

attention 的计算

给定 $q, k, v: \{i=1 \dots N\}$

$op_1 = \text{BMM}$

$$S_i = \text{dot}(q, k_i) \quad \{i=1 \dots N\}$$

$$m(x) = \max_i S_i$$

$$f(x) = [\exp(S_1 - m(x)) \dots \exp(S_N - m(x))]$$

为了数值计算的稳定性, 要把未归一化的值全部计算完后, 求 max value, 并减过去, 这会导致 $N \times M$ 的中间结果

$$l(x) = \sum_i f(x)_i$$

改善思路 ① 归一化因子 $l(x)$ 可以移动到整个 attention 计算的最后

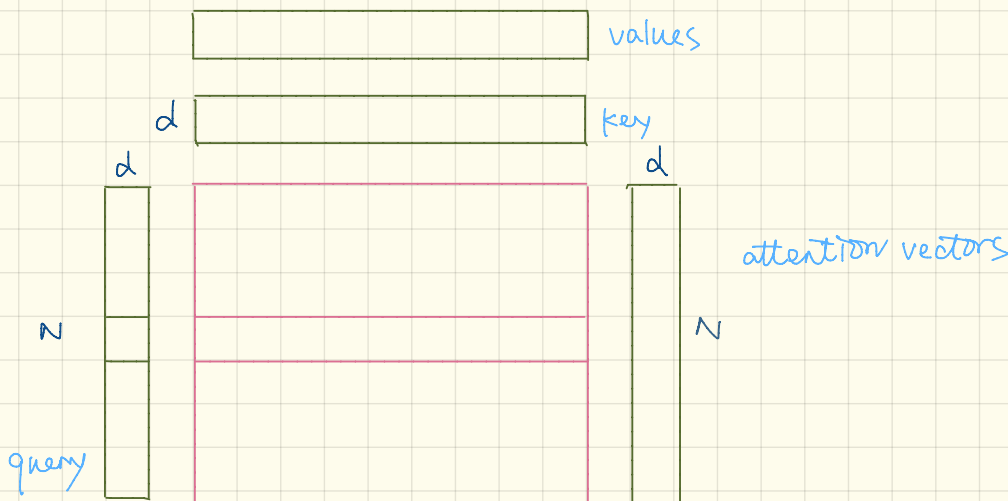
$op_2 = \text{softmax}$

$$\text{softmax} = \frac{f(x)}{l(x)}$$

② 每次计算一小块, 只考虑前几行的数值稳定性和归一化, 再计算新的一块时重新归一化

$$\text{attention}(q, k, v) = \frac{\sum_i f(x)_i v_i}{l(x)}$$

$op_3 = \text{BMM}$



我们来考虑一种最简单的情儿, 把 q, k, v 分成两块来计算 attention

$x = [a, b]_{2N}$ a 和 b 是一个维度为 N 的向量.

$$q = [q_a, q_b] \quad k = [k_a, k_b] \quad v = [v_a, v_b]$$

$$a_i = \text{dot}(q_a, k_{a_i})$$

计算第一块

$$m(a) = \max_i a_i$$

$$f(a) = [\exp(a_1 - m(a)) \dots \exp(a_N - m(a))]$$

$$l(x) = \sum_i f(a)_i$$

$$o_a = \text{attention}(q_a, k_a, v_a) = \frac{\sum_i f(a)_i v_i}{\sum_i f(a)_i}$$

$$b_i = \text{dot}(q_b, k_{b_i})$$

计算第二块

$$m(b) = \max_i b_i$$

$$f(b) = [\exp(b_1 - m(b)) \dots \exp(b_N - m(b))]$$

$$l(b) = \sum_{i=1}^N f(b)_i$$

$$m(x) = \max(m(a), m(b)) \leftarrow \text{需要开括号进行重新归一化}$$

$$\left\{ \begin{array}{l} e^{m(a) - m(x)} \\ e^{m(b) - m(x)} \end{array} \right\} \begin{array}{l} \text{是两块的重新归一化因子, 很容易验证} \\ e^{m(a) - m(x)} \quad e^{m(a) - m(x) + a_1 - m(a)} \end{array} = e^{a_1 - m(x)}$$

相对每一个新块都累加式地重新归一化

归一化

于是有如果考虑到重新归一化，第二块的计算变为：

$$f(x) = \left[e^{m(a)-m(x)} f(a) \quad e^{m(b)-m(x)} f(b) \right]$$

$$\begin{aligned} l(x) &= e^{m(a)-m(x)} \sum_i^N f(a_i) + e^{m(b)-m(x)} \sum_i^N f(b_i) \\ &= e^{m(a)-m(x)} l(a) + e^{m(b)-m(x)} l(b) \end{aligned}$$

$$\begin{aligned} \text{attention}(g, k, v) &= \frac{e^{m(a)-m(x)} f(a) V_a + e^{m(b)-m(x)} f(b) V_b}{l(x)} \\ &= \frac{O_a \cdot l(a) \cdot e^{m(a)-m(x)} + e^{m(b)-m(x)} f(b) V_b}{l(x)} \end{aligned}$$

从第二块的计算过程我们可以得到以下观测结果：

存储上一块的 max value: $m(a)$ 和归一化因子 $l(a)$ 。在新
一块计算时持续得更新 $m(x)$ 和 $l(x)$ ，就能够多次每次
只计算一小块的方式计算 softmax。

相当于可以将 softmax 与后面的矩阵乘法进一步进行融合。

上面只考虑了两块的情况，NV 的论文中同样已经证明
了分 N 块时，同样成立。只同维护 $m(x)$ 、 $l(x)$ 持续累加
结果到输出 O。