来自专辑 卖萌屋@自然语言处理



作者: 苏剑林(来自追一科技,人称"苏神")

等;也有改进模型结构的,比如Transformer-XL等;有改进训练方式的,比如ALBERT的参数共享 等... 以上的这些改动,都是在Attention外部进行改动的,也就是说它们都默认了Attention的合理性,没有 对Attention本身进行改动。而本文我们则介绍关于两个新结果:它们针对Multi-Head Attention中可 能存在建模瓶颈,提出了不同的方案来改进Multi-Head Attention。两篇论文都来自Google,并且做 了相当充分的实验,因此结果应该是相当有说服力的了 。

再小也不能小key\_size

第一个结果来自文章《Low-Rank Bottleneck in Multi-head Attention Models》,它明确地指出了 Multi-Head Attention里边的表达能力瓶颈,并提出通过增大key\_size的方法来缓解这个瓶颈。

## 首先简单回顾一下Multi-Head Attention,Multi-Head Attention的基础是自然是Single-Head Attention, 也叫Scaled-Dot Attention, 定义如下:

其中 $Q \in \mathbb{R}^{n \times d_k}$ , $K \in \mathbb{R}^{m \times d_k}$ , $V \in \mathbb{R}^{m \times d_v}$ 。而Multi-Head Attention,就是将Q,K,V分别用h个不 同的投影矩阵投影h次,然后分别做h次Single-Head Attention,最后把结果拼接起来,即

**Multi-Head Attention** 

在实际使用中,Q, K, V一般具有相同的特征维度 $d_k = d_v = d$ (即hidden\_size),比如BERT

Base里边是768;
$$h$$
一般选择12、16、24等,比如BERT base里边是12;确定了 $d,h$ 之后,通常的选择是让投影矩阵 $\mathbf{W} \in \mathbb{R}^{d \times (d/h)}$ ,也就是说,每个Attention Head里边,是将原始的 $d$ 维投影到 $d/h$ 维,然后在进行Attention运算,输出也是 $d/h$ 维,最后把 $h$ 个 $d/h$ 维的结果拼接起来,得到一个 $d$ 维的

秩瓶颈而减少h的做法可能得不偿失;如果增加d的话,那自然是能够增强模型整体表达能力的,但整 个模型的规模与计算量也会剧增,似乎也不是一个好选择。 那没有其他办法了吗?有!当我们用投影矩阵将Q, K, V都投影到低维时,前面都是将它们投影到 d/h维,但其实它们的维度不一定要相等,而是只需要保证Q, K的维度相等就行了(因为要做内 积),为了区别,我们通常称Q, K的维度为key\_size,V的维度才叫head\_size,改变key\_size的大

增大key\_size这个想法很简单,也容易实现,但是否真的有效呢? 我们来看看原论文的实验结果,其 实验都是以BERT为baseline的,实验结果图表很多,大家直接看原论文为好,这里只分享比较有代表 性的一个:

和head数,使得参数量可以跟原始的BERT设计一致,但是效果更优!所以,增加key\_size确实是有 意义的,哪怕将总体参数量重新调整到原来的一样大,也能一定程度上提升模型的效果。这无疑对我 们设计新的Transformer模型(尤其是小规模的模型)有重要的指导作用。 最后,附上我们预训练的两个增大了key size的RoBERTa小模型,欢迎大家使用(我们称之为

因此softmax之后无法很好地建议完整的二元分布。为了缓解这个问题,除了增大key\_size之外,还有 没有其他方法呢?有,比如这篇论文使用的混合分布思路。

所谓混合分布,就是多个简单分布的叠加(比如加权平均),它能极大地增强原分布的表达能力。典

在前一篇论文里边,我们提到了低秩瓶颈,也就是由于key\_size太小所以 $oldsymbol{Q}^{(i)}oldsymbol{K}^{(i)}$ 表达能力不足,

从单一分布到混合分布

RoBERTa<sup>+</sup>):

## $egin{aligned} egin{aligned} m{J}^{(1)} \ m{J}^{(2)} \ m{J}^{(2)} \ m{J}^{(h)} \end{aligned} = egin{pmatrix} \lambda_{11} & \lambda_{12} & \cdots & \lambda_{1h} \ \lambda_{21} & \lambda_{22} & \cdots & \lambda_{2h} \ m{\vdots} & m{\vdots} & \ddots & m{\vdots} \ \lambda_{h1} & \lambda_{h2} & \cdots & \lambda_{hh} \end{aligned} egin{pmatrix} m{\hat{J}}^{(1)} \ m{\hat{J}}^{(2)} \ m{\hat{J}}^{(h)} \end{aligned}$

