



Transformer相关——（10）Transformer 代码分析

📅 发表于 2021-08-19 | 🕒 更新于 2021-11-09 | 📁 深度学习
| 📖 字数总计: 6,956 | ⌚ 阅读时长: 34分钟 | 👁 阅读量: 2443

Transformer相关——（10） Transformer代码分析

引言

原理是原理，道理大概都懂了，代码也不能落下。这篇就把Transformer代码拿出来分析一下。代码来源：Pytorch编写完整的Transformer，我进一步做了一些修改和补充。

和之前一样，先把每个小模块分析一下，然后再把它们串起。来构建一整个model。

主要包括以下几个部分：

- 1 "input" embedding

- 2 position encoding
- 3 Multi-Head attention
- 4 Add&Norm

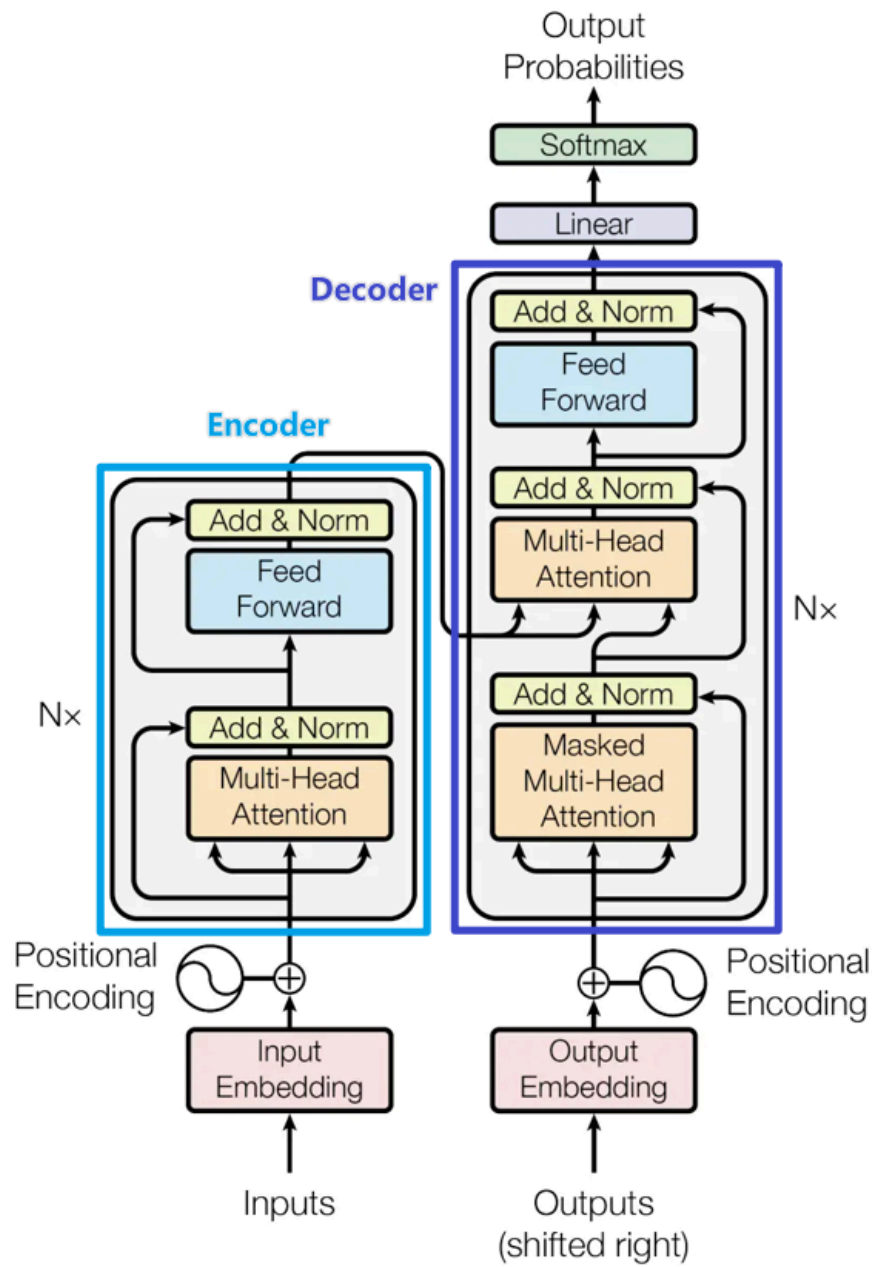


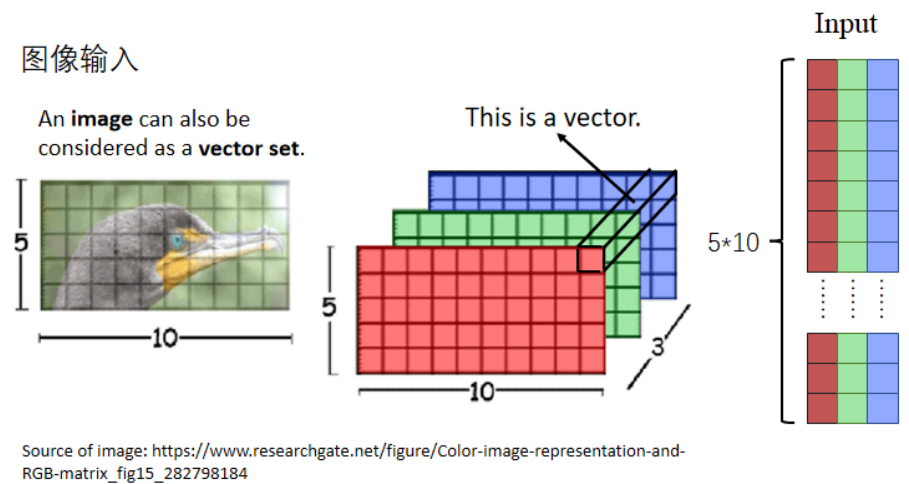
Figure 1: The Transformer - model architecture.

子模块

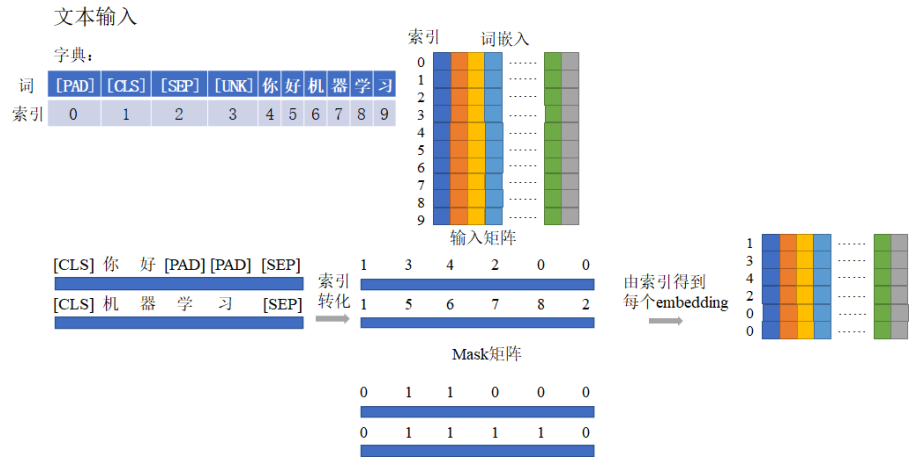
"input" embedding

这一步其实就是之前说的“将真实问题转为数学问题”，比如输入的数据是图像、语音、文本、用户ID等等，把这些内容转成数值矩阵。

像图像这种自带数值输入数据，可以经过转换直接作为input embedding，如下图所示：



对于语音、文本这一类数据，就需要用特殊的技巧转成数字信号输入。比如在文本中，就是使用**词嵌入**。示意图如下所示。



那么词嵌入（word embedding）怎么获得呢？

我觉得总的来说可以有三种方式：

- 1 根据字典构建**one-hot**编码，假设字典大小为 N 的话，这就是一个 $N \times N$ 的矩阵，词对应索引位置的向量值为1，其他位置为0，每个词的表示向量相互垂直，没有语义信息，字典太大这个矩阵也会很大很稀疏。
- 2 随机一个 $N \times M$ 矩阵， M 为自定义的每个词的特征维度大小。在网络中不断学习或者用一些方法求每个词在训练语料中的语义特

