# Transformer

（Sequence to Sequence Model，□□）

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Context□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## □□□□□□

1. **Word Embedding（□□□□□）**

- □□□□□□□□□□ID□□□□□□□□□□ 512□□
- □□□□□□□ `nn.Embedding` □□□□

2. **Positional Encoding（□□□□□）**

- □□ Transformer □□ **RNN** □□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□
    - □□□□□□□□□□□□□□□□
    - □□□□□□□□□□□□□□□□□

## Encoder（□□□□）

□□ Encoder □□□□□□□□□□**block**□□□□□ 6 □□□

□□□□**block**□□□□□□□□□□

1. **□□□□□□□（Multi-Head Self Attention）**
    - □□□□□□□ → □□"□ □ □□"□□□□□□□□□□

2. **□□□□（Feed Forward Network）**
    - □□□□□□□□□□□□□□□□□ MLP
    - □□□□□□□□□

□□□□□□□□□

- □□□□
- **Layer Normalization**

## Decoder（□□□□）

□□□□□□□

- □□□□□□□□□□□□□□□□□□□□□
- （**Auto-regressive**，□□□□□）

□□□□

□□□ Decoder □□ **3 □□□ + □□□□ + LayerNorm）**

- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
- □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

1. **Masked Multi-Head Self-Attention（□□□□□□□□）**

    - □□□□□□□□□□□□"□□□□□□□□□"
    - □□□□□ Decoder □□□□"□□□□□" □ □□□□□

2. **Encoder-Decoder Attention（□□□□□□）**

    - □□□□ Decoder □□□ Encoder □□□□□□□

- - - **Query 来自 Decoder，Key 和 Value 来自 Encoder 的输出：**
    - - 让 Decoder 的"当前"位置去关注源句子中所有位置，从而让 Decoder"知道"应该翻译/生成/关注哪些内容

  3. **Feed Forward Network（前馈神经网络）**
     - - 两层全连接神经网络 + 激活函数（**ReLU/GeLU**）

  4. **残差连接 + LayerNorm**

     每个子层后面都有：
     - - 残差连接： `output = input + Sublayer(input)`
     - LayerNorm（层归一化）：稳定训练过程

- 编码器和解码器都由**n**个相同的层堆叠而成


image-20250702225431058

# Train（训练过程）

  1. **Encoder：**
     - - 输入完整的源语言句子（比如一句中文），一次性并行处理
     - 输出每个词的上下文表示（编码后的向量）

  2. **Decoder：**
     - - 输入是目标语言句子（比如英文），也就是 label ，训练时已知
     - 但是是逐词预测的，比如看到"我 爱 ___"预测"你"

  3. **Teacher Forcing：**
     - - 训练时，解码器的输入不是它自己预测出来的**token**，而是"标准答案"，即真实目标句子的上一个token
     - 这样能让Decoder 更快收敛，不会因为早期错误而偏离
     - 举个例子：哪怕它预测错 了 一个词，下一个词照样正确输入

# Attention：

- **Decoder 中使用 Self-Attention（ Mask 掉未来的词，防止偷看答案），让模型理解上下文**
- **Encoder-Decoder Attention：** Decoder 的每个位置会"关注" Encoder 的所有输出，来决定生成哪个词（对源句子进行对齐）

输出预测

- 解码器最后输出 → **softmax** → 每个词的概率分布（预测）
- 与真实词的 one-hot 标签对比 → 计算 **Cross Entropy Loss**

损失函数

模型目标是最大化预测正确词的概率

实际优化方式：最小化交叉熵损失（负对数似然）

# Teaching Forcing

在训练时，Decoder 的输入是"真实的目标句子"，而不是模型自己"上一步预测的结果"，这种方式叫 `Teacher Forcing` （教师强制法）。

比如 Decoder 第一部分

- 刚开始我们只有一个起始标记 `<BOS>` 作为第一个
- 我们期望通过 这个标记推理出第一个字"我"作为下一个
- 我们将期待的输出"我 爱"
- ...一次类推

即 每一 部我们的期待输出和下一部的输入是一致的

这种方法叫做 **Teacher Forcing**

## Residual Connection（残差连接）

是什么？

残差连接是一种神经网络结构设计，用于解决"深度"和"梯度消失"的问题。它的做法是把输入x直接加到F(x)（经过一层或几层变换）上：x + F(x)

