

位置编码

相对位置编码

首先区分下绝对位置编码和相对位置编码：

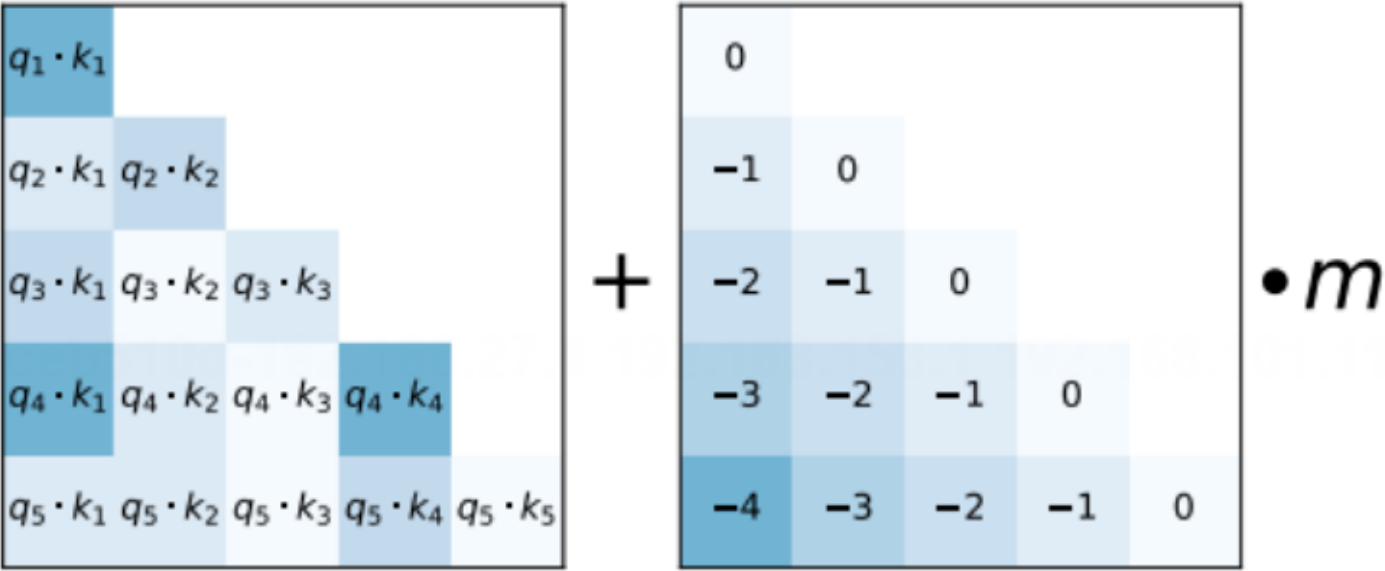
- 绝对位置编码直接将token的绝对位置信息添加到对应的 q_t, k_s 中,例如循环位置编码。

$$\begin{aligned} A_{t,s} &= q_t^T k_s = (x_t + p_t)^T W_Q^T W_K (x_s + p_s) \\ A_{t,s} &= x_t^T W_Q^T W_K x_s + x_t^T W_Q^T W_K p_s + p_t^T W_Q^T W_K x_s \\ &\quad + p_t^T W_Q^T W_K p_s \end{aligned} \tag{2}$$
$$p_t = \left[\cdots, \sin \frac{t}{10000^{2n/d}}, \cos \frac{t}{10000^{2n/d}}, \cdots \right]^T \in \mathbb{R}^d$$

- 相对位置编码则将两个token的相对位置信息添加到对应的attention值中。

1. Attention with linear biases enables input length extrapolation (ALiBi)

在计算attention时，对前边位置的分数进行惩罚，如图所示：



- 传统的绝对位置编码（如正弦/余弦编码）在训练时会为每个位置分配一个固定的向量，模型可能会过度拟合这些特定长度的模式。而ALiBi通过在注意力分数计算中直接使用线性偏置，减少了对模型对特定序列长度的依赖，从而提高了对未见过的序列长度的泛化能力。