征，最后输出每个词的embedding；经典的方法比如DeepWalk、Word2Vec等，深度学系的方法比如图神经网络等等都可以用于学习词嵌入。

- 3 从预训练好的词嵌入中取。像现在NLP还非常流行的范式是pre-trained+fine tuning，就是从一些已经训练好的词嵌入中取对应词的embedding，然后进行微调作不同场景下的下游任务。

词嵌入在Pytorch里基于 `torch.nn.Embedding` 实现，实例化时需要设置的参数为词表的大小和被映射的向量的维度，`embed = nn.Embedding(vocab_size, embed_dim)`。 `vocab_size` 是词表的大小，词表中包括了一些特殊作用的字符（比如用于padding的[PAD]；表示句子开头的[CLS]；表示句子结束的[SEP]；表示字典中没有出现过的未知字符[UNK]等等）。

代码

```
1 import torch
2 import torch.nn as nn
3 X = torch.zeros((2,4),dtype=torch.long)
4 embed = nn.Embedding(10,8)
5 print(embed(X).shape)
```

position encoding位置编码

位置编码用以表达元素在序列中的位置特征，比如名词经常出现在句子开头。

位置编码直接与元素的embedding相加。

代码中需要注意：X_只是初始化的矩阵；完成位置编码之后会加一个 **dropout**。另外，位置编码是最后加上去的，因此输入输出形状不变。

这里使用的是GPT-3中的相对位置编码：

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

其中, t 为绝对位置, w_i 为:

$$w_i = \frac{1}{10000^{2i/d_{model}}}$$

代码

```
1  Tensor = torch.Tensor
2  def positional_encoding(X, num_features,
3  dropout_p=0.1, max_len=512) -> Tensor:
4      r'''
5          给输入加入位置编码
6          参数:
7              - embed_dim: 词嵌入维度
8              - dropout_p: dropout的概率, 当其为非零时执行 dropout
9              - max_len: 句子的最大长度, 默认512
10             形状:
11                 - 输入: [batch_size, seq_length,
12                     embed_dim]
13                 - 输出: [batch_size, seq_length,
14                     embed_dim]
15                 - seq_length表示这个batch中句子的长度 (每个 batch都做了截断或补齐操作)
16             例子:
17                 >>> X = torch.randn((2,4,10))
18                 >>> X = positional_encoding(X, 10)
19                 >>> print(X.shape)
20                 >>> torch.Size([2, 4, 10])
21             '''
22             dropout = nn.Dropout(dropout_p)
23             P = torch.zeros((1,max_len,num_features))
24             p_t =
25                 torch.arange(max_len,dtype=torch.float32).reshape(-1,
26                     / torch.pow(
27                         10000,
28                     torch.arange(0,num_features,2,dtype=torch.float32)
29                     /num_features) #每个位置的t*w_k
30             P[:, :, 0::2] = torch.sin(p_t) #0::2偶数位
31             P[:, :, 1::2] = torch.cos(p_t) #0::1奇数位
32             X = X + P[:, :, X.shape[1],:].to(X.device) #位置
```

码和初始embedding直接相加
29 return dropout(X)

```
1    # 位置编码例子
2    X = torch.randn((2,4,10)) #
   (batch_size,seq_length,embed_dim)
3    X = positional_encoding(X, 10)
4    print(X.shape)
```

Multi-Head attention

我们先来分析attention机制的代码再转到Multi-Head attention。

attention机制

Transformer用缩放点积相关性计算attention score：

$$\alpha_{i,j} = \frac{(q^i \cdot k^j)}{\sqrt{d}}$$

当矩阵特征向量以行向量形式表示时，attention的输出矩阵可以按照下述公式计算（以缩放点积相关性+softmax为例）：

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

代码

```
1    from typing import Optional, Tuple, Any
2    Tensor = torch.Tensor
3    def _scaled_dot_product_attention(
4        q: Tensor,
5        k: Tensor,
6        v: Tensor,
7        attn_mask: Optional[Tensor] = None,
8        dropout_p: float = 0.0,
9    ) -> Tuple[Tensor, Tensor]:
10        r'''
11            在query, key, value上计算点积注意力，若有注意力
           遮盖则使用，并且应用一个概率为dropout_p的dropout
12            参数：
```

```

13         - q: shape: `(B, Nt, E)` B代表batch
size, Nt是目标语言序列长度, E是嵌入后的特征维度
14         - key: shape: `(B, Ns, E)` Ns是源语言序列
长度
15         - value: shape: `(B, Ns, E)` 与key形状一样
16         - attn_mask: 要么是3D的tensor, 形状
为: `(B, Nt, Ns)` 或者2D的tensor, 形状如: `(Nt, Ns)`
17         - Output: attention values: shape: `(B,
Nt, E)`, 与q的形状一致; attention weights: shape: `(B,
Nt, Ns)`

```

```

18
19     例子:
20         >>> q = torch.randn((2,3,6))
21         >>> k = torch.randn((2,4,6))
22         >>> v = torch.randn((2,4,6))
23         >>> out =
scaled_dot_product_attention(q, k, v)
24         >>> out[0].shape, out[1].shape
25         >>> torch.Size([2, 3, 6])
torch.Size([2, 3, 4])
26     '''
27     B, Nt, E = q.shape #每个词的特征向量是行向量
28     q = q / math.sqrt(E)
29     # (B, Nt, E) x (B, E, Ns) -> (B, Nt, Ns)
30     attn = torch.bmm(q, k.transpose(-2,-1))
#Q·K^T
31     if attn_mask is not None:
32         attn = attn + attn_mask
33     # attn意味着目标序列的每个词对源语言序列做注意力
34     attn = F.softmax(attn, dim=-1)
#softmax(Q·K^T)
35     if dropout_p:
36         attn = F.dropout(attn, p=dropout_p)
37     # (B, Nt, Ns) x (B, Ns, E) -> (B, Nt, E)
38     output = torch.bmm(attn, v)
#softmax(Q·K^T)·V
39     return output, attn

```

获得Q、K、V

之前提到序列向量输入attention之前还需要经过三个线性变换分别生成Q、K、V。

$$Q = W^q X + b^q \quad (1)$$

$$K = W^k X + b^k \quad (2)$$

$$V = W^v X + b^v \quad (3)$$

利用 `nn.functional.linear` 函数实现线性变换，与 `nn.Linear` 不同的是，前者可以提供权重矩阵和偏置，执行 $y = xW^T + b$ ，而后者是可以自由决定输出的维度（比如下游分类任务是回归任务时，可设定输出维度为1）。

代码

把三个Q、K、V的权重拼在一起，用一个大的权重参数矩阵进行线性变换（也就是 $W=[W1;W2;W3]$ ），生成Q、K、V使用的大权重参数矩阵不同的分块。

```

1  def _in_projection_packed(
2      q: Tensor,
3      k: Tensor,
4      v: Tensor,
5      w: Tensor,
6      b: Optional[Tensor] = None,
7  ) -> List[Tensor]:
8      r"""
9      用一个大的权重参数矩阵进行线性变换
10     参数:
11         q, k, v: 对自注意来说，三者都是src；对于
seq2seq模型，k和v是一致的tensor。
12         但它们的最后一维(num_features或者
叫做embed_dim)都必须保持一致。
13         w: 用以线性变换的大矩阵，按照q,k,v的顺序压在一个tensor里面。
14         b: 用以线性变换的偏置，按照q,k,v的顺序压在一个tensor里面。
15
16     形状:
17     输入:
18         - q: shape: `(..., E)`，E是词嵌入的维度（下面出现的E均为此意）。
19         - k: shape: `(..., E)`
20         - v: shape: `(..., E)`
21         - w: shape: `(E * 3, E)`
22         - b: shape: `E * 3`