为什么有用？

- 缓解梯度消失：由于 $y = F(x) + x$，反向传播时
- 反向传播时的梯度，即 ∂L/∂x 永远至少有一部分是

$$ \frac{\partial{L}}{\partial{y}} \cdot\frac{\partial{y}}{\partial{x}} = \frac{\partial{L}}{\partial{y}} \cdot (\frac{\partial{F(x)}}{\partial{x}} + 1) $$

- 即使 $\frac{\partial{F(x)}}{\partial{x}}$ 很小，加了1也不会梯度消失的太快

## Layer Normalization

为什么使用层归一化而不是批归一化？

- 在 BatchNorm 里，每个特征是在一个批次上做归一化的，依赖于整个批次
- 在 NLP 任务中，用 BatchNorm 会遇到麻烦，因为句子长度不同，batch 内部统计不稳定
- **Layer Normalization，就是对每个feature内不同dimension归一化，而Batch Normalization就是对不同feature同一个dimension归一化**

## Explanation

这部分我们将重新复述self-attention的运行过程

有四个向量


image-20250427121123206

> 1. 通过向量相乘我们得到四个向量$a^1,a^2,a^3,a^4$
> 2. 通过乘以变换矩阵我们得到以下矩阵
>
> $$ \begin{aligned} Q = \begin{bmatrix}q^1 \ q^2 \ q^3 \ q^4\end{bmatrix} &= \begin{bmatrix}a^1 \ a^2 \ a^3 \ a^4\end{bmatrix}W^q \[1em] K = \begin{bmatrix}k^1 \ k^2 \ k^3 \ k^4\end{bmatrix} &= \begin{bmatrix}a^1 \ a^2 \ a^3 \ a^4\end{bmatrix}W^k \[1em] V = \begin{bmatrix}v^1 \ v^2 \ v^3 \ v^4\end{bmatrix} &= \begin{bmatrix}a^1 \ a^2 \ a^3 \ a^4\end{bmatrix}W^v \ \end{aligned} $$
>
> 2. 计算**attention score** $$ \begin{aligned} A = \begin{bmatrix}\alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4}\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} \ \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} \ \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4}\end{bmatrix} &=Q \cdot K^T\ &=

$$\begin{bmatrix}q^1 \ q^2 \ q^3 \ q^4\end{bmatrix} \cdot \begin{bmatrix}k^1 & k^2 & k^3 & k^4\end{bmatrix} \end{aligned} $$

3. 除以$\sqrt{d_k}$，再做**softmax**。$d_k$是指$key/querey$的维度大小。

$$ \begin{aligned} &\quad\quad\quad\quad\quad\quad\text{注意力 } A \xrightarrow{\text{softmax}} \text{注意力 } A^{'} \ &\frac{1}{\sqrt{d_k}}\begin{bmatrix}\alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4}\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} \ \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} \ \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4}\end{bmatrix} \quad \Rightarrow \quad \begin{bmatrix}\alpha^{'}{1,1} & \alpha^{'}{1,2} & \alpha^{'}{1,3} & \alpha^{'}{1,4}\ \alpha^{'}{2,1} & \alpha^{'}{2,2} & \alpha^{'}{2,3} & \alpha^{'}{2,4} \ \alpha^{'}{3,1} & \alpha^{'}{3,2} & \alpha^{'}{3,3} & \alpha^{'}{3,4} \ \alpha^{'}{4,1} & \alpha^{'}{4,2} & \alpha^{'}{4,3} & \alpha^{'}{4,4}\end{bmatrix} \end{aligned} $$

4. 乘以归一化后的注意力矩阵

$$ [b^1,b^2,b^3,b^4] = \begin{bmatrix}\alpha^{'}{1,1} & \alpha^{'}{1,2} & \alpha^{'}{1,3} & \alpha^{'}{1,4}\ \alpha^{'}{2,1} & \alpha^{'}{2,2} & \alpha^{'}{2,3} & \alpha^{'}{2,4} \ \alpha^{'}{3,1} & \alpha^{'}{3,2} & \alpha^{'}{3,3} & \alpha^{'}{3,4} \ \alpha^{'}{4,1} & \alpha^{'}{4,2} & \alpha^{'}{4,3} & \alpha^{'}{4,4}\end{bmatrix} \cdot \begin{bmatrix}v^1 \ v^2 \ v^3 \ v^4\end{bmatrix} $$

