

# 一文读懂「Attention is All You Need」| 附代码实现

原创：苏剑林 PaperWeekly 1月10日

作者 | 苏剑林

单位 | 广州火焰信息科技有限公司

研究方向 | NLP, 神经网络

个人主页 | kexue.fm

## 前言

2017 年中，有两篇类似同时也是笔者非常欣赏的论文，分别是 FaceBook 的 **Convolutional Sequence to Sequence Learning** 和 Google 的 **Attention is All You Need**，它们都算是 Seq2Seq 上的创新，本质上来说，都是抛弃了 RNN 结构来做 Seq2Seq 任务。

在本篇文章中，笔者将对 **Attention is All You Need** 做一点简单的分析。当然，这两篇论文本身就比较火，因此网上已经有很多解读了（不过很多解读都是直接翻译论文的，鲜有自己的理解），因此这里尽可能多自己的文字，尽量不重复网上各位大佬已经说过的内容。

## 序列编码

深度学习做 NLP 的方法，基本上都是先将句子分词，然后每个词转化为对应的词向量序列。这样一来，每个句子都对应的是一个矩阵  $X=(x_1, x_2, \dots, x_t)$ ，其中  $x_i$  都代表着第  $i$  个词的词向量（行向量），维度为  $d$  维，故  $X \in \mathbb{R}^{n \times d}$ 。这样的话，问题就变成了编码这些序列了。

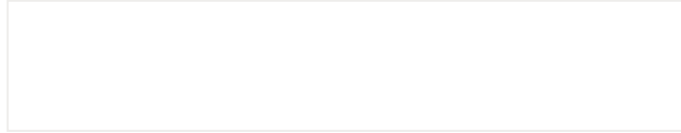
第一个基本的思路是 RNN 层，RNN 的方案很简单，递归式进行：

$$y_t = f(y_{t-1}, x_t)$$

不管是已经被广泛使用的 LSTM、GRU 还是最近的 SRU，都并未脱离这个递归框架。RNN 结构本身比较简单，也很适合序列建模，但 RNN 的明显缺点之一就是无法并行，因此速度较慢，这是递归的天然缺陷。

另外我个人觉得 **RNN 无法很好地学习到全局的结构信息，因为它本质是一个马尔科夫决策过程。**

**第二个思路是 CNN 层**，其实 CNN 的方案也是很自然的，窗口式遍历，比如尺寸为 3 的卷积，就是：

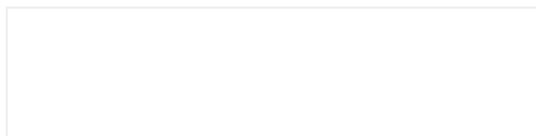


在 FaceBook 的论文中，纯粹使用卷积也完成了 Seq2Seq 的学习，是卷积的一个精致且极致的使用案例，热衷卷积的读者必须得好好读读这篇文论。

**CNN 方便并行，而且容易捕捉到一些全局的结构信息，笔者本身是比较偏爱 CNN 的，在目前的工作或竞赛模型中，我都已经尽量用 CNN 来代替已有的 RNN 模型了，并形成了自己的一套使用经验，这部分我们以后再谈。**

Google 的大作提供了**第三个思路：纯 Attention，单靠注意力就可以。**

RNN 要逐步递归才能获得全局信息，因此一般要双向 RNN 才比较好；CNN 事实上只能获取局部信息，是通过层叠来增大感受野；Attention 的思路最为粗暴，它一步到位获取了全局信息，它的解决方案是：