```



```

23
24         输出:
25         - 输出列表 : `[q', k', v']`, q,k,v经过线性
变换前后的形状都一致。
26         """"
27         E = q.size(-1)
28         # 若为自注意, 则q = k = v = src, 因此它们的引用
变量都是src
29         # 即k is v和q is k结果均为True
30         # 若为seq2seq, k = v, 因而k is v的结果是True,
但q!=k, 也就是cross attention
31         if k is v:
32             if q is k: # 自注意力
33                 return F.linear(q, w, b).chunk(3,
dim=-1)
34             else:
35                 # seq2seq模型
36                 w_q, w_kv = w.split([E, E * 2])
37                 if b is None:
38                     b_q = b_kv = None
39                 else:
40                     b_q, b_kv = b.split([E, E *
2])
41                 return (F.linear(q, w_q, b_q),) +
F.linear(k, w_kv, b_kv).chunk(2, dim=-1)
42         else: #q!=k!=v
43             w_q, w_k, w_v = w.chunk(3)
44             if b is None:
45                 b_q = b_k = b_v = None
46             else:
47                 b_q, b_k, b_v = b.chunk(3)
48             return F.linear(q, w_q, b_q),
F.linear(k, w_k, b_k), F.linear(v, w_v, b_v)

1   q, k, v = _in_projection_packed(query, key,
value, in_proj_weight, in_proj_bias)

```

获得 W^q, W^k, W^v (参数初始化)

由获得Q、K、V自然而然引出, 如何获得 $W^q, W^k, W^v, b^q, b^k, b^v$, 也就是如何进行参数初始化。

代码

在PyTorch的源码里，`projection`代表是一种线性变换的意思，
`in_proj_bias` 是指一开始的线性变换的偏置 b^q ，这里用
`q_proj_weight` 表示 W^q ，其他同理。

用 `torch.empty` 按照所给的形状形成对应的tensor，相当于预留了一个位置，其填充的值还未初始化，类比`torch.randn`（标准正态分布），这是一种初始化的方式。

在PyTorch中，`tensor`变量类型是无法修改值的，而`Parameter()`函数可以看作为一种类型转变函数，将不可改值的`tensor`转换为可训练可修改的模型参数（也就是可以随着模型一起学习训练），即与
`model.parameters`绑定在一起，`register_parameter`的意思是是否将这个参数放到`model.parameters`，`None`表示没有这个参数。

`if`判断句用于判断`q,k,v`的最后一维是否一致，若一致，则将这个大的权重矩阵全部乘然后分割出来，若不是，则各初始化各的，初始化不会改变原来的形状（如 $q = qW_q + b_q$ ，见注释）。

```
1     if self._qkv_same_embed_dim is False:
2         # 初始化前后形状维持不变
3         # (seq_length x embed_dim) x (embed_dim x
4         # embed_dim) ==> (seq_length x embed_dim)
5         self.q_proj_weight =
6         Parameter(torch.empty((embed_dim, embed_dim)))
7         self.k_proj_weight =
8         Parameter(torch.empty((embed_dim, self.kdim)))
9         self.v_proj_weight =
10        Parameter(torch.empty((embed_dim, self.vdim)))
11        self.register_parameter('in_proj_weight',
12        None)
13    else:
14        self.in_proj_weight =
15        Parameter(torch.empty((3 * embed_dim, embed_dim)))
16        self.register_parameter('q_proj_weight',
17        None)
18        self.register_parameter('k_proj_weight',
19        None)
20        self.register_parameter('v_proj_weight',
21        None)
```

```

13
14     if bias:
15         self.in_proj_bias =
Parameter(torch.empty(3 * embed_dim))
16     else:
17         self.register_parameter('in_proj_bias',
None)
18     # 后期会将所有头的注意力拼接在一起然后乘上权重矩阵输出
19     # out_proj是为了后期准备的
20     self.out_proj = nn.Linear(embed_dim,
embed_dim, bias=bias)
21     self._reset_parameters()

```

然后用 `_reset_parameters()` 函数初始化参数数值。`xavier_uniform` 是从连续型均匀分布里面随机取样出值来作为初始化的值，`xavier_normal` 取样的分布是正态分布。正因为初始化值在训练神经网络的时候很重要，所以才需要这两个函数。

`constant_`意思是用所给值来填充输入的向量。

```

1     def _reset_parameters(self):
2         if self._qkv_same_embed_dim:
3             xavier_uniform_(self.in_proj_weight)
4         else:
5             xavier_uniform_(self.q_proj_weight)
6             xavier_uniform_(self.k_proj_weight)
7             xavier_uniform_(self.v_proj_weight)
8         if self.in_proj_bias is not None:
9             constant_(self.in_proj_bias, 0.)
10            constant_(self.out_proj.bias, 0.)

```

给注意力机制补充Mask机制

Encoder和Decoder多头注意力机制不同有两处：

- 1 decoder的第一个多头注意力机制需要masked；
- 2 decoder的第二个多头注意力机制是cross attention。

这两点在 `_scaled_dot_product_attention` 中分别由 `attn_mask` 和 判断q与k是否相等 时进行了考虑，我们分析一下 `attn_mask` 这一输

入要怎么写。

这里再回忆一下为什么要进行mask:

- 1 sequence mask: 因为在decoder解码的时候, 只能看该位置和它之前的, 如果看后面就相当于提前知道答案了, 所以需要attn_mask遮挡住;
- 2 padding mask: 而key_padding_mask用于不计算较短句子用"[PAD]"字符补齐的部分。

代码

接下来会使用到两个函数:

- 1 logical_or, 输入两个tensor, 并对这两个tensor里的值做 逻辑或 运算, 只有当两个值均为0的时候才为 False, 其他时候均为 True;

```
1 a =  
  torch.tensor([0,1,10,0],dtype=torch.int8)  
2 b =  
  torch.tensor([4,0,1,0],dtype=torch.int8)  
3 print(torch.logical_or(a,b))  
4 # tensor([ True,  True,  True, False])
```

- 2 masked_fill, 输入是一个mask, 和用以填充的值。mask由1, 0组成, 0的位置值维持不变, 1的位置用新值填充。

```
1 r = torch.tensor([[0,0,0,0],[0,0,0,0]])  
2 mask = torch.tensor([[1,1,1,1],  
  [0,0,0,0]])  
3 print(r.masked_fill(mask,1))  
4 # tensor([[1, 1, 1, 1],  
5 #          [0, 0, 0, 0]])
```

sequence mask

对于attn_mask来说, 若为2D, 形状如(L, S), L和S分别代表着目标语言和源语言序列长度, 若为3D, 形状如(N * num_heads, L, S), N代表着batch_size, num_heads代表注意力头的数目。若为attn_mask的dtype为ByteTensor, 非0的位置会被忽略不做注意力; 若为BoolTensor, True对应的位置会被忽略; 若为数值, 则会直接加到attn_weights。

```
1 def generate_square_subsequent_mask(sz:
  int) -> Tensor:
2     r'''产生关于序列的mask, 被遮住的区域赋值`-inf`,
    未被遮住的区域赋值为`0`'''
3     mask = (torch.triu(torch.ones(sz, sz)) ==
4               1).transpose(0, 1) #返回(sz, sz)大小方阵的上三角部分, 其
    余部分定义为0, 转置后为下三角, 这个就是sequence mask。
5     mask = mask.float().masked_fill(mask == 0,
    float('-inf')).masked_fill(mask == 1, float(0.0)) #
    这里将有效的部分mask矩阵设置为0, 无效的部分设置为-inf
6     return mask
```

注意, 为什么这里我们把有效的部分mask为0呢, 注意前面attention机制中有这样的代码:

```
1 if attn_mask is not None:
2     attn = attn + attn_mask
```

