

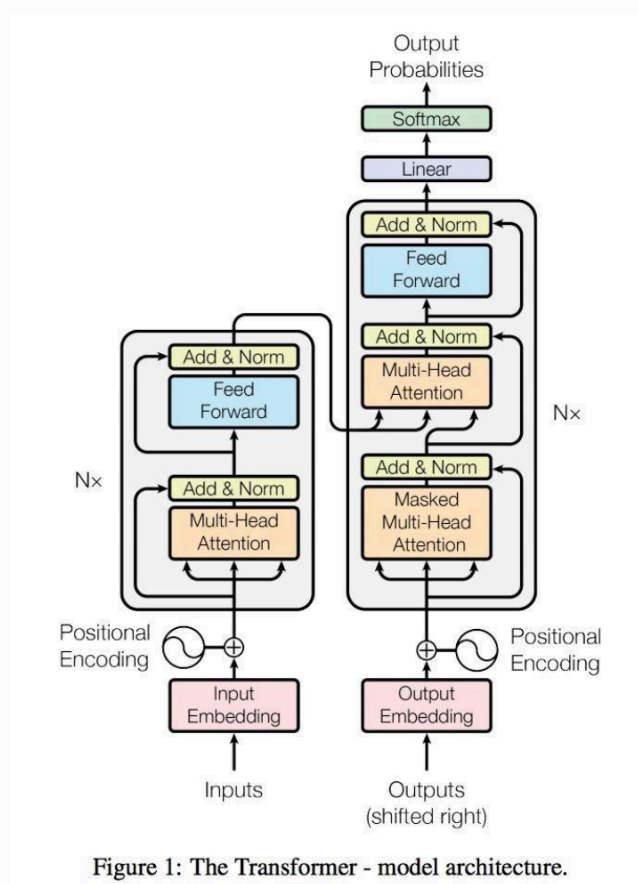
Transformer

特点

- CNN 编码局部信息，没有全局上下文
- RNN编码全局信息，不能并行
- Transformer能并行，编码全局信息

图的左半边用 $N \times$ 框出来的，就是我们的encoder的一层。encoder一共有6层这样的结构。

图的右半边用 $N \times$ 框出来的，就是我们的decoder的一层。decoder一共有6层这样的结构。



Encoder

encoder由6层相同的层组成，每一层分别由两部分组成：

- 第一部分是一个**multi-head self-attention mechanism**
- 第二部分是一个**position-wise feed-forward network**，是一个全连接层

两个部分，都有一个 **残差连接(residual connection)**，然后接着一个**Layer Normalization**。

Decoder

decoder由6个相同的层组成，每一个层包括以下3个部分：

- 第一个部分是**multi-head self-attention mechanism**
- 第二部分是**multi-head context-attention mechanism**
- 第三部分是一个**position-wise feed-forward network**

还是和encoder类似，上面三个部分的每一个部分，都有一个**残差连接**，后接一个**Layer Normalization**。

Self Attention

到目前为止，对Attention层的描述都是一般化的，我们可以落实一些应用。比如，如果做阅读理解的话，Q可以是篇章的词向量序列，取K=V为问题的词向量序列，那么输出就是所谓的Aligned Question Embedding。

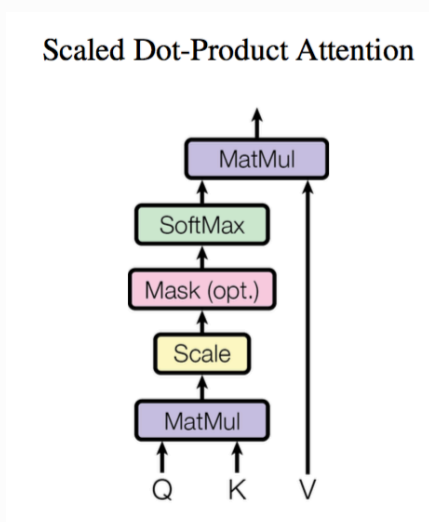
而在Google的论文中，大部分的Attention都是Self Attention，即“自注意力”，或者叫内部注意力。

所谓Self Attention，其实就是 $Attention(X, X, X)$ ，X就是前面说的输入序列。也就是说，在序列内部做Attention，寻找序列内部的联系。Google论文的主要贡献之一是它表明了**内部注意力在机器翻译（甚至是一般的Seq2Seq任务）的序列编码上是相当重要的，而之前关于Seq2Seq的研究基本都只是把注意力机制用在解码端**。类似的事情是，目前SQUAD阅读理解的榜首模型R-Net也加入了自注意力机制，这也使得它的模型有所提升。