即 $$ Attention(Q.K,V) = softmax(\frac{QK^T}{\sqrt{d_k}})V $$

多头注意力

区别在于每一个头都有各自的权重矩阵$W^{q,1},W^{q,2}$

多头注意力的权重矩阵计算稍有不同：

- 将输入序列（$a^i,a^j,a^m,a^n$）分别和$W^q、W^k、W^v$三个权重矩阵相乘，得到查询、键、值$Q、K、V$

- 初始化$W^q、W^k、W^v$三个权重矩阵

$$ Q = \begin{bmatrix} a^i\ a^j\a^m\a^n \end{bmatrix}W^q\ K = \begin{bmatrix} a^i\ a^j\a^m\a^n \end{bmatrix}W^k\ V = \begin{bmatrix} a^i\ a^j\a^m\a^n \end{bmatrix}W^v $$

- 初始化多头权重矩阵

将上一步得到的初始化权重矩阵继续分别和$Q、K、V$相乘。 $$ \begin{align} Q_1 = QW^{q,1} \quad Q_2 = QW^{q,2} \ K_1 = KW^{q,1} \quad K_2 = KW^{q,2}\ V_1 = VW^{q,1} \quad V_2 = VW^{q,2} \end{align} $$

- 得到多头信息

$$ \begin{align} Q_1 &= \begin{bmatrix} q^{i,1}\ q^{j,1}\q^{m,1}\q^{n,1} \end{bmatrix} Q_2 = \begin{bmatrix} q^{i,2}\ q^{j,2}\q^{m,2}\q^{n,2} \end{bmatrix} \ \end{align} $$

$$ K_1 = \begin{bmatrix} k^{i,1}\ k^{j,1}\k^{m,1}\k^{n,1} \end{bmatrix} K_2 = \begin{bmatrix} k^{i,2}\ k^{j,2}\k^{m,2}\k^{n,2} \end{bmatrix} $$

$$ V_1 = \begin{bmatrix} v^{i,1}\ v^{j,1}\v^{m,1}\v^{n,1} \end{bmatrix} V_2 = \begin{bmatrix} v^{i,2}\ v^{j,2}\v^{m,2}\v^{n,2} \end{bmatrix} $$

- 分别计算每个头的结果：

$$ head_1= softmax(\frac{Q_1K_1^T}{\sqrt{d_k}})V_1\ head_2= softmax(\frac{Q_2K_2^T}{\sqrt{d_k}})V_2 $$

- 拼接每个头：

$$ multihead = \begin{bmatrix} head_1 & head_2 \end{bmatrix} $$

- 乘以权重：

$$ output = multihead\cdot W^O $$

# Self-attention

给定一段向量序列，如何输出一段新的向量序列？

一种做法是FC层，但是单纯的FC无法考虑上下文的关系，效果并不会好

### Sequence Labeling（序列标注问题）

一种朴素想法

- 结合上下文信息，即将前后几个向量一起串起来丢进网络（即开一个FC的窗口）

#### 整体思路

其整体思路是：对于每个向量，都要考虑到整个句子的信息，并以此决定最终的输出

其中，有几个概念：

- 查询向量$(Query, Q)$
- 键向量$(Key, K)$
- 值向量$(Value, V)$

其中计算：

- 对于每个输入向量，都有三个对应的向量，即$Q、K、V$三个
    - 其计算方法是每个向量乘以权重：
        - $Q = X × W^Q$
        - $K = X × W^K$
        - $V = X × W^V$

        其中$W^Q、W^K、W^V$都是可学习的参数

- 计算相关性：用查询向量(Q)和其他向量的键向量(K)，计算得到attention score，即相关性分数
    - dot product：点积，即对应元素相乘再相加，代表的是它们之间的相似度，点积越大，代表这两个向量方向越接近，相关性越强
        - 查询向量(Q)代表的**"我想找什么信息"**
        - 键向量(K)代表的**"我有什么信息可以被找"**

- 这是一种对**"该从句子中提取哪些信息"的注意力机制**

- 将每一个向量经过点乘后得到的分数进行**softmax**操作，使所有注意力分数之和为一

- 使用softmax的原因是相对于RELU等激活函数它能够有更好的表现

- 将注意力分数与对应的值向量(V)相乘，再将所有的attention score进行加权求和

  - 值向量(V)包含该位置的实际语义内容
  - Q-K的相似度是**"应该关注多少"**的权重
  - V则是被加权的**"实际要提取的信息"的语义内容**

## Multi-head Self-attention（多头自注意力）

不同的相关性使得$Q$有多个，因此对应的自注意力也有多个

- 对于$a^i$乘上两个不同的矩阵
  - 对于$b^{i,1}$
    - 使用$q^{i,1}、k^{i,1}、k^{j,1}$得到了$b^{i,1}$
  - 对于$b^{i,2}$
    - 使用$q^{i,2}、k^{i,2}、k^{j,2}$得到了$b^{i,2}$
  - 对于$b^{i}$
    - 使用$b^{i,1}$、$b^{i,2}$乘上一个矩阵得到了$b^{i}$，即为**attention score**
  - 其余位置的计算同理