 $oldsymbol{O}^{(1)} = oldsymbol{P}^{(1)} oldsymbol{V}^{(1)}, \quad oldsymbol{O}^{(2)} = oldsymbol{P}^{(2)} oldsymbol{V}^{(2)}, \quad , \cdots, \quad oldsymbol{O}^{(h)} = oldsymbol{P}^{(h)} oldsymbol{V}^{(h)}$ 

 $oldsymbol{O} = \left[oldsymbol{O}^{(1)}, oldsymbol{O}^{(2)}, \ldots, oldsymbol{O}^{(h)}
ight]$ 写起来很复杂,事实上很简单,**就是在"m{Q}m{K}^ op之后、softmax之前"用一个参数矩阵m{\lambda}将各个m{Q}m{K}^ op** 的结果叠加一下而已。这样就把原本是孤立的各个Attention Head联系了起来,即做了一个简单的

对上述公式,做两点补充说明: 1、简单起见,上述公式中笔者省去了缩放因子 $\sqrt{d_k}$ ,如果有需要,读者自行补充上去即可;

Talking-heads 这个结果显示,使用Talking-Head Attention情况下,保持hidden\_size不变,head数目越大(相应地 key\_size和head\_size都越小),效果越优。这看起来跟前一篇增大key\_size的结论矛盾,但事实上

12

12

24

24

 $^{24}$ 

12

12

12

12

12

12

24

 $^{24}$ 

12

24 24\*

12

12

24

24

12

12

12

12

12

12

24\*

12

12

12

12

12

12

12

12

24\*

24\*

插曲:神奇的论文画风

64

64

32

32

32

64

64

64

64

64

64

32

32

32

64

64

64

32

32

32

64

64

64

64

64

64

32

32

32

64

每个超参组合都做了实验,

73.16

71.98

73.98

75.80

72.70

73.16

74.55

73.85

74.51

73.61

71.98

73.98

73.90

74.65

72.53

0.416

81.38

80.52

81.35

81.96

79.45

81.38

81.19

80.73

81.62

80.61

80.52

81.35

82.72

82.60

76.85

0.365

81.13

86.31

84.24

90.60

78.80

81.13

86.19

82.39

86.06

79.72

86.31

84.24

84.20

84.76

94.37

3.237

该调的参数我们都调了,就我们的Talking-Heads Attention最好~

Multi-head

Talking-heads

Talking-heads

Talking-heads

Talking-heads

Talking-heads

Talking-heads

Talking-heads

Talking-heads

这正说明了混合分布对分布拟合能力明显提升作用,能够将key\_size缩小时本身变弱的单一分布, 叠加成拟合能力更强大的分布。当然,这不能说明就直接设key\_size=1就好了,因为key\_size=1时 计算量会远远大于原始的BERT base,应用时需要根据实际情况平衡效果和计算量。

也没什么反响,估计也跟这个任性的风格有关系~

• ACL2020 | 对话数据集Mutual: 论对话逻辑, BERT还差的很远 • 如何优雅地编码文本中的位置信息? 三种positioanl encoding方法简述

本质上的提升还是来源于增大key\_size, 所以本文只强调了增大key\_size这一步。此外, 如果同 时增大key\_size和head\_size,那么会导致计算量和显存消耗都明显增加,而只增大key\_size的

补充说明:事实上原论文考虑的是同时增大key\_size和head\_size、然后Multi-Head Attention的 输出拼接之后再用一个变换矩阵降维,但笔者认为由于拼接降维这一步只是一个线性变换,所以

话,增加的资源消耗就小很多了。

来看看实验结果~

https://github.com/ZhuiyiTechnology/pretrained-models 再缺也不能缺Talking

那么"多个"低秩分布哪里来呢?不是有Multi-Head嘛,每个head都带有一个低秩分布,就直接用它们 叠加就行了,这就是Talking-Heads Attention了。具体来说,它的形式是:  $\hat{oldsymbol{J}}^{(1)} = oldsymbol{Q}^{(1)} oldsymbol{K}^{(1)}^{ op}, \quad \hat{oldsymbol{J}}^{(2)} = oldsymbol{Q}^{(2)} oldsymbol{K}^{(2)}^{ op}, \quad \cdots, \quad \hat{oldsymbol{J}}^{(h)} = oldsymbol{Q}^{(h)} oldsymbol{K}^{(h)}^{ op}$ 