Context-attention

它是**encoder和decoder之间的attention**！所以，你也可以称之为**encoder-decoder attention**！

Attention定义



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

其中 $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$, $\mathbf{K} \in \mathbb{R}^{m \times d_k}$, $\mathbf{V} \in \mathbb{R}^{m \times d_v}$, 于是我们可以认为：这是一个Attention层，将 $n \times d_k$ 的序列Q编码成了一个新的 $n \times d_v$ 的序列。

首先说明一下我们的K、Q、V是什么：

- 在encoder的self-attention中，Q、K、V都来自同一个地方（相等），他们是上一层encoder的输出。对于第一层encoder，它们就是word embedding和positional encoding相加得到的输入。
- 在decoder的self-attention中，Q、K、V都来自于同一个地方（相等），它们是上一层decoder的输出。对于第一层decoder，它们就是word embedding和positional encoding相加得到的输入。但是对于decoder，我们不希望它能获得下一个time step（即将来的信息），因此我们需要进行**sequence masking**。
- 在encoder-decoder attention中，Q来自于decoder的上一层的输出，K和V来自于encoder的输出，K和V是一样的。
- Q、K、V三者的维度一样，即 $d_q = d_k = d_v$ 。

Scaled dot-product attention的实现

```
1 class ScaledDotProductAttention(nn.Module):
2     """Scaled dot-product attention mechanism."""
3
4     def __init__(self, attention_dropout=0.0):
5         super(ScaledDotProductAttention, self).__init__()
6         self.dropout = nn.Dropout(attention_dropout)
7         self.softmax = nn.Softmax(dim=2)
8
9     def forward(self, q, k, v, scale=None, attn_mask=None):
10        """前向传播.
11
12        Args:
13            q: Queries张量, 形状为[B, L_q, D_q]
14            k: Keys张量, 形状为[B, L_k, D_k]
15            v: Values张量, 形状为[B, L_v, D_v], 一般来说就是k
16            scale: 缩放因子, 一个浮点标量
17            attn_mask: Masking张量, 形状为[B, L_q, L_k]
18
19        Returns:
20            上下文张量和attention张量
21        """
22        attention = torch.bmm(q, k.transpose(1, 2))
23        if scale:
24            attention = attention * scale
25        if attn_mask:
26            # 给需要mask的地方设置一个负无穷
27            attention = attention.masked_fill_(attn_mask, -np.inf)
28        # 计算softmax
29        attention = self.softmax(attention)
30        # 添加dropout
31        attention = self.dropout(attention)
32        # 和v做点积
33        context = torch.bmm(attention, v)
```

```

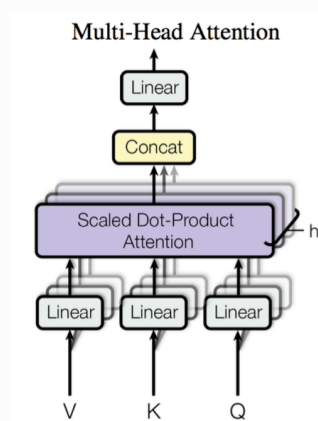
34         return context, attention
35

```

Multi-Head Attention

论文提到，他们发现将 Q 、 K 、 V 通过一个线性映射之后，分成 h 份，对每一份进行 **scaled dot-product attention** 效果更好。然后，把各个部分的结果合并起来，再次经过线性映射，得到最终的输出。这就是所谓的 **multi-head attention**。上面的超参数 h 就是 **heads** 数量。论文默认是 8。

值得注意的是，上面所说的分成 h 份是在 d_k 、 d_q 、 d_v 维度上面进行切分的。因此，进入到 scaled dot-product attention 的 d_k 实际上等于未进入之前的 D_K/h 。



就是把 Q, K, V 通过参数矩阵映射一下，然后再做 Attention，把这个过程重复做 h 次，结果拼接起来就行了，可谓“大道至简”了。

所谓“多头”（Multi-Head），就是只多做几次同样的事情（参数不共享），然后把结果拼接。

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)$$

论文里面， $d_{model} = 512$ ， $h = 8$ 。所以在 scaled dot-product attention 里面的

$$d_q = d_k = d_v = d_{model}/h = 512/8 = 64$$

Multi-head attention的实现

```

1  class MultiHeadAttention(nn.Module):
2
3      def __init__(self, model_dim=512, num_heads=8, dropout=0.0):
4          super(MultiHeadAttention, self).__init__()
5
6          self.dim_per_head = model_dim // num_heads
7          self.num_heads = num_heads
8          self.linear_k = nn.Linear(model_dim, self.dim_per_head *
num_heads)