image-20250420112828589

## Masked Multi-Head Self-Attention(掩码多头注意力)

掩码多头注意力的运用场景是当我们需要对未来的信息进行掩盖时使用

因为在预测的时候我们是不知道下文的内容的，所以需要对下文进行掩盖，假如此时有一个序列 "A B C D" ，当我们对第二个token做注意力分数时，需要掩盖后文信息

所以实际上我们掩盖的信息是对于K和V来说 $$ \begin{bmatrix} 1 & masked &masked&masked\\ 1 & 1 &masked & masked\\ 1 &1& 1 &masked\\ 1&1&1&1 \end{bmatrix} $$ image-20250702191639250

## Word Embedding（词嵌入）

对于输入的文字，我们可以用**One-Hot Encoding**来进行编码，但是这样做会使得编码过于稀疏，并且无法表示词与词之间的关系

- 所以我们需要对词进行词嵌入操作
- 词嵌入操作可以将词映射到一个低维空间

## Positional Encoding

- **Self-attention**在计算的过程中并没有"位置的信息"，所以我们需要对输入的词进行位置编码，使得模型能够知道词的位置
- 不管距离的远近，不管左边还是右边，对于自注意力机制来说都是一样的

位置编码的实现方式如下：

1. 为每个位置生成一个**positional vector**，这个向量的维度与词向量的维度相同
2. 将这个位置向量与词向量相加，得到最终的输入向量

**Final_embedding = Token_embedding + Positional_encoding**

□□□□□□

□□□□□□□□□□□□□□

□□□□□□□□□□□

□□□
$$ PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})\ PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}}) $$

1. $pos$□□□□□
2. $2i$□□□□□□$2i\le d_{model}$
3. $d_{model}$□□□□□□□□□□□□□□□□□□□□□$d_{model}$□□□□

- □□I am a kid
  - □□□□
    - I□pos□0
    - am□pos□1□□□□□□□
  - □□□□□□am□pos=1□
    - i=0
      - □0□□$PE_{(1,0)} = sin(1/10000^{0/d_{model}})$
      - □1□□$PE_{(1,1)}=cos(1/10000^{0/d_{model}})$
    - i=1
      - □2□□$PE_{(1,2)}=sin(1/10000^{2/d_{model}})$
      - □3□□$PE_{(1,3)}=cos(1/10000^{2/d_{model}})$

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

- □□□□□$w_i = 10000^{\frac{2i}{d_{model}}}$
  - □□□□□□□
    - □□□□□□$pos/w_i$□□□□□□□□□□10000□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
    - □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LSTM□□□□□□□□□□□
  - □□□□□□
    - □□□□□□□□□□□$d_{model}$□□□□□$w_i$□□□□□□□□□□□$10000^{2i/d_{model}}$□□□□□□□□□□□pos□□□$ pos/10000^{2i/d_{model}}$□□□□□□□□□□□□□□□
    - □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□

$$ \underbrace{\begin{bmatrix} PE_{(pos + \Delta,2i)} \ PE_{(pos + \Delta,2i+1)} \end{bmatrix}}_{□□(pos + \Delta)□□□□}$$

\underbrace{ \begin{bmatrix} cos(\Delta \theta_i) & sin(\Delta \theta_i)\ -sin(\Delta \theta_i) & cos(\Delta \theta_i) \end{bmatrix}}{□□□□□□} \underbrace{ \begin{bmatrix} PE{(pos ,2i)} \ PE_{(pos ,2i+1)} \end{bmatrix}}_{□□pos□□□□} $$

其中，$\Delta$为相对位置，并且$\theta_i = \frac{1}{10000^{\frac{2i}{d_{model}}}}$

故 $$ \underbrace{ \begin{bmatrix} cos(\Delta \theta_i) & sin(\Delta \theta_i)\\ -sin(\Delta \theta_i) & cos(\Delta \theta_i) \end{bmatrix}}_{旋转矩阵，\Delta\theta_i} $$

