models.

This also answers the question asked in the previous section that it shows the sufficiency to ask the question.

In practice, 709 is the upper bound of the exponent of the exponential function based  $e$  in lots of machines before the machines start to overflow. Likewise, on each machine, there exists a non-negative integer  $N$  such that the machine starts to overflow the result. Hence, an alternative algorithm is: firstly try  $\beta \in \{10^N, \frac{1}{10^N}\}$  for  $N = 0, 1, 2, \dots$ , then run a binary search to narrow down the range of the best scaling factor  $\beta$  during the training after either each iteration or batch.

A further meaning of this algorithm is shown in the next section by using an empirical example.

### 3 Empirical example

The task of the empirical example is to flip a given string of non-negative integers. In other words, the goal is to train a multihead in a transformer to learn a function that maps each index  $i$  of an integer in a string to the index  $(n - i) + 1$  in the new string that will be returned by the model. Since the model needs to learn to swap the first and last elements in the string, it could be a challenging task for recurrent neural nets if the length of the input string is large.

In this empirical example, when the string length is 20, the range of the non-negative integers is between 0 and 100, and the number of samples used in training is 15000, 1000 samples for validation, and 100000 samples for testing.

Then, without running the above alternative algorithm, i.e. when the scaling factor  $\beta$  is  $\frac{1}{\sqrt{d_k}}$  the validation accuracy is 1.51%, and the test accuracy is 1.50%. To have a stable test accuracy above 95%, the training sample size needs to be increased from 15000 to 25000.

However, by applying the above alternative algorithm, the scaling factor is in  $O(1)$ , a constant  $\sim 5$ , with training sample size at 15000, the validation accuracy and test accuracy can already reach above 95% stably. Thus, the above algorithms can not only compress the training set when the training samples are limited or rare, but also reduce the training time to reach the best accuracy.

Furthermore, the model that achieved the above accuracy did not achieve it by memorizing the training examples. Since a 100% accuracy is also achievable within the ensemble, and once achieved that 100% accuracy the model can flip any input string with the same length and the range of the values of its each word. However, if it was by memorizing, then the model had to memorize  $100^{20}$  strings with length 20 which greatly exceed the number of parameters of the model (less than 8000).

## 4 Discussion

The empirical results in the previous section shows that with dynamic scaling factor being implemented, the model can reach an accuracy higher than 96% with about half of the training data compared to the fixed scaling factor introduced in the self-attention.

Our empirical model only has one task. If it is a large model, the number of features needed to be learned by a large model such as LLMs are several orders of magnitude higher than the model that only needs to learn one task. Hence, the difference of the number of training samples needed between models that use dynamic or static scaling factor will be an exponential growth when the number of features that need to be learned grows.

## 5 Conclusion

## 6 Appendix: Deriving the scaling factor $\frac{1}{\sqrt{d_k}}$ used in [2]

We can derive the assumption in [2] for the variance of  $QK^T$  written in the hidden dimension  $d_k$ :

**Proposition 1.** *Let  $q_i$  and  $k_i$  be random variables with zero means, i.e.  $\mathbb{E}q_i = 0 = \mathbb{E}k_i, \forall i$ , and unit variances, i.e.  $\text{Var}q_i = 1 = \text{Var}k_i$ , where  $i \in \{1, \dots, d_k\}$ , and they are all independent in the following ways:*

$$\text{Cov}(q_i, k_i) = 0, \forall i,$$

$$\text{Cov}(q_i, q_j) = 0, \forall i \neq j, \text{ and}$$

$$\text{Cov}(k_i, k_j) = 0, \forall i \neq j.$$

Then

$$\text{Var}(QK^T) = d_k.$$

*Proof.*

$$\begin{aligned} & \text{Var}(QK^T) \\ &= \text{Var}\left(\sum_{i=1}^{d_k} q_i k_i\right) \\ &= \sum_i \text{Var}(q_i k_i) \\ &= \sum_i \text{Var}(q_i) \text{Var}(k_i) \\ &= \sum_i 1 \cdot 1 \\ &= d_k. \end{aligned}$$

□

The first equality is because in [2],  $Q$  and  $K$  are row vectors, and  $QK^T$  is taking their scalar product. The second equality is because  $\text{Cov}(q_i, q_j) = 0, \forall i \neq j$ , and  $\text{Cov}(k_i, k_j) = 0, \forall i \neq j$ . Then, the variance of the sum is the sum of variances. The third equality is because of the conditions  $\mathbb{E}q_i = 0 = \mathbb{E}k_i, \forall i$  and  $\text{Cov}(q_i, k_i) = 0, \forall i$ , hence we have the variance of the product is the product of variances.

The authors in [2] want to scale the result of  $QK^T$  using the standard deviation of the scalar product, so  $QK^T$  is multiplied by a factor

$$\frac{1}{\sqrt{\text{Var}(QK^T)}} = \frac{1}{\sqrt{d_k}}.$$

## References

- [1] Bec Crew, *Google scholar reveals its most influential papers for 2020*, Nature Index **13** (2020).

- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, *Attention is all you need*, Advances in neural information processing systems **30** (2017).