

From Weighted Mean to FlashAttention

liuzhenhai93@outlook.com

May 31, 2025

In this note, self-attention is approached from the view of weighted mean. Then from this insight, flash attention and ring attention are deciphered in a unified way.

1 The Self-Attention

The computation of self attention in it's simplest form can be described as:

$$O = softmax(QK^T)V \quad (1)$$

which can be further factorized to 3 stage:

$$S = QK^T \quad (2)$$

$$P = softmax(S) \quad (3)$$

$$O = PV \quad (4)$$

For the convenience of subsequent discussions, two extra matrices are defines beforehand.

$$W = exp(S) = P * rowsum(P) \quad (5)$$

$$W_{safe} = exp(S - rowmax(S)) = W * exp(-rowmax(S)) \quad (6)$$

2 Weighted Mean

Formally, the weighted mean of a non-empty finite tuple of data (x_1, x_2, \dots, x_n) , with the corresponding non-negative weights (w_1, w_2, \dots, w_n) is

$$\bar{x} = \frac{\sum_1^n w_k * x_k}{\sum_1^n w_k} \quad (7)$$

2.1 Online Weighted Mean

Weighted mean can be calculated in a streaming way, with the data and weight feed one by one (or block by block). And there are two algorithms (1, 2).

3 Weighted Mean and Self-Attention

Since the sum of each row of P is 1. Each row of O can be considered as a weighted mean of the rows of V , with the corresponding row of P as weights.

Mathematically, If the weights are scaled by a common factor, the weighted mean remains the same. Hence, each row of O can also be considered as a weighted mean of the rows of V , with the corresponding row of W or W_{safe} as weights.

4 Flash Attention: Online Weighted Mean

With numerical stability taken into consideration, each row of O can be regarded as a weighted mean of the rows of V , with the corresponding row of W_{safe} as weights. Resembling the two algorithms for online weighted mean, there are two algorithms (3, 4) for computing O , corresponding to flash attention 1 and flash attention 2 respectively.

5 Ring Attention: Hierarchy Weighted Mean

If we define

$$\bar{x}_{(i,j)} = \frac{\sum_i^j w_k * x_k}{\sum_i^j w_k} \quad (8)$$

$$d_{(i,j)} = \sum_i^j w_k \quad (9)$$

then the global weighted mean can be computed as a weighted mean of two local weighted mean

$$\bar{x} = \frac{d_{(1,m)}}{d_{(1,n)}} \bar{x}_{(1,m)} + \frac{d_{(m+1,n)}}{d_{(1,n)}} \bar{x}_{(m+1,n)} \quad (10)$$

The equation above assumes that the elements are split into two groups. However, it can be extrapolated to any group:

$$\bar{x} = \frac{d_{(1,m_1)}}{d_{(1,n)}} \bar{x}_{(1,m_1)} + \frac{d_{(m_1+1,m_2+1)}}{d_{(1,n)}} \bar{x}_{(m_1+1,m_2+1)} + \dots + \frac{d_{(m_k+1,n)}}{d_{(1,n)}} \bar{x}_{(m_k+1,n)} \quad (11)$$

This is the mathematics behind ring attention: V are split into blocks; each block attention results in a partial result, and the final result is a weighted mean of the partial results.

Algorithm 1 one pass weighted mean 1

```
1:  $\bar{x}_0 = 0$ 
2:  $d_0 = 0$ 
3: for  $i = 1, 2, \dots, N$  do
4:    $d_i = d_{i-1} + w_i$ 
5:    $\bar{x}_i = \frac{d_{i-1}}{d_i} \bar{x}_{i-1} + \frac{w_i}{d_i} x_i$ 
6: end for
7:  $\bar{x} = \bar{x}_N$ 
```

Algorithm 2 one pass weighted mean 2

```
1:  $\bar{x}_0 = 0$ 
2:  $d_0 = 0$ 
3: for  $i = 1, 2, \dots, N$  do
4:    $d_i = d_{i-1} + w_i$ 
5:    $\bar{x}_i = \bar{x}_{i-1} + w_i x_i$ 
6: end for
7:  $\bar{x} = \frac{\bar{x}_N}{d_N}$ 
```

Algorithm 3 flash attention 1

```
1:  $m_0 = -\infty$ 
2:  $d'_0 = 0$ 
3:  $o'_0 = \text{row vector of } 0$ 
4: for  $i = 1, 2, \dots, N$  do
5:    $x_i = Q[k, :] K^T[:, i]$ 
6:    $m_i = \max(m_{i-1}, x_i)$ 
7:    $d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$ 
8:    $o'_i = \frac{d'_{i-1} e^{m_{i-1} - m_i}}{d'_i} o'_{i-1} + \frac{e^{x_i - m_i}}{d'_i} V[i, :]$ 
9: end for
10:  $O[k, :] = o'_N$ 
```

Algorithm 4 flash attention 2

```
1:  $m_0 = -\infty$ 
2:  $d'_0 = 0$ 
3:  $o'_0 = \text{row vector of } 0$ 
4: for  $i = 1, 2, \dots, N$  do
5:    $x_i = Q[k, :] K^T[:, i]$ 
6:    $m_i = \max(m_{i-1}, x_i)$ 
7:    $d'_i = d'_{i-1} e^{m_{i-1} - m_i} + e^{x_i - m_i}$ 
8:    $o'_i = o'_{i-1} + e^{x_i - m_i} V[i, :]$ 
9: end for
10:  $O[k, :] = \frac{o'_N}{d'_N}$ 
```