这样我们在相加的时候, 可以不改变attn有效部分的值, 而把需要被遮盖的地方设置为了-inf, 从而实现了掩膜。

```
1 if attn_mask is not None:
2     if attn_mask.dtype == torch.uint8:
3         warnings.warn("Byte tensor for
    attn_mask in nn.MultiheadAttention is deprecated.
    Use bool tensor instead.")
4     attn_mask = attn_mask.to(torch.bool)
5     else:
6         assert attn_mask.is_floating_point()
    or attn_mask.dtype == torch.bool, \
7         f"Only float, byte, and bool types
    are supported for attn_mask, not {attn_mask.dtype}"
```

```

8      # 对不同维度的形状判定
9      if attn_mask.dim() == 2:
10         correct_2d_size = (tgt_len, src_len)
11         if attn_mask.shape != correct_2d_size:
12             raise RuntimeError(f"The shape of
the 2D attn_mask is {attn_mask.shape}, but should
be {correct_2d_size}.")
13         attn_mask = attn_mask.unsqueeze(0)
14     elif attn_mask.dim() == 3:
15         correct_3d_size = (bsz * num_heads,
tgt_len, src_len)
16         if attn_mask.shape != correct_3d_size:
17             raise RuntimeError(f"The shape of
the 3D attn_mask is {attn_mask.shape}, but should
be {correct_3d_size}.")
18     else:
19         raise RuntimeError(f"attn_mask's
dimension {attn_mask.dim()} is not supported")

```

padding mask

与 `attn_mask` 不同的是，`key_padding_mask` 是用来遮挡住key里面的值，详细来说应该是 [PAD]，被忽略的情况与 `attn_mask` 一致。

```

1  # 将key_padding_mask值改为布尔值
2  if key_padding_mask is not None and
key_padding_mask.dtype == torch.uint8:
3      warnings.warn("Byte tensor for
key_padding_mask in nn.MultiheadAttention is
deprecated. Use bool tensor instead.")
4      key_padding_mask =
key_padding_mask.to(torch.bool)

```

其实 `attn_mask` 和 `key_padding_mask` 有些时候对象是一致的，所以有时候可以合起来看。`-inf` 做softmax之后值为0，即被忽略。

```

1  if key_padding_mask is not None:
2      assert key_padding_mask.shape == (bsz,
src_len), \
3          f"expecting key_padding_mask shape of
{(bsz, src_len)}, but got {key_padding_mask.shape}"
4      key_padding_mask =

```

```

key_padding_mask.view(bsz, 1, 1,
src_len).expand(-1, num_heads, -1, -1).reshape(bsz
* num_heads, 1, src_len)
5         # 若attn_mask为空, 直接用key_padding_mask
6         if attn_mask is None:
7             attn_mask = key_padding_mask
8         elif attn_mask.dtype == torch.bool:
9             attn_mask =
attn_mask.logical_or(key_padding_mask)
10        else:
11            attn_mask =
attn_mask.masked_fill(key_padding_mask, float("-
inf"))
12
13    # 若attn_mask值是布尔值, 则将mask转换为float
14    if attn_mask is not None and attn_mask.dtype
== torch.bool:
15        new_attn_mask =
torch.zeros_like(attn_mask, dtype=torch.float)
16        new_attn_mask.masked_fill_(attn_mask,
float("-inf"))
17        attn_mask = new_attn_mask

```

Multi-Head attention

接下来分析一下Multi-Head attention的代码。

代码

之前在Attention机制中提到：

多头注意力机制往往满足：self-attention的隐藏层维度（ $hidden_emb_dim$ ） $\times n_heads$ = 输入特征（经过线性变换）的维度（ emb_dim ），然后 W^O 的维度为 emb_dim 。

所以令每个head的Q、K、Vd的

$head\ dim = embed_dim = emb_dim / n_heads$ 。Multi-Head attention核心代码如下。

```

1    import torch
2    Tensor = torch.Tensor

```

```

3  def multi_head_attention_forward(
4      query: Tensor,
5      key: Tensor,
6      value: Tensor,
7      num_heads: int,
8      in_proj_weight: Tensor,
9      in_proj_bias: Optional[Tensor],
10     dropout_p: float,
11     out_proj_weight: Tensor,
12     out_proj_bias: Optional[Tensor],
13     training: bool = True,
14     key_padding_mask: Optional[Tensor] = None,
15     need_weights: bool = True,
16     attn_mask: Optional[Tensor] = None,
17     use_separate_proj_weight = None,
18     q_proj_weight: Optional[Tensor] = None,
19     k_proj_weight: Optional[Tensor] = None,
20     v_proj_weight: Optional[Tensor] = None,
21 ) -> Tuple[Tensor, Optional[Tensor]]:
22     r'''
23     形状:
24         输入:
25         - query: `(L, N, E)`
26         - key: `(S, N, E)`
27         - value: `(S, N, E)`
28         - key_padding_mask: `(N, S)`
29         - attn_mask: `(L, S)` or `(N *
num_heads, L, S)`
30         输出:
31         - attn_output: `(L, N, E)`
32         - attn_output_weights: `(N, L, S)`
33     '''
34     tgt_len, bsz, embed_dim = query.shape
35     src_len, _, _ = key.shape
36     head_dim = embed_dim // num_heads
37     q, k, v = _in_projection_packed(query,
key, value, in_proj_weight, in_proj_bias)
38
39     if attn_mask is not None:
40         if attn_mask.dtype == torch.uint8:
41             warnings.warn("Byte tensor for
attn_mask in nn.MultiheadAttention is deprecated.
Use bool tensor instead.")

```



```

42         attn_mask =
attn_mask.to(torch.bool)
43     else:
44         assert
attn_mask.is_floating_point() or attn_mask.dtype ==
torch.bool, \
45         f"Only float, byte, and bool
types are supported for attn_mask, not
{attn_mask.dtype}"
46
47         if attn_mask.dim() == 2:
48             correct_2d_size = (tgt_len,
src_len)
49             if attn_mask.shape !=
correct_2d_size:
50                 raise RuntimeError(f"The shape
of the 2D attn_mask is {attn_mask.shape}, but
should be {correct_2d_size}.")
51             attn_mask = attn_mask.unsqueeze(0)
52             elif attn_mask.dim() == 3:
53                 correct_3d_size = (bsz *
num_heads, tgt_len, src_len)
54                 if attn_mask.shape !=
correct_3d_size:
55                     raise RuntimeError(f"The shape
of the 3D attn_mask is {attn_mask.shape}, but
should be {correct_3d_size}.")
56             else:
57                 raise RuntimeError(f"attn_mask's
dimension {attn_mask.dim()} is not supported")
58
59         if key_padding_mask is not None and
key_padding_mask.dtype == torch.uint8:
60             warnings.warn("Byte tensor for
key_padding_mask in nn.MultiheadAttention is
deprecated. Use bool tensor instead.")
61             key_padding_mask =
key_padding_mask.to(torch.bool)
62
63         if key_padding_mask is not None:
64             assert key_padding_mask.shape == (bsz,
src_len), \
65                 f"expecting key_padding_mask shape

```

```

of {(bsz, src_len)}, but got
{key_padding_mask.shape}"
66         key_padding_mask =
key_padding_mask.view(bsz, 1, 1,
src_len).expand(-1, num_heads, -1, -1).reshape(bsz
* num_heads, 1, src_len)
67         if attn_mask is None:
68             attn_mask = key_padding_mask
69         elif attn_mask.dtype == torch.bool:
70             attn_mask =
attn_mask.logical_or(key_padding_mask)
71         else:
72             attn_mask =
attn_mask.masked_fill(key_padding_mask, float("-
inf"))
73         # 若attn_mask值是布尔值, 则将mask转换为float
74         if attn_mask is not None and
attn_mask.dtype == torch.bool:
75             new_attn_mask =
torch.zeros_like(attn_mask, dtype=torch.float)
76             new_attn_mask.masked_fill_(attn_mask,
float("-inf"))
77             attn_mask = new_attn_mask
78
79         # reshape q,k,v将Batch放在第一维以适合点积注意
力
80         # 同时为多头机制, 将不同的头拼在一起组成一层
81         q = q.contiguous().view(tgt_len, bsz *
num_heads, head_dim).transpose(0, 1)
82         k = k.contiguous().view(-1, bsz *
num_heads, head_dim).transpose(0, 1)
83         v = v.contiguous().view(-1, bsz *
num_heads, head_dim).transpose(0, 1)
84
85         # 若training为True时才应用dropout
86         if not training:
87             dropout_p = 0.0
88         # Attention计算
89         attn_output, attn_output_weights =
_scaled_dot_product_attention(q, k, v, attn_mask,
dropout_p)
90         attn_output = attn_output.transpose(0,
1).contiguous().view(tgt_len, bsz, embed_dim)