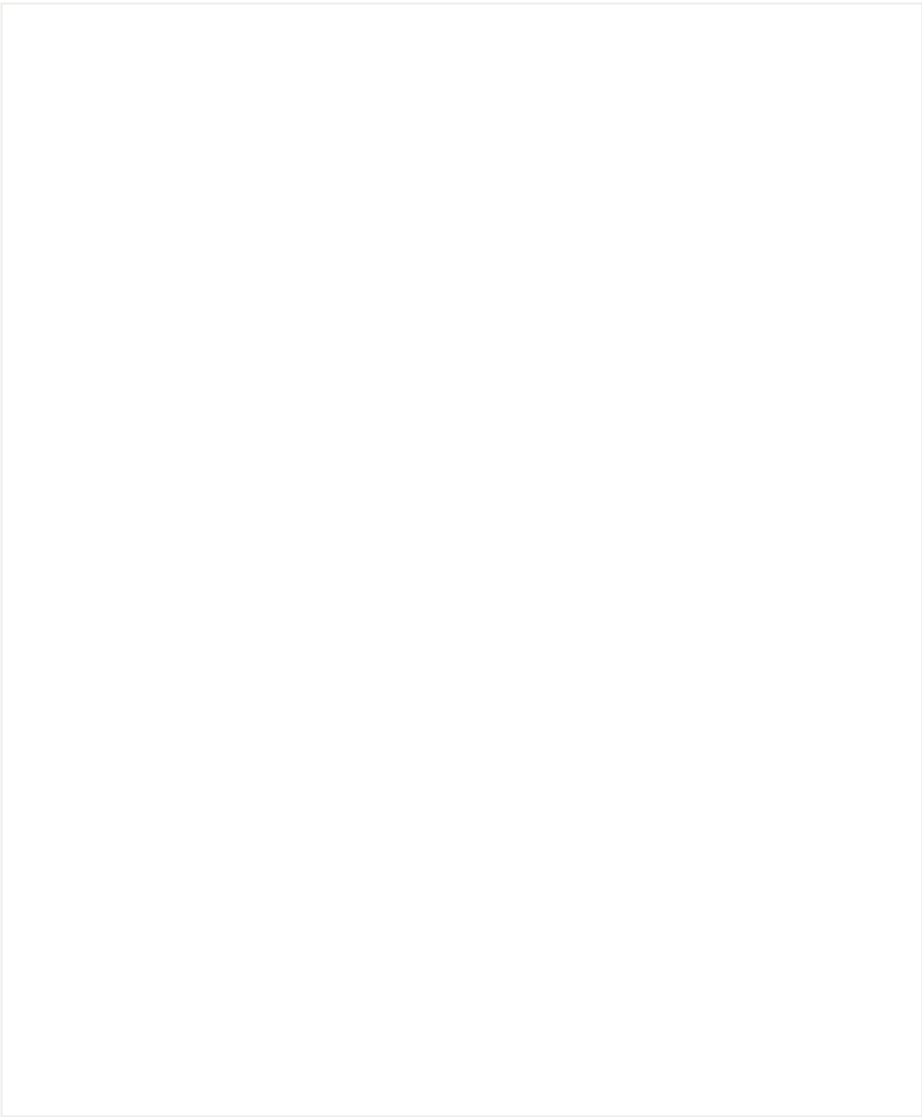


其中  $A, B$  是另外一个序列（矩阵）。如果都取  $A=B=X$ ，那么就称为 Self Attention，它的意思是直接将  $x_t$  与原来的每个词进行比较，最后算出  $y_t$ 。

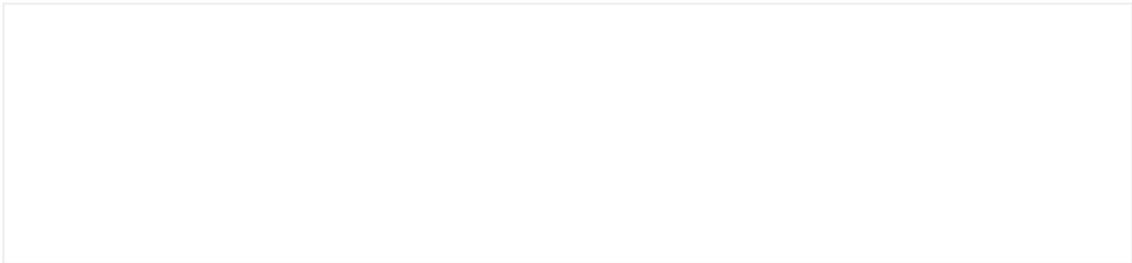
## Attention 层

### Attention 定义

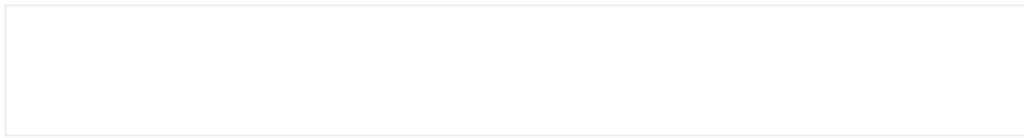
Google 的一般化 Attention 思路也是一个编码序列的方案，因此我们也可以认为它跟 RNN、CNN 一样，都是一个序列编码的层。