 $oldsymbol{P}^{(1)} = softmax \Big(oldsymbol{J}^{(1)}\Big), oldsymbol{P}^{(2)} = softmax \Big(oldsymbol{J}^{(2)}\Big), \ldots, oldsymbol{P}^{(h)} = softmax \Big(oldsymbol{J}^{(h)}\Big)$ 

 $d_k = d_v$ 

64

128

64

32

16

12

heads

12

6

12

24

48

64

96

192

384

768

MNLI

82.6

83.4

83.6

83.6

83.4

83.8

83.6

83.9

83.9

84.2

WSC

SQuAD1.1 (f1)

88.51

88.8

89.2

89.4

89.5

89.9

89.3

89.8

90.5

90.5

F1 EM $d_k$  $d_v$ Acc Acc Acc F1 EMF1 Acc Acc Acc Average multi-head 12812872.0078.8478.8387.50 70.0075.0136.9472.0871.3283.0369.1279.8112 64 89.15multi-head 6473.5980.3189.2973.0075.1337.5773.9272.9483.7569.2877.88multi-head 243232 73.9881.3584.2491.0770.0075.9839.2474.4073.3783.0371.0078.8516 multi-head 48 16 71.8580.34 77.9087.50 67.0075.2637.6772.3171.3284.1269.4477.88talking-heads 128 72.5783.37 76.76 75.3774.4883.39 12889.2965.0040.0866.3080.7780.83talking-heads 12 12 12 64 81.1389.2969.0077.3640.5075.6383.39 65.8382.696473.1681.3876.4532 75.80talking-heads 24 243290.6094.6475.0076.5639.7777.2276.3785.2081.732481.9670.0616 talking-heads 48 48 16 76.3982.9487.46 91.0774.0078.1142.7177.5176.68 69.5987.5084.84multi-head 2464 74.0877.9087.50 73.0077.2039.3576.8976.1183.3969.4480.776481.8012 768768 73.3583.46 89.2976.92general bilinear 81.4771.0076.6939.2476.8076.0285.9269.59talking-heads 128 128 72.5783.37 89.2965.0076.76 40.0875.3774.4883.39 66.3080.7780.83talking-heads 246 12812873.7181.0780.92 87.50 68.0076.4740.8274.7773.8985.9270.0679.8132 talking-heads 246 243273.5680.9283.5289.2975.0075.6437.1574.5873.7381.9571.3276.92talking-heads 24 128 3274.2987.62 74.0076.2337.0476.6782.692480.9591.0775.8383.39 68.65talking-heads 24 32 128 246 76.3781.7788.9792.8676.0077.6342.8176.7275.8886.2867.7187.5094.64talking-heads 2424323275.8090.602481.9675.0076.5639.7777.2276.3785.2070.0681.7332 75.98 74.4073.37 multi-head 2432 73.9884.2470.0039.2483.03 78.8581.3591.0771.0024 3232 project logits 72.6383.1589.2971.0076.9839.3575.0074.2185.2069.7579.81 $^{24}$ 81.47project weights 3224243274.0581.9981.9687.50 73.0077.3541.0376.6275.7485.2068.0378.852432talking-heads 24 32 75.8090.6094.64 75.0039.7777.2281.732481.9676.5676.3785.2070.06multi-head 24323273.9884.2491.0770.0075.9839.2474.4073.3783.03 71.0078.8581.3532 TH-enc-self 2424\*3273.9082.7284.2089.2969.0078.1843.5576.0775.3485.9269.7579.81TH-dec-self 24\*32 32 24\* $^{24}$ 72.5479.9387.50 67.0038.0973.7272.9483.03 69.4479.8180.9276.4932 TH-encdec  $^{24}$ 24\*3273.4478.6787.50 74.0076.4139.2475.3574.4783.39 24\*80.8370.5380.773232 talking-heads 2424 $^{24}$ 75.8081.9690.6094.6475.0076.5639.7777.2276.3785.2070.0681.73multi-head 12 646473.5980.3189.1589.2973.0075.1337.5773.9272.9483.7569.2877.88