```

```

91         attn_output =
nn.functional.linear(attn_output, out_proj_weight,
out_proj_bias)
92         if need_weights:
93             # average attention weights over heads
94             attn_output_weights =
attn_output_weights.view(bsz, num_heads, tgt_len,
src_len)
95             return attn_output,
attn_output_weights.sum(dim=1) / num_heads
96         else:
97             return attn_output, None

```

再结合参数初始化函数，完整代码如下：

```

1  class MultiheadAttention(nn.Module):
2      r'''
3      参数:
4          embed_dim: 词嵌入的维度
5          num_heads: 平行头的数量
6          batch_first: 若`True`，则为(batch, seq,
feature)，若为`False`，则为(seq, batch, feature)
7
8      例子:
9          >>> multihead_attn =
MultiheadAttention(embed_dim, num_heads)
10          >>> attn_output, attn_output_weights =
multihead_attn(query, key, value)
11          '''
12      def __init__(self, embed_dim, num_heads,
dropout=0., bias=True,
13                  kdim=None, vdim=None,
batch_first=False) -> None:
14          # factory_kwargs = {'device': device,
'dtype': dtype}
15          super(MultiheadAttention,
self).__init__()
16          self.embed_dim = embed_dim
17          self.kdim = kdim if kdim is not None
else embed_dim
18          self.vdim = vdim if vdim is not None
else embed_dim
19          self._qkv_same_embed_dim = self.kdim

```

```

== embed_dim and self.vdim == embed_dim
20
21         self.num_heads = num_heads
22         self.dropout = dropout
23         self.batch_first = batch_first
24         self.head_dim = embed_dim // num_heads
25         assert self.head_dim * num_heads ==
self.embed_dim, "embed_dim must be divisible by
num_heads"
26
27         if self._qkv_same_embed_dim is False:
28             self.q_proj_weight =
Parameter(torch.empty((embed_dim, embed_dim)))
29             self.k_proj_weight =
Parameter(torch.empty((embed_dim, self.kdim)))
30             self.v_proj_weight =
Parameter(torch.empty((embed_dim, self.vdim)))
31
self.register_parameter('in_proj_weight', None)
32         else:
33             self.in_proj_weight =
Parameter(torch.empty((3 * embed_dim, embed_dim)))
34
self.register_parameter('q_proj_weight', None)
35
self.register_parameter('k_proj_weight', None)
36
self.register_parameter('v_proj_weight', None)
37
38         if bias:
39             self.in_proj_bias =
Parameter(torch.empty(3 * embed_dim))
40         else:
41
self.register_parameter('in_proj_bias', None)
42             self.out_proj = nn.Linear(embed_dim,
embed_dim, bias=bias)
43
44             self._reset_parameters()
45
46         def _reset_parameters(self):
47             if self._qkv_same_embed_dim:
48

```

```

xavier_uniform_(self.in_proj_weight)
49         else:
50
xavier_uniform_(self.q_proj_weight)
51
xavier_uniform_(self.k_proj_weight)
52
xavier_uniform_(self.v_proj_weight)
53
54         if self.in_proj_bias is not None:
55             constant_(self.in_proj_bias, 0.)
56             constant_(self.out_proj.bias, 0.)
57
58     def forward(self, query: Tensor, key:
Tensor, value: Tensor, key_padding_mask:
Optional[Tensor] = None,
59                 need_weights: bool = True,
attn_mask: Optional[Tensor] = None) ->
Tuple[Tensor, Optional[Tensor]]:
60         if self.batch_first:
61             query, key, value =
[x.transpose(1, 0) for x in (query, key, value)]
62
63         if not self._qkv_same_embed_dim:
64             attn_output, attn_output_weights =
multi_head_attention_forward(
65                 query, key, value,
self.num_heads,
66                 self.in_proj_weight,
self.in_proj_bias,
67                 self.dropout,
self.out_proj.weight, self.out_proj.bias,
68                 training=self.training,
69
key_padding_mask=key_padding_mask,
need_weights=need_weights,
70                 attn_mask=attn_mask,
use_separate_proj_weight=True,
71
q_proj_weight=self.q_proj_weight,
k_proj_weight=self.k_proj_weight,
72
v_proj_weight=self.v_proj_weight)

```

```

73         else:
74             attn_output, attn_output_weights =
multi_head_attention_forward(
75                 query, key, value,
self.num_heads,
76                 self.in_proj_weight,
self.in_proj_bias,
77                 self.dropout,
self.out_proj.weight, self.out_proj.bias,
78                 training=self.training,
79                 key_padding_mask=key_padding_mask,
need_weights=need_weights,
80                 attn_mask=attn_mask)
81         if self.batch_first:
82             return attn_output.transpose(1,
0), attn_output_weights
83         else:
84             return attn_output,
attn_output_weights

```

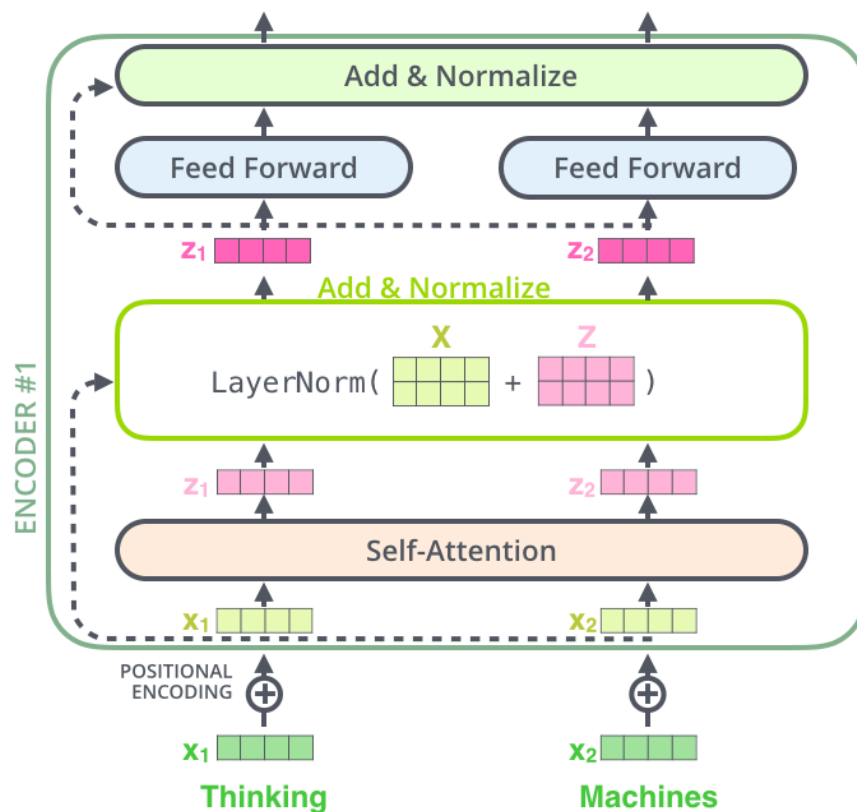
Add&Norm

残差模块就是多加了self-attention的输入，如下图所示。

```

1     x=x+z

```



LayerNorm用Pytorch的 `nn.LayerNorm` 实现:

```
1 norm = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
```