前面给出的是一般化的框架形式的描述，事实上 Google 给出的方案是很具体的。首先，它先把 Attention 的定义给了出来：



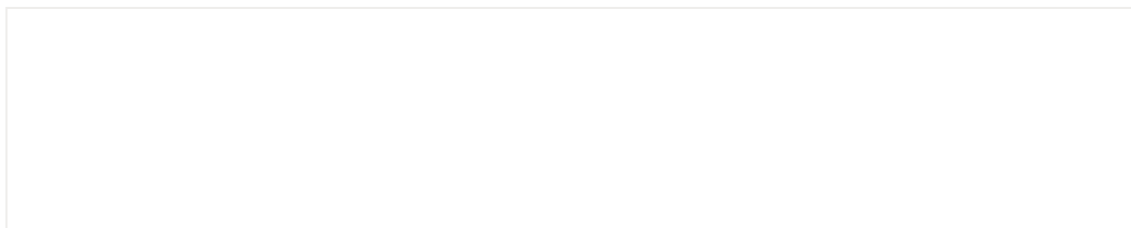
这里用的是跟 Google 的论文一致的符号，其中：



如果忽略激活函数 softmax 的话，那么事实上它就是三个  $n \times dk, dk \times m, m \times dv$  的矩阵相乘，最后的结果就是一个  $n \times dv$  的矩阵。

于是我们可以认为：这是一个 Attention 层，将  $n \times dk$  的序列 Q 编码成了一个新的  $n \times dv$  的序列。

那怎么理解这种结构呢？我们不妨逐个向量来看。



其中  $Z$  是归一化因子。事实上  $q, k, v$  分别是 query, key, value 的简写， $K, V$  是一一对应的，它们就像是 key-value 的关系，那么上式的意思就是通过  $q_t$  这个 query，通过与各个  $k_s$  内积的并 softmax 的方式，来得到  $q_t$  与各个  $v_s$  的相似度，然后加权求和，得到一个  $dv$  维的向量。

其中因子  $\frac{1}{Z}$  起到调节作用，使得内积不至于太大（太大的话 softmax 后就非 0 即 1 了，不够“soft”了）。

事实上这种 Attention 的定义并不新鲜，但由于 Google 的影响力，我们可以认为现在是更加正式地提出了这个定义，并将其视为一个层地看待。

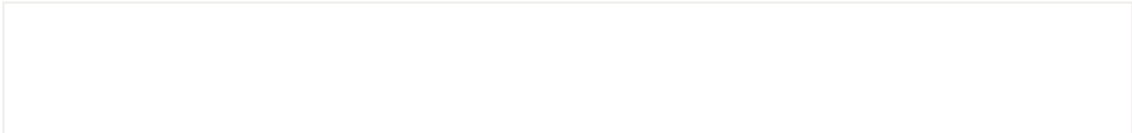
此外这个定义只是注意力的一种形式，还有一些其他选择，比如 query 跟 key 的运算方式不一定是点乘（还可以是拼接后再内积一个参数向量），甚至权重都不一定要归一化，等等。

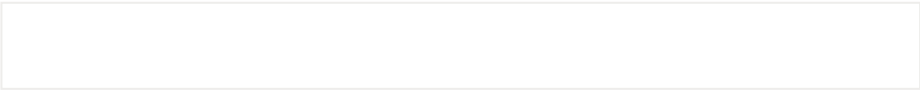
## Multi-Head Attention

这个是 Google 提出的新概念，是 Attention 机制的完善。



不过从形式上看，它其实就再简单不过了，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。具体来说：



这里 ，然后：



最后得到一个  $n \times (h \cdot d_v)$  的序列。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。

## 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  $\text{Attention}(X, X, X)$ ， $X$  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。

当然，更准确来说，Google 所用的是 Self Multi-Head Attention：



## Position Embedding

然而，只要稍微思考一下就会发现，这样的模型并不能捕捉序列的顺序。换句话说，如果将  $K, V$  按行打乱顺序（相当于句子中的词序打乱），那么 Attention 的结果还是一样的。

这就表明了，**到目前为止，Attention 模型顶多是一个非常精妙的“词袋模型”而已。**

这问题就比较严重了，大家知道，对于时间序列来说，尤其是对于 NLP 中的任务来说，顺序是很重要的信息，它代表着局部甚至是全局的结构，学习不到顺序信息，那么效果将会大打折扣（比