- ALiBi的位置偏差随距离线性增长，这种设计让模型在处理不同长度的序列时，可以自然地根据距离调整注意力权重，无需显式学习位置编码的复杂周期性结构。
- 由于线性偏置的引入直接与序列中元素的位置相关，没有固定大小的编码矩阵限制，理论上模型可以更容易地处理任意长度的序列，从而展现出良好的长度外推性能。
- 通过直接在注意力分数上施加与距离相关的线性惩罚，ALiBi鼓励模型关注更近的位置，同时不完全排除远处的依赖，从而在一定程度上平衡了局部和全局依赖的学习，这对于处理长序列尤其有利。

2. XLENT

从绝对位置编码出发：

$$q_i k_j^T = x_i W_Q W_K^T x_j^T + x_i W_Q W_K^T p_j^T + p_i W_Q W_K^T x_j^T + p_i W_Q W_K^T p_j^T \quad (7)$$

XLENT将 p_j 替换成相对位置向量 R_{i-j} ， p_i 替换成可训练的向量 u 和 v ，

$$x_i W_Q W_K^T x_j^T + x_i W_Q W_K^T R_{i-j}^T + u W_Q W_K^T x_j^T + v W_Q W_K^T R_{i-j}^T \quad (8)$$

$$r_{t-s} = \left[\cdots, \sin \frac{t-s}{10000^{2n/d}}, \cos \frac{t-s}{10000^{2n/d}}, \cdots \right]^T \in \mathbb{R}^d$$

3. T5

(7) 式中的每一项可以理解为“输入-输入”、“输入-位置”、“位置-输入”、“位置-位置”四项注意力的组合，由于输入信息与位置信息应该是独立（解耦）的，它们不应该有过多的交互，所以“输入-位置”、“位置-输入”两项Attention可以删掉，“位置-位置”实际上是一个依赖于 $(s,)$ 的一个标量。

此外，通过固定的桶函数 $b(t-s)$ ，将 $t-s$ 从 $[-128, 128]$ 压缩至 $[0, 31]$ ，再对每个 $b(t-s)$ 训练对应的偏移量 $r_{b(t-s)}$

$$A_{t,s} = x_t^T W_Q^T W_K x_s + r_{b(t-s)} \quad (4)$$

4. DeBERTa

与T5相反，扔掉“位置-位置”一项只保留剩下三项，通过通过 $\delta(t, s)$ 将 $t-s$ 直接截断在区间 $(-k, k]$ 内

$$A_{t,s} = x_t^T W_Q^T W_K x_s + x_t^T W_Q^T W_K P_{\delta(t,s)} + x_s^T W_K^T W_Q P_{\delta(s,t)} \quad (7)$$

DeBERTa在softmax时校正系数为 $\sqrt{3d}$ ，不是默认的 \sqrt{d} 。此外，指出NLP的大多数任务可能都只需要相对位置信息，但确实有些场景下绝对位置信息更有帮助，于是它将整个模型分为两部分来理

解。以Base版的MLM预训练模型为例，它一共有13层，前11层只是用相对位置编码，这部分称为Encoder，后面2层加入绝对位置信息，这部分它称之为Decoder；至于下游任务的微调截断，则是使用前11层的Encoder加上1层的Decoder来进行。

RoPE

RoPE实现了绝对位置编码和相对位置编码的统一，它通过绝对位置编码的形式，实现了相对位置编码的效果。

RoPE将输入序列的位置信息通过旋转操作嵌入到self-attention的计算中，不同位置的 token 可以有不同的旋转角度，从而嵌入位置信息，从而增强模型对长序列和相对位置的处理能力。RoPE是可学习的，在自注意力公式中结合了明确的相对位置依赖性。

RoPE保持了序列长度的灵活性、随相对距离的增加而衰减的token间依赖性。其原理如下图， x'_m 表示经过RoPE后的结果：

从内积角度 证明目标: 在 x_m 表达式为 (1) 的情况下，内积能够表示成位置的函数，即表达式 (3) 成立

$$x'_m = W_q x_m e^{im\theta} = (W_q x_m) e^{im\theta} = q_m e^{im\theta} \quad (1) \quad \text{假设词向量只有2维}$$

$$x'_n = W_k x_n e^{in\theta} = (W_k x_n) e^{in\theta} = k_n e^{in\theta} \quad (2)$$

$$x_m^T x'_n = (q_m^1 \ q_m^2) \begin{pmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix} \quad (3)$$