故旋转矩阵是固定的，只和相对位置$\Delta\theta_i$有关，而与绝对位置**pos**无关。



故对于固定偏移量 $$ \begin{align} \begin{bmatrix} PE_{(pos + \Delta,2i)} \\ PE_{(pos + \Delta,2i+1)} \end{bmatrix}&=\begin{bmatrix} sin((pos+\Delta)\cdot\theta_i) \\ cos((pos+\Delta)\cdot\theta_i) \end{bmatrix}\\ &=\begin{bmatrix} sin(pos\cdot\theta_i)cos(\Delta\cdot\theta_i)+cos(pos\cdot\theta_i)sin(\Delta\cdot\theta_i) \\ cos(pos\cdot\theta_i)cos(\Delta\cdot\theta_i)-sin(pos\cdot\theta_i)sin(\Delta\cdot\theta_i) \end{bmatrix}\\ &=\begin{bmatrix} cos(\Delta \theta_i) & sin(\Delta \theta_i)\\ -sin(\Delta \theta_i) & cos(\Delta \theta_i) \end{bmatrix} \begin{bmatrix} sin(pos\cdot\theta_i)\\ cos(pos\cdot\theta_i) \end{bmatrix}\\ &=\begin{bmatrix} cos(\Delta \theta_i) & sin(\Delta \theta_i)\\ -sin(\Delta \theta_i) & cos(\Delta \theta_i) \end{bmatrix} \begin{bmatrix} PE_{(pos,2i)}\\ PE_{(pos,2i+1)} \end{bmatrix}

\end{align} $$

**Transformer**的位置编码选择合适的波长？
对于： $$ \begin{align}

w_i&= 10000^{2i/d_{model}} \\ \theta_i &= pos/10000^{2i/d_{model}} \\ &= pos / w_i

\end{align} $$

- i越大，$w_i$越大，波长越大（周期越大），$\theta_i$随pos增大变化越慢，频率越低；
- i越小，$w_i$越小，波长越小（周期越小），$\theta_i$随pos增大变化越快，频率越高。当趋近于0时，对于$\Delta pos=1000$，$\frac{1000}{10000} = 0.1$，基本上没有变化。
- 如果$w_i$是10000，$sin(pos/10000)$在pos到达$20000\pi$时，才能完整变化一个周期，意味着位置编码的分辨率很低，在很长的距离范围内，位置编码的值都几乎相同，难以区分不同位置。

相对位置编码

- 绝对位置编码，是在输入前就加上位置编码，然后在做$QK^T$的注意力计算时，得到了相对位置编码的形式；相对位置编码是直接在注意力的计算中，加入相对位置信息。
- 相对位置编码不是作用在输入上，而是作用在注意力**(query与key的计算)**上，因为直接作用在**注意力**上，$Q$和$K$之间具有相对位置信息。

对于绝对位置编码，在融入相对位置信息的时候，对于$Q$和$K$的注意力： $$ q_i = W_q(x_i+p_i)\\ k_j=W_k(x_j+p_j) $$ 注意力为： $$ q_i^Tk_j = (x_i+p_i)^TW_q^TW_k(x_j+p_j)\\ =\underbrace{x_i^TW_q^TW_kx_j}_{纯内容交互}+\underbrace{x_i^TW_q^TW_kp_j+p_i^TW_q^TW_kx_j}_{内容-位置交互}+\underbrace{p_i^TW_q^TW_kp_j}_{纯位置相关的} $$ 相对位置编码的做法，就是去除掉绝对位置编码

$$ q_i^Tk_j =x_i^TW_q^TW_kx_j + \underbrace{p_i^TW_q^TW_kp_j}_{纯位置相关的} \\ = x_i^TW_q^TW_kx_j +\underbrace{\beta_{i-j}}_{纯位置相关的} $$ 其$\beta_{i-j}$是一个只和相对位置相关的标量，并且可能每个注意力头$head_h$都有独立的可训练的$\beta_{i-j}^h$(why??/)

**T5**的相对位置编码
添加一个可训练的标量，这个标量只和相对位置有关，并且是直接在注意力上的添加

旋转位置编码

**ALibi**（线性偏差注意力）？？？

所以把$QK$矩阵的内积，沿着对角线偏移进行对应的添加，其中$m$


image-20250702192315616

其中的m是一个常数，具体为$m_h=2^{-\frac{8\times h}{n\_head}}$，$n_{head}$是注意力头的个数 $$ m=\begin{bmatrix}2^{-\frac{8\times 1}{n\_head}}\\2^{-\frac{8\times 2}{n\_head}}\\ \cdots\\2^{-\frac{8\times n_{head}}{n\_head}}\end{bmatrix} $$ 因为是偏移的位置，所以相对位置是$D_{ij} = -(i - j)$