如机器翻译中，有可能只把每个词都翻译出来了，但是不能组织成合理的句子）。

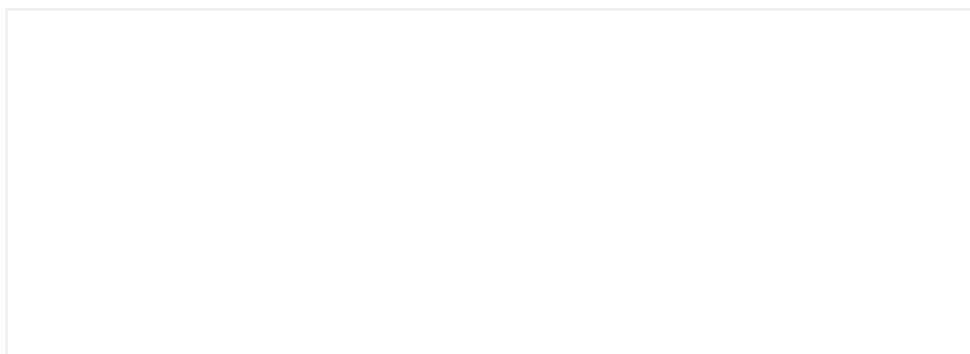
于是 Google 再祭出了一招——**Position Embedding**，也就是“位置向量”，将每个位置编号，然后每个编号对应一个向量，通过结合位置向量和词向量，就给每个词都引入了一定的位置信息，这样 Attention 就可以分辨出不同位置的词了。

Position Embedding 并不算新鲜的玩意，在 FaceBook 的 **Convolutional Sequence to Sequence Learning** 也用到了这个东西。但在 Google 的这个作品中，它的 Position Embedding 有几点区别：

1. 以前在 RNN、CNN 模型中其实都出现过 Position Embedding，但在那些模型中，Position Embedding 是锦上添花的辅助手段，也就是“有它会更好、没它也就差一点点”的情况，因为 RNN、CNN 本身就能捕捉到位置信息。

但是在这个纯 Attention 模型中，Position Embedding 是位置信息的唯一来源，因此它是模型的核心成分之一，并非仅仅是简单的辅助手段。

2. 在以往的 Position Embedding 中，基本都是根据任务训练出来的向量。而 Google 直接给出了一个构造 Position Embedding 的公式：



这里的意思是将 id 为  $p$  的位置映射为一个  $d_{pos}$  维的位置向量，这个向量的第  $i$  个元素的数值就是  $PE_i(p)$ 。

Google 在论文中说到他们比较过直接训练出来的位置向量和上述公式计算出来的位置向量，效果是接近的。因此显然我们更乐意使用公式构造的 Position Embedding 了。

3. Position Embedding 本身是一个绝对位置的信息，但在语言中，相对位置也很重要，Google 选择前述的位置向量公式的一个重要原因如下：

由于我们有  $\sin(\alpha+\beta)=\sin\alpha \cos\beta+\cos\alpha \sin\beta$  以及  $\cos(\alpha+\beta)=\cos\alpha \cos\beta-\sin\alpha \sin\beta$ ，这表明位置  $p+k$  的向量可以表明位置  $p$  的向量的线性变换，这提供了表达相对位置信息的可能性。

结合位置向量和词向量有几个可选方案，**可以把它们拼接起来作为一个新向量，也可以把位置向量定义为跟词向量一样大小，然后两者加起来。**

FaceBook 的论文用的是前者，而 Google 论文中用的是后者。直觉上相加会导致信息损失，似乎不可取，但 Google 的成果说明相加也是很好的方案。看来我理解还不够深刻。

## 一些不足之处

到这里，Attention 机制已经基本介绍完了。**Attention 层的好处是能够一步到位捕捉到全局的联系，因为它直接把序列两两比较（代价是计算量变为  $O(n^2)$ ，当然由于是纯矩阵运算，这个计算量相当也不是很严重）。**

**相比之下，RNN 需要一步步递推才能捕捉到，而 CNN 则需要通过层叠来扩大感受野，这是 Attention 层的明显优势。**

Google 论文剩下的工作，就是介绍它怎么用到机器翻译中，这是个应用和调参的问题，我们这里不特别关心它。当然，Google 的结果表明将纯注意力机制用在机器翻译中，能取得目前最好的效果，这结果的确是辉煌的。