从内积的角度推导：

$$x'_m = W_q x_m e^{im\theta} = (W_q x_m) e^{im\theta} = q_m e^{im\theta}$$

$q_m = \begin{pmatrix} W_q^{11} & W_q^{12} \\ W_q^{21} & W_q^{22} \end{pmatrix} \begin{pmatrix} x_m^1 \\ x_m^2 \end{pmatrix} = \begin{pmatrix} q_m^1 \\ q_m^2 \end{pmatrix}$

二维向量可以表示成虚数形式
应为向量的一个维度表示一个坐标
类比到实数坐标和虚数坐标

$$q_m = q_m^1 + i q_m^2$$

$e^{im\theta} = \cos(m\theta) + i \sin(m\theta) \quad \text{欧拉公式}$

$$i^2 = -1$$

$$= q_m e^{im\theta} = (q_m^1 + i q_m^2)(\cos(m\theta) + i \sin(m\theta))$$

$$= (q_m^1 \cos(m\theta) - q_m^2 \sin(m\theta)) + i(q_m^2 \cos(m\theta) + q_m^1 \sin(m\theta))$$

再写成向量形式

$$= [q_m^1 \cos(m\theta) - q_m^2 \sin(m\theta), q_m^2 \cos(m\theta) + q_m^1 \sin(m\theta)] = \begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix} \begin{pmatrix} q_m^1 \\ q_m^2 \end{pmatrix}$$

带入到内积中

$$x_m^T x'_n = \left(\begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix} \begin{pmatrix} q_m^1 \\ q_m^2 \end{pmatrix} \right)^T \begin{pmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix}$$

$$= (q_m^1 \ q_m^2) \begin{pmatrix} \cos(m\theta) & \sin(m\theta) \\ -\sin(m\theta) & \cos(m\theta) \end{pmatrix} \begin{pmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix}$$

$$\begin{aligned}
x_m^T x_n' &= \left(\begin{pmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{pmatrix} \begin{pmatrix} q_m^1 \\ q_m^2 \end{pmatrix} \right)^T \begin{pmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix} \\
&= (q_m^1 \ q_m^2) \begin{pmatrix} \cos(m\theta) & \sin(m\theta) \\ -\sin(m\theta) & \cos(m\theta) \end{pmatrix} \begin{pmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix} \\
&= (q_m^1 \ q_m^2) \begin{pmatrix} \cos(m\theta)\cos(n\theta) + \sin(m\theta)\sin(n\theta) & -\cos(m\theta)\sin(n\theta) + \sin(m\theta)\cos(n\theta) \\ -\sin(m\theta)\cos(n\theta) + \cos(m\theta)\sin(n\theta) & \sin(m\theta)\sin(n\theta) + \cos(m\theta)\cos(n\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix} \\
&= (q_m^1 \ q_m^2) \begin{pmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{pmatrix} \begin{pmatrix} k_n^1 \\ k_n^2 \end{pmatrix}
\end{aligned}$$

$$\begin{aligned}
\sin(m\theta + n\theta) &= \sin(m\theta)\cos(n\theta) + \cos(m\theta)\sin(n\theta) \\
\sin(m\theta - n\theta) &= \sin(m\theta)\cos(n\theta) - \cos(m\theta)\sin(n\theta) \\
\cos(m\theta + n\theta) &= \cos(m\theta)\cos(n\theta) - \sin(m\theta)\sin(n\theta) \\
\cos(m\theta - n\theta) &= \cos(m\theta)\cos(n\theta) + \sin(m\theta)\sin(n\theta)
\end{aligned}$$

对于多维的旋转位置编码，可以简化为以下形式：

$$\begin{pmatrix} \cos m\theta_0 & -\sin m\theta_0 & 0 & 0 & \dots & 0 & 0 \\ \sin m\theta_0 & \cos m\theta_0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \cos m\theta_1 & -\sin m\theta_1 & \dots & 0 & 0 \\ 0 & 0 & \sin m\theta_1 & \cos m\theta_1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos m\theta_{d/2-1} & -\sin m\theta_{d/2-1} \\ 0 & 0 & 0 & 0 & \dots & \sin m\theta_{d/2-1} & \cos m\theta_{d/2-1} \end{pmatrix} \begin{pmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_{d-2} \\ q_{d-1} \end{pmatrix}$$