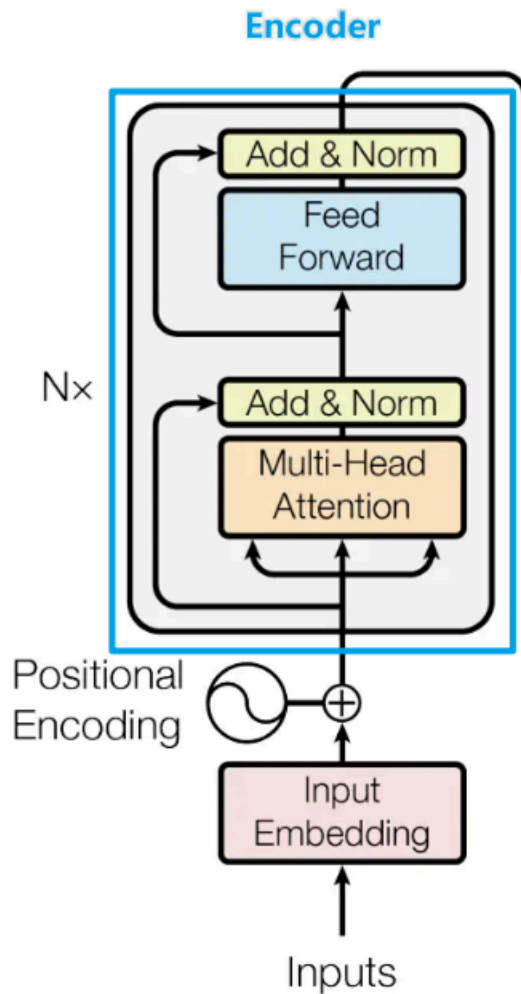
Feed forward

用Pytorch的 `nn.Linear` 实现:

```
1 linear1 = nn.Linear(input_dim, output_dim1)
2 linear2 = nn.Linear(output_dim1, output_dim2)
3 dropout=nn.Dropout(0.01)
4 activation=F.relu
5 res=linear2(dropout(activation(linear1(src))))
```

Encoder

Encoder结构如下图所示:



EncoderLayer

先写EncoderLayer根据示意图把之前将各个模块连起来。

代码

```

1  class TransformerEncoderLayer(nn.Module):
2      r'''
3      参数:
4          embed_dim: 词嵌入的维度 (必备)
5          nhead: 多头注意力中平行头的数目 (必备)
6          dim_feedforward: 全连接层的神经元的数目, 又称
            入的维度 (Default = 2048)
7          dropout: dropout的概率 (Default = 0.1)
8          activation: 两个线性层中间的激活函数, 默认relu
9          lay_norm_eps: layer normalization中的微小量, 默认
            母为0 (Default = 1e-5)
10         batch_first: 若`True`, 则为(batch, seq, feature), 若
            为`False`, 则为(seq, batch, feature) (Default: False)
11

```



```

12     例子:
13         >>> encoder_layer =
TransformerEncoderLayer(embed_dim=512, nhead=8)
14         >>> src = torch.randn((32, 10, 512))
15         >>> out = encoder_layer(src)
16     '''
17
18     def __init__(self, embed_dim, nhead,
dim_feedforward=2048, dropout=0.1, activation=F.relu,
19                 layer_norm_eps=1e-5, batch_first
-> None:
20         super(TransformerEncoderLayer, self).__i
21         self.self_attn = MultiheadAttention(embed
nhead, dropout=dropout, batch_first=batch_first)
22         self.linear1 = nn.Linear(embed_dim,
dim_feedforward)
23         self.dropout = nn.Dropout(dropout)
24         self.linear2 = nn.Linear(dim_feedforward
embed_dim)
25
26         self.norm1 = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
27         self.norm2 = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
28         self.dropout1 = nn.Dropout(dropout)
29         self.dropout2 = nn.Dropout(dropout)
30         self.activation = activation
31
32
33     def forward(self, src: Tensor, src_mask:
Optional[Tensor] = None, src_key_padding_mask:
Optional[Tensor] = None) -> Tensor:
34         src = positional_encoding(src, src.shape
置编码
35         src2 = self.self_attn(src, src, src,
attn_mask=src_mask,
36         key_padding_mask=src_key_padding_mask)[0
力, (attn_output, attn_output_weights)
37         src = src + self.dropout1(src2) #Add
38         src = self.norm1(src) #Norm
39         src2 =
self.linear2(self.dropout(self.activation(self.linear
#Feed forward

```

```

40         src = src + self.dropout(src2)#Add
41         src = self.norm2(src) #Norm
42         return src
43

```

```

1  # 用小例子看一下
2  encoder_layer =
TransformerEncoderLayer(embed_dim=512, nhead=8)
3  src = torch.randn((32, 10, 512))
4  out = encoder_layer(src)
5  print(out.shape)
6  # torch.Size([32, 10, 512])

```

EncoderLayer组成Encoder

Encoder层由多个EncoderLayer串联。

```

1  class TransformerEncoder(nn.Module):
2      r'''
3      参数:
4          encoder_layer (必备)
5          num_layers: encoder_layer的层数 (必备)
6          norm: 归一化的选择 (可选)
7
8      例子:
9          >>> encoder_layer =
TransformerEncoderLayer(embed_dim=512, nhead=8)
10         >>> transformer_encoder =
TransformerEncoder(encoder_layer, num_layers=6)
11         >>> src = torch.randn((10, 32, 512))
12         >>> out = transformer_encoder(src)
13         ...
14
15     def __init__(self, encoder_layer,
num_layers, norm=None):
16         super(TransformerEncoder,
self).__init__()
17         self.layer = encoder_layer
18         self.num_layers = num_layers
19         self.norm = norm
20

```

```

21         def forward(self, src: Tensor, mask:
Optional[Tensor] = None, src_key_padding_mask:
Optional[Tensor] = None) -> Tensor:
22             output = positional_encoding(src,
src.shape[-1])
23             for _ in range(self.num_layers):
24                 output = self.layer(output,
src_mask=mask,
src_key_padding_mask=src_key_padding_mask)
25
26             if self.norm is not None:
27                 output = self.norm(output)
28
29             return output

```

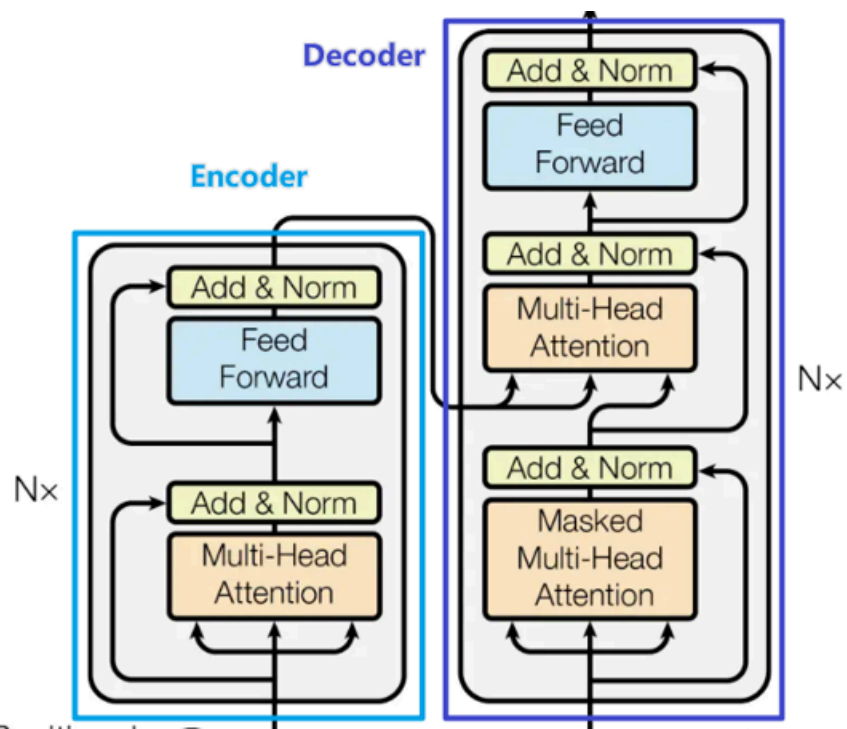
```

1  # 例子
2  encoder_layer =
TransformerEncoderLayer(embed_dim=512, nhead=8)
3  transformer_encoder =
TransformerEncoder(encoder_layer, num_layers=6)
4  src = torch.randn((10, 32, 512))
5  out = transformer_encoder(src)
6  print(out.shape)
7  # torch.Size([10, 32, 512])