然而，我还是想谈谈这篇论文本身和 Attention 层自身的一些不足的地方。

1. 论文标题为 **Attention is All You Need**，因此论文中刻意避免出现了 RNN、CNN 的字眼，但我觉得这种做法过于刻意了。

事实上，论文还专门命名了一种 Position-wise Feed-Forward Networks，事实上它就是窗口大小为 1 的一维卷积，因此有种为了不提卷积还专门换了个名称的感觉，有点不厚道。（也有可能我过于臆测了）。



2. Attention 虽然跟 CNN 没有直接联系，但事实上充分借鉴了 CNN 的思想，比如 Multi-Head Attention 就是 Attention 做多次然后拼接，这跟 CNN 中的多个卷积核的思想是一致的；还有论文用到了残差结构，这也源于 CNN 网络。

3. 无法对位置信息进行很好地建模，这是硬伤。尽管可以引入 Position Embedding，但我认为这只是一个缓解方案，并没有根本解决问题。

举个例子，用这种纯 Attention 机制训练一个文本分类模型或者是机器翻译模型，效果应该都还不错，但是用来训练一个序列标注模型（分词、实体识别等），效果就不怎么好了。

那为什么在机器翻译任务上好？我觉得原因是机器翻译这个任务并不特别强调语序，因此 Position Embedding 所带来的位置信息已经足够了，此外翻译任务的评测指标 BLEU 也并不特别强调语序。

4. 并非所有问题都需要长程的、全局的依赖的，也有很多问题只依赖于局部结构，这时候用纯 Attention 也不大好。

事实上，Google 似乎也意识到了这个问题，因此论文中也提到了一个 restricted 版的 Self-Attention（不过论文正文应该没有用到它）。

它假设当前词只与前后  $r$  个词发生联系，因此注意力也只发生在这  $2r+1$  个词之间，这样计算量就是  $O(nr)$ ，这样也能捕捉到序列的局部结构了。但是很明显，这就是卷积核中的卷积窗口的概念。

通过以上讨论，我们可以体会到，把 Attention 作为一个单独的层来看，跟 CNN、RNN 等结构混合使用，应该能更充分融合它们各自的优势，而不必像 Google 论文号称 Attention is All You Need，那样实在有点“矫枉过正”了（“口气”太大），事实上也做不到。

就论文的工作而言，也许降低一下身段，称为 Attention is All Seq2Seq Need（事实上也这标题的“口气”也很大），会获得更多的肯定。

## 代码实现

最后，为了使得本文有点实用价值，笔者试着给出了论文的 Multi-Head Attention 的实现代码。有需要的读者可以直接使用，或者参考着修改。

注意的是，Multi-Head 的意思虽然很简单——重复做几次然后拼接，但事实上不能按照这个思路来写程序，这样会非常慢。因为 TensorFlow 是不会自动并行的，比如：

```
a = tf.zeros((10, 10))
b = a + 1
c = a + 2
```

其中 b,c 的计算是串联的，尽管 b,c 没有相互依赖。因此我们必须把 Multi-Head 的操作合并到一个张量来运算，因为单个张量的乘法内部则会自动并行。

此外，我们要对序列做 Mask 以忽略填充部分的影响。一般的 Mask 是将填充部分置零，但 Attention 中的 Mask 是要在 softmax 之前，把填充部分减去一个大整数（这样 softmax 之后就非常接近 0 了）。这些内容都在代码中有对应的实现。

## TensorFlow 版

[https://github.com/bojone/attention/blob/master/attention\\_tf.py](https://github.com/bojone/attention/blob/master/attention_tf.py)

## Keras 版

[https://github.com/bojone/attention/blob/master/attention\\_keras.py](https://github.com/bojone/attention/blob/master/attention_keras.py)