89.29

91.07

91.07

94.64

87.50

89.29

91.07

89.29

92.86

89.29

91.07

91.07

89.29

91.07

94.64

2.560

69.00

66.00

70.00

75.00

72.00

69.00

74.00

74.00

72.00

70.00

66.00

70.00

69.00

69.00

70.00

2.741

77.36

76.47

75.98

76.56

74.68

77.36

76.20

75.50

75.48

75.36

76.47

75.98

78.18

77.32

67.64

0.716

并给出实验结果。如此强大的实验阵容,

40.50

37.67

39.24

39.77

35.15

40.50

41.45

38.51

38.82

38.20

37.67

39.24

43.55

42.29

28.75

1.011

76.45

73.24

74.40

77.22

72.39

76.45

75.45

73.72

75.54

75.19

73.24

74.40

76.07

77.88

70.84

0.370

75.63

72.42

73.37

76.37

71.70

75.63

74.63

72.93

74.70

74.41

72.42

73.37

75.34

76.99

69.90

0.379

83.39

82.31

83.03

85.20

80.14

83.39

84.84

83.03

83.03

83.39

82.31

83.03

85.92

84.48

74.73

1.228

65.83

67.40

71.00

70.06

68.34

65.83

68.81

68.50

68.65

67.55

67.40

71.00

69.75

71.16

67.71

0.850

基本上也就只有

82.69

75.00

78.85

81.73

82.69

82.69

82.69

81.73

82.69

83.65

75.00

78.85

79.81

79.81

77.88

2.029

话说回来,笔者在Arxiv上首次刷到《Talking-Heads Attention》这篇论文时,第一感觉是一篇垃圾论
文。为啥? 因为它的画风是这样的:
3 Review of Attention Algorithms
3.1 Dot-Product Attention
Simple dot-product attention can be described by the pseudocode below. The logits L are computed as the dot-products of the query-vectors and the memory-vectors. For each query, the logits are passed through a softmax function to produce weights, and the different memory-vectors are averaged together, weighted by those weights. In this code, we show the case where there are $n$ different queries all attending to the same $m$ memory-vectors. If there is only one query, the code is identical except that the "n" dimension is removed from all tensors.
<pre>def DotProductAttention(     X[n, d], # n query-vectors with dimensionality d     M[m, d]): # m memory-vectors with dimensionality d     L[n, m] = einsum(X[n, d], M[m, d]) # Attention logits     W[n, m] = softmax(L[n, m], reduced_dim=m) # Attention weights     Y[n, d] = einsum(W[n, m], M[m, d])     return Y[n, d]</pre> 3.2 Dot-Product Attention With Projections
[Vaswani et al., 2017] propose a dimensionality-reduction to reduce the computational complexity of the
attention algorithm. In this version, instead of computing the attention algorithm directly on the inputs $X$ and $M$ , we first project the inputs using the learned linear projections $P_q$ , $P_k$ and $P_v$ , to produce lower-dimensional query-vectors, key-vectors and value-vectors $Q$ , $K$ and $V$ . We use a fourth learned linear projection, $P_o$ , to produce the output.
<pre>def DotProductAttentionWithProjections(     X[n, d_X],     # n vectors with dimensionality d_X     M[m, d_M],     # m vectors with dimensionality d_M     P_q[d_X, d_k],    # learned linear projection to produce queries     P_k[d_M, d_k],     # learned linear projection to produce keys     P_v[d_M, d_v],     # learned linear projection to produce values     P_o[d_Y, d_v]):    # learned linear projection of output     Q[n, d_k] = einsum(X[n, d_X], P_q[d_X, d_k])     # queries     K[m, d_k] = einsum(M[m, d_M], P_k[d_M, d_k])     # keys     V[m, d_v] = einsum(M[m, d_M], P_v[d_M, d_v])     # values     L[n, m] = einsum(Q[n, d_k], K[m, d_k])     # Attention logits</pre>

- LayerNorm是Transformer的最优解吗? • 知乎搜索框背后的Query理解和语义召回技术 • 详解ERNIE-Baidu进化史及应用场景
- 深入解析GBDT二分类算法(附代码实现) 夕小瑶的卖萌屋
- 关注&星标小夕,带你解锁AI秘籍 订阅号主页下方「撩一下」有惊喜哦

···· 点击查看精选留言

前言

(公式可以左右滑动哦)

 $m{Q}^{(h)} = m{Q}m{W}_O^{(h)}, m{K}^{(h)} = m{K}m{W}_K^{(h)}, m{V}^{(h)} = m{V}m{W}_V^{(h)}, m{O}^{(h)} = Attention(m{Q}, m{K}, m{V})$  $oldsymbol{O} = \left[oldsymbol{O}^{(1)}, oldsymbol{O}^{(2)}, \ldots, oldsymbol{O}^{(h)}
ight]$ Attention里有个瓶颈

输出。这里的d/h我们通常称为head\_size。

在Attention中,关键的一步是  $m{P} = softmaxigg(rac{m{Q}m{K}^{ op}}{\sqrt{d_k}}igg)$  ( 2 ) 这一步是描述了Q与K的两两向量之间的联系,我们可以将P看成一个二元联合分布(实际上是n个

一元分布,不过这个细节并不重要),如果序列长度都为n,也就是每个元有n个可能的取值,那么这 个分布共有 $n^2$ 个值。 但是,我们将Q, K分别投影到低维后,各自的参数量只有 $n \times (d/h)$ ,总的参数量是2nd/h,所以 式(2)就相当于用2nd/h的参数量去逼近一个本身有 $n^2$ 个值的量,而我们通常有 $2nd/h \ll n^2$ ,尤其 是h比较大时更是如此,因此**这种建模有点"强模型所难",这就是原论文中的"低秩瓶颈(Low-Rank** Bottleneck)"的含义. 不妨试试增大key\_size?