```

Decoder

Decoder结构如下图紫色框所示，因为涉及到和Encoder的交互，这里把Encoder的图示也放上。



DecoderLayer

同理把各个模块串成DecoderLayer。

注意：

- 1 第一个self-attention需要输入掩膜(tgt_mask)。
- 2 第二个multihead_attn输入的Q为上一层输出tgt，K=V为Encoder最后一层输出memory。

```

1 class TransformerDecoderLayer(nn.Module):
2     r'''
3     参数:
4         embed_dim: 词嵌入的维度 (必备)
5         nhead: 多头注意力中平行头的数目 (必备)
6         dim_feedforward: 全连接层的神经元的数目, 又称
          的维度 (Default = 2048)
7         dropout: dropout的概率 (Default = 0.1)
8         activation: 两个线性层中间的激活函数, 默认relu
9         lay_norm_eps: layer normalization中的微小值
          为0 (Default = 1e-5)
10        batch_first: 若`True`, 则为(batch, seq, fe
          为`False`, 则为(seq, batch, feature) (Default: False)
11
12        例子:

```

```

13         >>> decoder_layer =
TransformerDecoderLayer(embed_dim=512, nhead=8)
14         >>> memory = torch.randn((10, 32, 512))
15         >>> tgt = torch.randn((20, 32, 512))
16         >>> out = decoder_layer(tgt, memory)
17         '''
18     def __init__(self, embed_dim, nhead,
dim_feedforward=2048, dropout=0.1, activation=F.relu,
19                 layer_norm_eps=1e-5, batch_first
> None:
20         super(TransformerDecoderLayer, self).__i
21         self.self_attn = MultiheadAttention(embed
nhead, dropout=dropout, batch_first=batch_first)
22         self.multihead_attn =
MultiheadAttention(embed_dim, nhead, dropout=dropout,
batch_first=batch_first)
23
24         self.linear1 = nn.Linear(embed_dim,
dim_feedforward)
25         self.dropout = nn.Dropout(dropout)
26         self.linear2 = nn.Linear(dim_feedforward
embed_dim)
27
28         self.norm1 = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
29         self.norm2 = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
30         self.norm3 = nn.LayerNorm(embed_dim,
eps=layer_norm_eps)
31         self.dropout1 = nn.Dropout(dropout)
32         self.dropout2 = nn.Dropout(dropout)
33         self.dropout3 = nn.Dropout(dropout)
34
35         self.activation = activation
36
37     def forward(self, tgt: Tensor, memory: Tensor
tgt_mask: Optional[Tensor] = None,
38                 memory_mask: Optional[Tensor] =
None, tgt_key_padding_mask: Optional[Tensor] = None,
memory_key_padding_mask: Optional[Tensor] = None) ->
39         r'''
40         参数:
41         tgt: 目标语言序列 (必备)

```

```

42         memory: 从最后一个encoder_layer跑出的句
43         tgt_mask: 目标语言序列的mask (可选)
44         memory_mask (可选)
45         tgt_key_padding_mask (可选)
46         memory_key_padding_mask (可选)
47         ...
48         tgt_mask = self.self_attn.
49         tgt2 = self.self_attn(tgt, tgt, tgt,
attn_mask=tgt_mask,
50
key_padding_mask=tgt_key_padding_mask)[0] #Masked自注
51         tgt = tgt + self.dropout1(tgt2) #Add
52         tgt = self.norm1(tgt) #Norm
53         tgt2 = self.multihead_attn(tgt, memory, i
attn_mask=memory_mask, key_padding_mask=memory_key_pac
[0] #cross attention, Q为上一层输出tgt, K=V为Encoder最后-
54         tgt = tgt + self.dropout2(tgt2) #Add
55         tgt = self.norm2(tgt) #Norm
56         tgt2 =
self.linear2(self.dropout(self.activation(self.linear
#Feed Forward
57         tgt = tgt + self.dropout3(tgt2) #Add
58         tgt = self.norm3(tgt) #Norm
59         return tgt

```

```

1  # 例子
2  decoder_layer =
nn.TransformerDecoderLayer(embed_dim=512, nhead=8)
3  memory = torch.randn((10, 32, 512))
4  tgt = torch.randn((20, 32, 512))
5  out = decoder_layer(tgt, memory)
6  print(out.shape)
7  # torch.Size([20, 32, 512])

```

DecoderLayer组成Decoder

```

1  class TransformerDecoder(nn.Module):
2      r'''
3      参数:
4          decoder_layer (必备)
5          num_layers: decoder_layer的层数 (必备)

```

```

6         norm: 归一化选择
7
8     例子:
9         >>> decoder_layer
=TransformerDecoderLayer(embed_dim=512, nhead=8)
10         >>> transformer_decoder =
TransformerDecoder(decoder_layer, num_layers=6)
11         >>> memory = torch.rand(10, 32, 512)
12         >>> tgt = torch.rand(20, 32, 512)
13         >>> out = transformer_decoder(tgt,
memory)
14         '''
15         def __init__(self, decoder_layer,
num_layers, norm=None):
16             super(TransformerDecoder,
self).__init__()
17             self.layer = decoder_layer
18             self.num_layers = num_layers
19             self.norm = norm
20
21         def forward(self, tgt: Tensor, memory:
Tensor, tgt_mask: Optional[Tensor] = None,
22                     memory_mask: Optional[Tensor]
= None, tgt_key_padding_mask: Optional[Tensor] =
None,
23                     memory_key_padding_mask:
Optional[Tensor] = None) -> Tensor:
24             output = tgt
25             for _ in range(self.num_layers):
26                 output = self.layer(output,
memory, tgt_mask=tgt_mask,
27
memory_mask=memory_mask,
28
tgt_key_padding_mask=tgt_key_padding_mask,
29
memory_key_padding_mask=memory_key_padding_mask)
30             if self.norm is not None:
31                 output = self.norm(output)
32
33         return output

```

```

1  # 例子
2  decoder_layer
=TransformerDecoderLayer(embed_dim=512, nhead=8)
3  transformer_decoder =
TransformerDecoder(decoder_layer, num_layers=6)
4  memory = torch.rand(10, 32, 512)
5  tgt = torch.rand(20, 32, 512)
6  out = transformer_decoder(tgt, memory)
7  print(out.shape)
8  # torch.Size([20, 32, 512])

```

Transformer

完整的Transformer模型，结合input embedding、Encoder、Decoder、最后输出的线性层和Softmax层。

代码

```

1  class Transformer(nn.Module):
2      r'''
3      参数:
4          embed_dim: 词嵌入的维度 (必备)
          (Default=512)
5          nhead: 多头注意力中平行头的数目 (必备)
          (Default=8)
6          num_encoder_layers: 编码层层数
          (Default=8)
7          num_decoder_layers: 解码层层数
          (Default=8)
8          dim_feedforward: 全连接层的神经元的数目, 又
          称经过此层输入的维度 (Default = 2048)
9          dropout: dropout的概率 (Default = 0.1)
10         activation: 两个线性层中间的激活函数, 默认
          relu或gelu
11         custom_encoder: 自定义encoder
          (Default=None)
12         custom_decoder: 自定义decoder
          (Default=None)
13         lay_norm_eps: layer normalization中的微
          小量, 防止分母为0 (Default = 1e-5)
14         batch_first: 若`True`, 则为(batch, seq,

```


feature), 若为`False`, 则为(seq, batch, feature)
(Default: False)