```

```

9         self.linear_v = nn.Linear(model_dim, self.dim_per_head *
num_heads)
10        self.linear_q = nn.Linear(model_dim, self.dim_per_head *
num_heads)
11
12        self.dot_product_attention =
ScaledDotProductAttention(dropout)
13        self.linear_final = nn.Linear(model_dim, model_dim)
14        self.dropout = nn.Dropout(dropout)
15        # multi-head attention之后需要做layer norm
16        self.layer_norm = nn.LayerNorm(model_dim)
17
18        def forward(self, key, value, query, attn_mask=None):
19            # 残差连接
20            residual = query
21
22            dim_per_head = self.dim_per_head
23            num_heads = self.num_heads
24            batch_size = key.size(0)
25
26            # linear projection
27            key = self.linear_k(key)
28            value = self.linear_v(value)
29            query = self.linear_q(query)
30
31            # split by heads
32            key = key.view(batch_size * num_heads, -1, dim_per_head)
33            value = value.view(batch_size * num_heads, -1,
dim_per_head)
34            query = query.view(batch_size * num_heads, -1,
dim_per_head)
35
36            if attn_mask:
37                attn_mask = attn_mask.repeat(num_heads, 1, 1)
38            # scaled dot product attention
39            scale = (key.size(-1) // num_heads) ** -0.5
40            context, attention = self.dot_product_attention(
41                query, key, value, scale, attn_mask)
42            # concat heads
43            context = context.view(batch_size, -1, dim_per_head *
num_heads)
44            # final linear projection
45            output = self.linear_final(context)
46            # dropout
47            output = self.dropout(output)
48            # add residual and norm layer: layer_norm(wx+F(x))
49            output = self.layer_norm(residual + output)
50            return output, attention
51
52        # Layer normalization: 在x的最后一个维度求方差，最后一个维度就是隐层的维度
53        # Normalization有很多种，但是它们都有一个共同的目的，那就是把输入转化成均值为
0方差为1的数据。我们在把数据送入激活函数之前进行normalization（归一化），因为
我们不希望输入数据落在激活函数的饱和区。

```

Mask

Transformer模型里面涉及两种mask。分别是padding mask和sequence mask。其中后者我们已经在decoder的self-attention里面见过啦！

其中，padding mask在所有的scaled dot-product attention里面都需要用到，而sequence mask只有在decoder的self-attention里面用到。

所以，我们之前ScaledDotProductAttention的forward方法里面的参数attn_mask在不同的地方会有不同的含义。这一点我们会在后面说明。

1. Padding mask, 常见
2. Sequence mask

sequence mask是为了使得decoder不能看见未来的信息。也就是对于一个序列，在time_step为t的时刻，我们的解码输出应该只能依赖于t时刻之前的输出，而不能依赖t之后的输出。

那么具体怎么做呢？也很简单：产生一个上三角矩阵，上三角的值全为1，下三角的值全为0，对角线也是0。把这个矩阵作用在每一个序列上，就可以达到我们的目的啦。

attn_mask 参数有几种情况？分别是什么意思？

- 对于decoder的self-attention，里面使用到的scaled dot-product attention，同时需要padding mask 和 sequence mask 作为 attn_mask，具体实现就是两个mask相加作为attn_mask。
- 其他情况，attn_mask 一律等于padding mask。

Position Embedding

论文的实现很有意思，使用正余弦函数。公式如下：

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

公式中的 d_{model} 是模型的维度，论文默认是512，pos是指词语在序列中的位置。可以看出，在偶数位置，使用正弦编码，在奇数位置，使用余弦编码。

Position-wise Feed-Forward network

这就是一个全连接网络，包含两个线性变换和一个非线性函数（实际上就是ReLU）。公式如下：

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

这个线性变换在不同的位置都表现地一样，并且在不同的层之间使用不同的参数。

论文提到，这个公式还可以用两个核大小为1的一维卷积来解释，卷积的输入输出都是 $d_{model} = 512$ ，中间层的维度是 $d_{ff} = 2048$ 。