## 代码测试

在 Keras 上对 IMDB 进行简单的测试（不做 Mask）：

```
from __future__ import print_function
from keras.preprocessing import sequence
from keras.datasets import imdb

max_features = 20000
maxlen = 80
batch_size = 32

print('Loading data...')
(x_train, y_train), (x_test, y_test) = imdb.load_data(num_words=max_features)
print(len(x_train), 'train sequences')
print(len(x_test), 'test sequences')
```

```
print('Pad sequences (samples x time)')
x_train = sequence.pad_sequences(x_train, maxlen=maxlen)
x_test = sequence.pad_sequences(x_test, maxlen=maxlen)
print('x_train shape:', x_train.shape)
print('x_test shape:', x_test.shape)

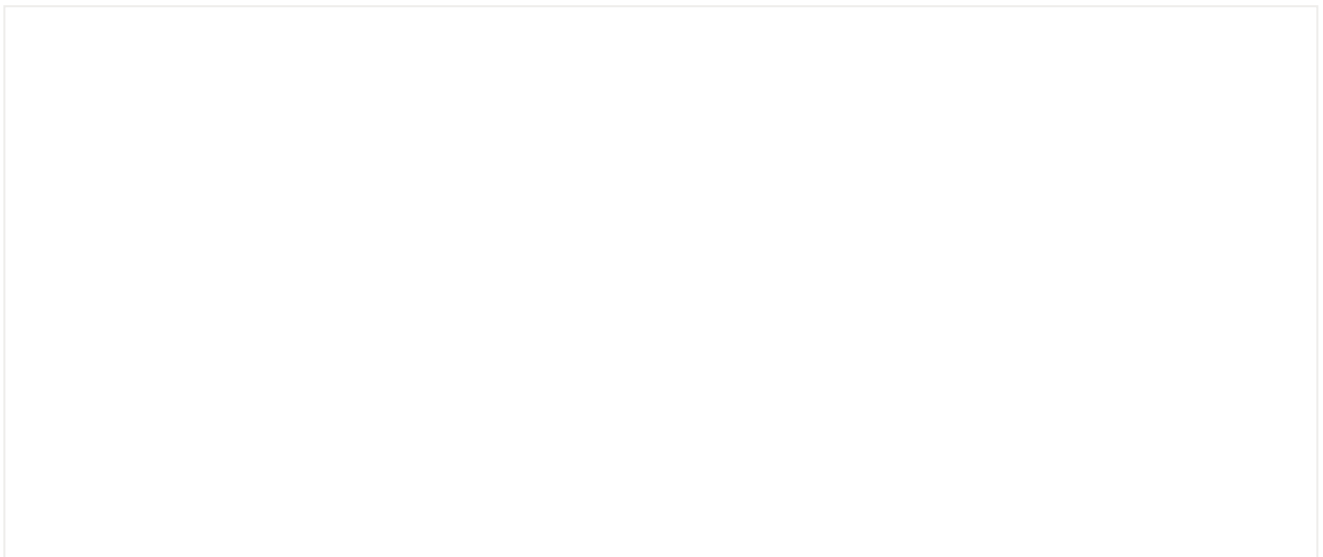
from keras.models import Model
from keras.layers import *

S_inputs = Input(shape=(None, ), dtype='int32')
embeddings = Embedding(max_features, 128)(S_inputs)
#embeddings = Position_Embedding()(embeddings) #增加Position_Embedding能轻微提高准确率
O_seq = Attention(8,16)([embeddings,embeddings,embeddings])
O_seq = GlobalAveragePooling1D()(O_seq)
O_seq = Dropout(0.5)(O_seq)
outputs = Dense(1, activation='sigmoid')(O_seq)

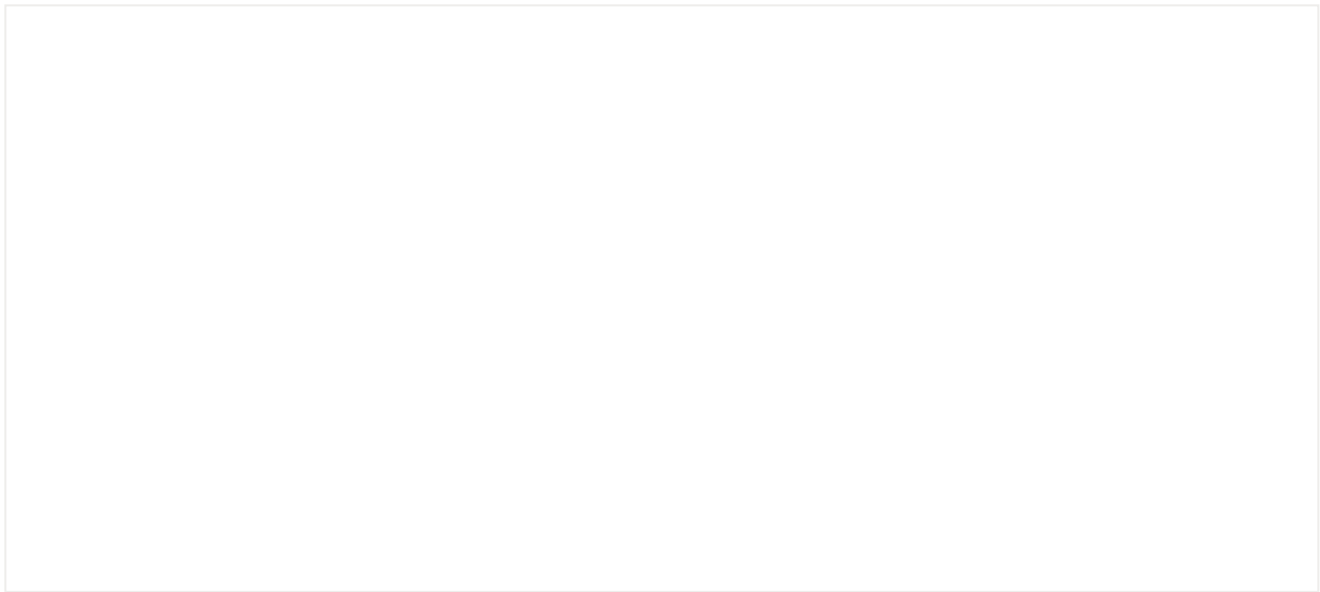
model = Model(inputs=S_inputs, outputs=outputs)
# try using different optimizers and different optimizer configs
model.compile(loss='binary_crossentropy',
              optimizer='adam',
              metrics=['accuracy'])

print('Train...')
model.fit(x_train, y_train,
        batch_size=batch_size,
        epochs=5,
        validation_data=(x_test, y_test))
```