可以看到矩阵中有很多元素是0，
如果用矩阵乘法，其实有很多计算是无用的 😊

$$\begin{pmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_{d-2} \\ q_{d-1} \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_0 \\ \cos m\theta_0 \\ \cos m\theta_1 \\ \cos m\theta_1 \\ \vdots \\ \cos m\theta_{d/2-1} \\ \cos m\theta_{d/2-1} \end{pmatrix} + \begin{pmatrix} -q_1 \\ q_0 \\ -q_3 \\ q_2 \\ \vdots \\ -q_{d-1} \\ q_{d-2} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_0 \\ \sin m\theta_0 \\ \sin m\theta_1 \\ \sin m\theta_1 \\ \vdots \\ \sin m\theta_{d/2-1} \\ \sin m\theta_{d/2-1} \end{pmatrix}$$

$$\theta_i = 10000^{-2i/d}, i \in [1, 2, \dots, \frac{d}{2}]$$

d 一定是偶数

结合代码来看，在chatglm中，因为内积计算与顺序无关，巧妙地将所有负数和正数分开

$$\theta_i = 10000^{-2i/d}, i \in [1, 2, \dots, \frac{d}{2}]$$

```
theta = 1. / (self.base ** (torch.arange(0, self.d, 2).float() / self.d)).to(x.device)
10000
```

获得一句话
中每个字的顺序

```
[0, 1, ..., seq_len - 1]
```

```
seq_idx = torch.arange(seq_len, device=x.device).float().to(x.device)
```

将 θ_i 与顺序融合



$m\theta_i$

```
idx_theta = torch.einsum('m,d->md', seq_idx, theta)
```

对于第m个字,

$m\theta_0, m\theta_1 \dots m\theta_{\frac{d}{2}}, m\theta_0, m\theta_1 \dots m\theta_{\frac{d}{2}}$

```
idx_theta2 = torch.cat([idx_theta, idx_theta], dim=1)
```

```
cos_cached = idx_theta2.cos()[:, None, None, :]
sin_cached = idx_theta2.sin()[:, None, None, :]
```

一定需要 按维度顺序两两组合
旋转么? 😊

不一定要严格按照相邻两两组合旋转
神经元是无序的, 不依赖维度顺序

$$\begin{pmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_{d/2-2} \\ q_{d/2-1} \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_0 \\ \cos m\theta_0 \\ \cos m\theta_1 \\ \cos m\theta_1 \\ \vdots \\ \cos m\theta_{d/2-2} \\ \cos m\theta_{d/2-1} \end{pmatrix} + \begin{pmatrix} -q_1 \\ q_0 \\ -q_3 \\ q_2 \\ \vdots \\ -q_{d/2-1} \\ q_{d/2-2} \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_0 \\ \sin m\theta_0 \\ \sin m\theta_1 \\ \sin m\theta_1 \\ \vdots \\ \sin m\theta_{d/2-2} \\ \sin m\theta_{d/2-1} \end{pmatrix}$$

Chatglm 中

```
def rotate_half(x):
    """Rotates half the hidden dims of the input."""
    x1 = x[..., : x.shape[-1] // 2]
    x2 = x[..., x.shape[-1] // 2 :]
    return torch.cat((-x2, x1), dim=-1)

x_rope, x_pass = x[..., :self.d], x[..., self.d:]
neg_half_x = rotate_half(x_rope)
x_rope = (x_rope * cos_cached[:x.shape[0]]) +
[neg_half_x * sin_cached[:x.shape[0]]]
```

$$\begin{pmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ \vdots \\ -q_{d/2} \\ \vdots \\ -q_{d-2} \\ -q_{d-1} \end{pmatrix} \begin{pmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_{d/2} \\ \vdots \\ q_{d-1} \\ q_{d-2} \end{pmatrix}$$