15

16 例子:

```
17       >>> transformer_model =  
Transformer(nhead=16, num_encoder_layers=12)  
18       >>> src = torch.rand((10, 32, 512))  
19       >>> tgt = torch.rand((20, 32, 512))  
20       >>> out = transformer_model(src, tgt)  
21       '''  
22       def __init__(self, embed_dim: int = 512,  
nhead: int = 8, num_encoder_layers: int = 6,  
23                   num_decoder_layers: int = 6,  
dim_feedforward: int = 2048, dropout: float = 0.1,  
24                   activation = F.relu,  
custom_encoder: Optional[Any] = None,  
custom_decoder: Optional[Any] = None,  
25                   layer_norm_eps: float = 1e-5,  
batch_first: bool = False) -> None:  
26           super(Transformer, self).__init__()  
27           if custom_encoder is not None:  
28               self.encoder = custom_encoder  
29           else:  
30               encoder_layer =  
TransformerEncoderLayer(embed_dim, nhead,  
dim_feedforward, dropout,  
31                   activation, layer_norm_eps, batch_first)  
32               encoder_norm =  
nn.LayerNorm(embed_dim, eps=layer_norm_eps)  
33               self.encoder =  
TransformerEncoder(encoder_layer,  
num_encoder_layers)  
34  
35           if custom_decoder is not None:  
36               self.decoder = custom_decoder  
37           else:  
38               decoder_layer =  
TransformerDecoderLayer(embed_dim, nhead,  
dim_feedforward, dropout,  
39                   activation, layer_norm_eps, batch_first)  
40               decoder_norm =
```

```

nn.LayerNorm(embed_dim, eps=layer_norm_eps)
41         self.decoder =
TransformerDecoder(decoder_layer,
num_decoder_layers, decoder_norm)
42
43         self._reset_parameters()
44
45         self.embed_dim = embed_dim
46         self.nhead = nhead
47
48         self.batch_first = batch_first
49
50     def forward(self, src: Tensor, tgt:
Tensor, src_mask: Optional[Tensor] = None,
tgt_mask: Optional[Tensor] = None,
51                 memory_mask: Optional[Tensor]
= None, src_key_padding_mask: Optional[Tensor] =
None,
52                 tgt_key_padding_mask:
Optional[Tensor] = None, memory_key_padding_mask:
Optional[Tensor] = None) -> Tensor:
53         r'''
54         参数:
55             src: 源语言序列 (送入Encoder) (必备)
56             tgt: 目标语言序列 (送入Decoder) (必备)
57             src_mask: (可选)
58             tgt_mask: (可选)
59             memory_mask: (可选)
60             src_key_padding_mask: (可选)
61             tgt_key_padding_mask: (可选)
62             memory_key_padding_mask: (可选)
63
64         形状:
65             - src: shape: `(S, N, E)`, `(N, S,
E)` if batch_first.
66             - tgt: shape: `(T, N, E)`, `(N, T,
E)` if batch_first.
67             - src_mask: shape: `(S, S)`.
68             - tgt_mask: shape: `(T, T)`.
69             - memory_mask: shape: `(T, S)`.
70             - src_key_padding_mask: shape: `(N,
S)`.
71             - tgt_key_padding_mask: shape: `(N,

```

```

T)` .
72             - memory_key_padding_mask:
shape: `(N, S)` .
73
74             [src/tgt/memory]_mask确保有些位置不被
看到，如做decode的时候，只能看该位置及其以前的，而不能看后面
的。
75             若为ByteTensor，非0的位置会被忽略不做注
意力；若为BoolTensor，True对应的位置会被忽略；
76             若为数值，则会直接加到attn_weights
77
78             [src/tgt/memory]_key_padding_mask
使得key里面的某些元素不参与attention计算，三种情况同上
79
80             - output: shape: `(T, N, E)` , `(N,
T, E)` if batch_first.
81
82             注意：
83             src和tgt的最后一维需要等于embed_dim,
batch的那一维需要相等
84
85             例子：
86             >>> output =
transformer_model(src, tgt, src_mask=src_mask,
tgt_mask=tgt_mask)
87             ...
88             memory = self.encoder(src,
mask=src_mask,
src_key_padding_mask=src_key_padding_mask)
89             output = self.decoder(tgt, memory,
tgt_mask=tgt_mask, memory_mask=memory_mask,
90
tgt_key_padding_mask=tgt_key_padding_mask,
91
memory_key_padding_mask=memory_key_padding_mask)
92
93             return output
94
95             def generate_square_subsequent_mask(self,
sz: int) -> Tensor:
96             r'''产生关于序列的mask，被遮住的区域赋值`-
inf`，未被遮住的区域赋值为`0`'''
97             mask = (torch.triu(torch.ones(sz, sz))

```

```

== 1).transpose(0, 1)
98         mask = mask.float().masked_fill(mask
== 0, float('-inf')).masked_fill(mask == 1,
float(0.0))
99         return mask
100
101     def _reset_parameters(self):
102         r'''用正态分布初始化参数'''
103         for p in self.parameters():
104             if p.dim() > 1:
105                 xavier_uniform_(p)

1  # 例子
2  transformer_model = Transformer(nhead=16, num_en
3  src = torch.rand((10, 32, 512)) # (seq_len, batch
4  tgt = torch.rand((20, 32, 512)) # (seq_len, batch
5  (seq_len, batch_size, embed_dim)=(20, 32, 512)
6  # 这里假设所有句子都不需要padding,
src_key_padding_mask=None, tgt_key_padding_mask=None
7  # 生成tgt_mask
8
tgt_mask=transformer_model.generate_square_subsequent
9  out = transformer_model(src, tgt, tgt_mask=tgt_m
10  print(out.shape)
11  # torch.Size([20, 32, 512])

```

基于Transformer的模型

上面的代码实际上只写到了Decoder输出结束，还未根据下游任务继续设计网络，以文本翻译任务为例的模型其实就和上面的图一样，在Transformer的Decoder输出后接上一个线性层和一个softmax层。

```

1  class MyModel(nn.Module):
2      def __init__(self, transformer_layer, output_
3          super(MyModel, self).__init__()
4          self.transformer_layer = transformer_lay
5
6          #这里根据下游任务继续设计层，以文本翻译为例，这里
7          self.linear = nn.Linear(self.transfor
8

```

```

9         def forward(self, src: Tensor, tgt: Tensor,
10                       memory_mask: Optional[Tensor] = None,
11                       tgt_key_padding_mask: Optional[Tensor] = None,
12                       output_mask: Optional[Tensor] = None):
13             output=self.transformer_layer(src=src, tgt=tgt,
14             src_mask=src_mask, tgt_mask=tgt_mask, memory_mask=memory_mask,
15             output_mask=output_mask)
16             output=F.softmax(self.linear(output))
17             return output

```

```

1  vocab_size = 154 # 这里是随便举的字典大小的例子
2  transformer_layer = Transformer(embed_dim, nhead, dropout, activation_fn)
3
4  model=MyModel(transformer_layer,output_dim = vocab_size)
5  src = torch.rand((10, 32, 512)) # (seq_len,batch_size,embed_dim)
6  tgt = torch.rand((20, 32, 512)) # (seq_len,batch_size,embed_dim)
7  (seq_len,batch_size,embed_dim)=(20, 32, 512)
8  # 这里假设所有句子都不需要padding, src_key_padding_mask=None
9  # 生成tgt_mask
10
11 tgt_mask=MyModel.transformer_layer.generate_square_subsequent_mask(tgt_seq_len)
12 out = model(src, tgt, tgt_mask=tgt_mask)
13 print(out.shape)
14 # torch.Size([20, 32, 154])
15 # 然后取概率分布最大的值对应索引作为输出
16 res=torch.max(out,2)[1]
17 print(res.shape)
18 # torch.Size([20, 32])

```

参考文献

Pytorch编写完整的Transformer

Transformer相关——（8）Transformer模型

Transformer相关——（7）Mask机制

Transformer相关——（6）Normalization方式

Transformer相关——（5）残差模块

Transformer相关——（4）Position encoding

Transformer相关——（3）Attention机制