无 Position Embedding 的结果:



有 Position Embedding 的结果:



貌似最高准确率比单层的 LSTM 准确率还高一点，另外还可以看到 Position Embedding 能提高准确率、减弱过拟合。

## 计算量分析

可以看到，事实上 Attention 的计算量并不低。比如 Self Attention 中，首先要对  $X$  做三次线性映射，这计算量已经相当于卷积核大小为 3 的一维卷积了，不过这部分计算量还只是  $O(n)$  的；然后还包含了两次序列自身的矩阵乘法，这两次矩阵乘法的计算量都是  $O(n^2)$  的，要是序列足够长，这个计算量其实是很难接受的。

这也表明，restricted 版的 Attention 是接下来的研究重点，并且将 Attention 与 CNN、RNN 混合使用，才是比较适中的道路。

## 结语

感谢 Google 提供的精彩的使用案例，让我等在大开眼界之余，还对 Attention 的认识更深一层。Google 的这个成果在某种程度上体现了“大道至简”的理念，的确是 NLP 中不可多得的精品。

本文围绕着 Google 的大作，班门弄斧一番，但愿能够帮助有需要的读者更好的理解 Attention。最后恳请大家建议和批评。

我是彩蛋

解锁新功能：热门职位推荐！

PaperWeekly小程序升级啦

今日arXiv√猜你喜欢√热门职位√

找全职找实习都不是问题

解锁方式

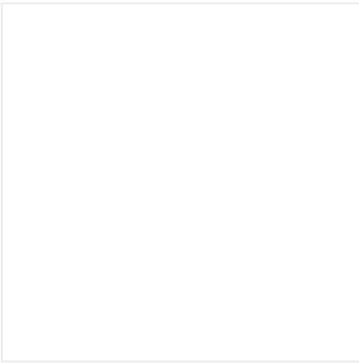
- 1. 识别下方二维码打开小程序
- 2. 用PaperWeekly社区账号进行登陆
- 3. 登陆后即可解锁所有功能

职位发布

请添加小助手微信（pwbot01）进行咨询

长按识别二维码，使用小程序

\*点击阅读原文即可注册



关于PaperWeekly

PaperWeekly 是一个推荐、解读、讨论、报道人工智能前沿论文成果的学术平台。如果你研究或从事 AI 领域，欢迎在公众号后台点击「交流群」，小助手将把你带入 PaperWeekly 的交流群里。



▽ 点击 | [阅读原文](#) | 加入社区

阅读原文