从公式的角度推理, 要计算m和n之间的距离:

$$g(x'_m, x'_n, m-n) = \text{Re}[(W_q x_m)(W_k x_n)^* e^{i(m-n)\theta}]$$

Re 表示, 取实数部分的结果

$$= \text{Re}[q_m k_n^* e^{i(m-n)\theta}]$$

$$* \text{表示, 共轭复数} \quad k_n = k_n^1 + i k_n^2 \quad k_n^* = k_n^1 - i k_n^2$$

$$= \text{Re}[(q_m^1 + i q_m^2)(k_n^1 - i k_n^2)(\cos((m-n)\theta) + i \sin((m-n)\theta))]$$

$$\text{欧拉公式} \quad e^{i(m-n)\theta} = \cos((m-n)\theta) + i \sin((m-n)\theta)$$

$$= \text{Re}[(q_m^1 k_n^1 + q_m^2 k_n^2 + i(q_m^2 k_n^1 - q_m^1 k_n^2))(\cos((m-n)\theta) + i \sin((m-n)\theta))]$$

$$= \text{Re}[(q_m^1 k_n^1 + q_m^2 k_n^2) \cos((m-n)\theta) - (q_m^2 k_n^1 - q_m^1 k_n^2) \sin((m-n)\theta) + i(q_m^2 k_n^1 - q_m^1 k_n^2) \cos((m-n)\theta) + i(q_m^1 k_n^1 + q_m^2 k_n^2) \sin((m-n)\theta)]$$

不要复数, 只保留实数部分



$$= (q_m^1 k_n^1 + q_m^2 k_n^2) \cos((m-n)\theta) - (q_m^2 k_n^1 - q_m^1 k_n^2) \sin((m-n)\theta)$$

总结下RoPE的流程: 首先计算得到q和k矩阵, 然后对q和k向量的元素按照两两一组应用RoPE, 例如:

$$\mathbf{f}(\mathbf{x}, m) = [(x_0 + i x_1) e^{i m \theta_0}, (x_2 + i x_3) e^{i m \theta_1}, \dots, (x_{d-2} + i x_{d-1}) e^{i m \theta_{d/2-1}}]^T$$

那么对于m和n之间的考虑位置距离的注意力就是:

$$a(m, n) = \text{Re} \langle \mathbf{f}(\mathbf{q}, m), \mathbf{f}(\mathbf{k}, n) \rangle$$

$$= \text{Re} \left[\sum_{j=0}^{d/2-1} (q_{2j} + i q_{2j+1})(k_{2j} - i k_{2j+1}) e^{i(m-n)\theta_j} \right]$$

$$= \sum_{j=0}^{d/2-1} (q_{2j} k_{2j} + q_{2j+1} k_{2j+1}) \cos((m-n)\theta_j) + (q_{2j} k_{2j+1} - q_{2j+1} k_{2j}) \sin((m-n)\theta_j)$$

得到的注意力随m与n之间的距离增大而减小, 注意对v矩阵不需要应用RoPE。

xPOS

RoPE能扩展到任意长度，但其外推性能较差：虽然RoPE可以拓展到任意长度，但对于语言建模等生成任务，无法在测试长序列性能时，维持其在训练长度序列上的表现。xPos在旋转的基础上，在旋转角度向量的每个维度上都包含了独特的指数衰减因子，以及blockwise causal attention，让模型忽略相距较远的语义：

$$\mathbf{A}_{t,s} = \text{Re} \left[\sum_{n=1}^{d/2} \left(\zeta_n^{-t} e^{it\theta_n} \tilde{q}_t^{(n)} \right)^* \zeta_n^s e^{is\theta_n} \tilde{k}_s^{(n)} \right] = \sum_{n=1}^{d/2} \text{Re} \left[\tilde{q}_t^{(n)*} \tilde{k}_s^{(n)} \zeta_n^{s-t} e^{i(s-t)\theta_n} \right] \quad (16)$$
$$\zeta_n = \frac{1 + \gamma}{2n/d + \gamma}, \quad \gamma = 0.4, \quad \theta_n = 10000^{-2n/d}$$

参考：[通俗易懂-大模型的关键技术之一：旋转位置编码

rope(https://www.bilibili.com/video/BV12x42127Pb/?spm_id_from=333.788&vd_source=4804e44ea59abe1f4c2b7dde30651898)]

[Long-Context LLM综述(https://blog.csdn.net/Cyril_KI/article/details/139573263)]

[Transformer位置编码（基础）(<https://zhuanlan.zhihu.com/p/631363482>)]

[Transformer升级之路：7、长度外推性与局部注意力(<https://zhuanlan.zhihu.com/p/602487515>)]