Transformer的实现

```
1 import torch
2 import torch.nn as nn
3
4
5 class EncoderLayer(nn.Module):
6     """Encoder的一层。"""
7
8     def __init__(self, model_dim=512, num_heads=8, ffn_dim=2018,
9 dropout=0.0):
10         super(EncoderLayer, self).__init__()
11
12         self.attention = MultiHeadAttention(model_dim, num_heads,
13 dropout)
14         self.feed_forward = PositionalWiseFeedForward(model_dim,
15 ffn_dim, dropout)
16
17     def forward(self, inputs, attn_mask=None):
18
19         # self attention
20         context, attention = self.attention(inputs, inputs,
21 inputs, padding_mask)
22
23         # feed forward network
24         output = self.feed_forward(context)
25
26         return output, attention
27
28
29 class Encoder(nn.Module):
30     """多层EncoderLayer组成Encoder。"""
31
32     def __init__(self,
33 vocab_size,
34 max_seq_len,
35 num_layers=6,
36 model_dim=512,
37 num_heads=8,
38 ffn_dim=2048,
39 dropout=0.0):
40         super(Encoder, self).__init__()
41
42         self.encoder_layers = nn.ModuleList(
43             [EncoderLayer(model_dim, num_heads, ffn_dim, dropout)
44 for _ in
45 range(num_layers)])
46
47         self.seq_embedding = nn.Embedding(vocab_size + 1,
48 model_dim, padding_idx=0)
```

```

43         self.pos_embedding = PositionalEncoding(model_dim,
max_seq_len)
44
45     def forward(self, inputs, inputs_len):
46         output = self.seq_embedding(inputs)
47         output += self.pos_embedding(inputs_len)
48
49         self_attention_mask = padding_mask(inputs, inputs)
50
51         attentions = []
52         for encoder in self.encoder_layers:
53             output, attention = encoder(output,
self_attention_mask)
54             attentions.append(attention)
55
56         return output, attentions
57
58
59 class DecoderLayer(nn.Module):
60
61     def __init__(self, model_dim, num_heads=8, ffn_dim=2048,
dropout=0.0):
62         super(DecoderLayer, self).__init__()
63
64         self.attention = MultiHeadAttention(model_dim, num_heads,
dropout)
65         self.feed_forward = PositionalWiseFeedForward(model_dim,
ffn_dim, dropout)
66
67     def forward(self,
68                 dec_inputs,
69                 enc_outputs,
70                 self_attn_mask=None,
71                 context_attn_mask=None):
72         # self attention, all inputs are decoder inputs
73         dec_output, self_attention = self.attention(
74             dec_inputs, dec_inputs, dec_inputs, self_attn_mask)
75
76         # context attention
77         # query is decoder's outputs, key and value are encoder's
inputs
78         dec_output, context_attention = self.attention(
79             enc_outputs, enc_outputs, dec_output, context_attn_mask)
80
81         # decoder's output, or context
82         dec_output = self.feed_forward(dec_output)
83
84         return dec_output, self_attention, context_attention
85
86
87 class Decoder(nn.Module):
88

```



```

89         def __init__(self,
90                       vocab_size,
91                       max_seq_len,
92                       num_layers=6,
93                       model_dim=512,
94                       num_heads=8,
95                       ffn_dim=2048,
96                       dropout=0.0):
97             super(Decoder, self).__init__()
98
99             self.num_layers = num_layers
100
101             self.decoder_layers = nn.ModuleList(
102                 [DecoderLayer(model_dim, num_heads, ffn_dim, dropout)
103                  for _ in
104                      range(num_layers)])
105
106             self.seq_embedding = nn.Embedding(vocab_size + 1,
107                                               model_dim, padding_idx=0)
108             self.pos_embedding = PositionalEncoding(model_dim,
109                                                       max_seq_len)
110
111             def forward(self, inputs, inputs_len, enc_output,
112                         context_attn_mask=None):
113                 output = self.seq_embedding(inputs)
114                 output += self.pos_embedding(inputs_len)
115
116                 self_attention_padding_mask = padding_mask(inputs, inputs)
117                 seq_mask = sequence_mask(inputs)
118                 self_attn_mask = torch.gt((self_attention_padding_mask +
119                                             seq_mask), 0)
120
121                 self_attentions = []
122                 context_attentions = []
123                 for decoder in self.decoder_layers:
124                     output, self_attn, context_attn = decoder(
125                         output, enc_output, self_attn_mask, context_attn_mask)
126                     self_attentions.append(self_attn)
127                     context_attentions.append(context_attn)
128
129                 return output, self_attentions, context_attentions
130
131 class Transformer(nn.Module):
132
133     def __init__(self,
134                 src_vocab_size,
135                 src_max_len,
136                 tgt_vocab_size,
137                 tgt_max_len,
138                 num_layers=6,
139                 model_dim=512,

```

```

136         num_heads=8,
137         ffn_dim=2048,
138         dropout=0.2):
139     super(Transformer, self).__init__()
140
141     self.encoder = Encoder(src_vocab_size, src_max_len,
num_layers, model_dim,
142                             num_heads, ffn_dim, dropout)
143     self.decoder = Decoder(tgt_vocab_size, tgt_max_len,
num_layers, model_dim,
144                             num_heads, ffn_dim, dropout)
145
146     self.linear = nn.Linear(model_dim, tgt_vocab_size,
bias=False)
147     self.softmax = nn.Softmax(dim=2)
148
149     def forward(self, src_seq, src_len, tgt_seq, tgt_len):
150         context_attn_mask = padding_mask(tgt_seq, src_seq)
151
152         output, enc_self_attn = self.encoder(src_seq, src_len)
153
154         output, dec_self_attn, ctx_attn = self.decoder(
155             tgt_seq, tgt_len, output, context_attn_mask)
156
157         output = self.linear(output)
158         output = self.softmax(output)
159
160         return output, enc_self_attn, dec_self_attn, ctx_attn
161
162

```