带入计算就是 $$ softmax(q_iK^T+m\cdot[-(i-1),\dots,-2,-1,0]) $$

**RoPE（旋转位置编码）**

核心：通过绝对位置编码的方式实现**Q、K的相对位置编码**

对于自注意力机制，最核心的一个运算就是内积 $$ Attention\_score = q_m\cdot k_n^T $$ 我们希望经过某个操作之后，q、k就带有了相对位置的信息 $$ q^{rot} = q_m\cdot R_m \ k^{rot} = k_n \cdot R_n $$ 那么在 $$ \begin{aligned} q^{rot}\cdot k^{rot} &= q_mR_mR_n^Tk_n^T \\ &= q_mR_mR_{-n}k_n^T \\ &= q_mR_{m-n}k_n^T \end{aligned} $$ 这样就可以完成在计算$Q、K$的时候添加了**RoPE**位置编码，那么应该怎么去进行设计这个操作呢？

1. 内积满足线性叠加性

   因此，任意偶数维的旋转编码，我们都可以表示为二维情形的拼接，所以我们只需要先研究二维的情况下，然后进行拼接。

2. 二维情况下

 假设在位置为 $m,n$ 的绝对位置编码，我们希望得到$m-n$

所以二维推导：

  - 嵌入向量（d维度），RoPE将其分为$d/2$个二维子空间，每个子空间有一个**d的旋转角度**

$$ R_{\theta,m} = \begin{bmatrix} \cos m\theta_{1} & -\sin m\theta_{1} & 0 & \cdots & 0 \\ \sin m\theta_{1} & \cos m\theta_{1} & 0 & \cdots & 0 \\ 0 & 0 & \cos m\theta_{2} & -\sin m\theta_{2} & 0 \\ 0 & 0 & \sin m\theta_{2} & \cos m\theta_{2} & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \end{bmatrix} $$

  - 对于每个子空间，旋转操作可以表示为

$$

# \begin{bmatrix} x_{i}' \ x_{i+1}' \end{bmatrix}

\begin{bmatrix} \cos m\theta_i & -\sin m\theta_i \\ \sin m\theta_i & \cos m\theta_i \end{bmatrix} \begin{bmatrix} x_i \ x_{i+1} \end{bmatrix}

$$

最终得到了$(m-n)$的相对编码

其中$\theta_i = \frac{1}{10000^{\frac{2i}{d_{model}}}}$ $$ R_{m-n} = \begin{bmatrix} cos((m-n) \theta_i) & sin((m-n) \theta_i)\\ -sin((m-n) \theta_i) & cos((m-n) \theta_i) \end{bmatrix} $$ 通过： $$ R_n^TR_m = R_{m-n} $$ 同时旋转矩阵还满足： $$ R_{-m} = R^T_m $$

可以变成相对的 $$ (R_nq)^TR_mk=q^TR_{m-n}k $$ 由此， $$ q_n^{rot} = R_nq\ k_m^{rot} = R_mk $$ 即旋转后的查询和键$Q，K$向量

## KV Cache

- 只存储前向传播中的
- 推理阶段用到的，
- 自回归时候只需要$Q，K，V$的值即可
- 历史的结果也用到了

因为自回归式的输出，让生成当前的内容依赖于之前生成的token，所以每次有重复计算的部分，把这些重复计算的部分，也就是$K，V$的值存储起来。

好处：

- 加速
- 低延迟
- 高吞吐量
- 显存的消耗

坏处：用KV Cache是用空间换时间，需要"存储"一定空间的$K，V$

### 没有的KV Cache



### 有的KV Cache



## Pre-Norm和Post-Norm

**Pre-Norm（先归一化）**和**Post-Norm（后归一化）**是Transformer架构中两种不同的层归一化位置策略，核心区别在于LN（Layer Normalization）的放置位置：

- 梯度传播稳定性：恒等路径上没有归一化操作，确保了梯度可以无衰减的传播
- 训练稳定性：由于梯度可以稳定的传播，因此训练Transformer的时候不需要采用复杂的学习率调整策略，即**学习率预热 (warm-up)**。
- 即插即用，可以方便的添加和删除对应网络层的操作