小而不改变head\_size的话,也不影响模型的hidden\_size。 所以,这篇论文提出来的解决方法就是增大模型的key\_size,它能增加Attention的表达能力,并且不 改变模型整体的hidden\_size, 计算量上也只是稍微增加了一点。

对Multi-Head Attention改进的第二个结果来自论文《Talking-Heads Attention》,这篇论文虽然没有 显式地指出它跟前一篇论文的联系,但笔者认为它们事实上在解决同一个问题,只不过思路不一样: 它指出当前的Multi-Head Attention每个head的运算是相互孤立的,而通过将它们联系(Talking) 起来,则可以得到更强的Attention设计,即标题的"Talking-Heads Attention"。

Talking.

2、更一般的Talking-Heads Attention允许可以在 $m{J}=m{\lambda}\hat{m{J}}$ 这一步进行升维,即叠加出多于h个混 再来看看实验结果

容:

dyn. proj.  $P_{Xl}$ dyn. proj.  $P_{Xw}$ dyn. proj.  $P_{Ml}$ dyn. proj.  $P_{Mw}$ dyn. proj. multi-head TH-enc-self DP-enc-self Raffel et al., 2019 ibid. stddev. 几乎每个任务。 Google能搞出来了,而且整篇论文明显是浓浓的"T5 Style"(还没看过T5论文的读者,可以去 Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer》感受一下),果不 其然,作者之一Noam Shazeer也正是T5的作者之一。 笔者只想说,这种庞大的实验轰炸,仿佛在向我们宣告着:

talking-heads

talking-heads

talking-heads

dyn. proj.

multi-head

dyn. proj.

谁能想象到,一篇如此强大的论文,里边居然没有一条数学公式,取而代之的全是伪代码!!其实伪 还是有任性的代价的,这篇实验阵容这么强大又这么有效的论文,发布至今也有一个多月了,但似乎

return Y[n, d\_Y]

arxiv访问慢的小伙伴也可以在订阅号后台回复关键词【0415】下载论文PDF。 • ACL2020 | FastBERT: 放飞BERT的推理速度

阅读原文

《Attention is All You Need》一文发布后,基于Multi-Head Attention的Transformer模型开始流行起 来,而去年发布的BERT模型更是将Transformer模型的热度推上了又一个高峰。当然,技术的探索是 无止境的,改进的工作也相继涌现:有改进预训练任务的,比如XLNET的PLM、ALBERT的SOP等; 有改进归一化的,比如Post-Norm向Pre-Norm的改变,以及T5中去掉了Layer Norm里边的beta参数

arxiv访问慢的小伙伴也可以在订阅号后台回复关键词【0415】下载论文PDF。

 $Attention(oldsymbol{Q}, oldsymbol{K}, oldsymbol{V}) = softmaxigg(rac{oldsymbol{Q} oldsymbol{K}^{ op}}{\sqrt{d_1}}igg)oldsymbol{V}$  ( 1 )  $m{Q}^{(1)} = m{Q}m{W}_O^{(1)}, m{K}^{(1)} = m{K}m{W}_K^{(1)}, m{V}^{(1)} = m{V}m{W}_V^{(1)}, m{O}^{(1)} = Attention(m{Q}, m{K}, m{V})$  $oldsymbol{Q}^{(2)} = oldsymbol{Q}oldsymbol{W}_{O}^{(2)}, oldsymbol{K}^{(2)} = oldsymbol{K}oldsymbol{W}_{K}^{(2)}, oldsymbol{V}^{(2)} = oldsymbol{V}oldsymbol{W}_{V}^{(2)}, oldsymbol{O}^{(2)} = Attention(oldsymbol{Q}, oldsymbol{K}, oldsymbol{V})$ 

增加hidden\_size大小d。但是更多的Attention Head本身也能增强模型的表达能力,所以为了缓解低

那么,解决办法是什么呢?直接的想法是让2nd/h增大,所以要不就是减少head的数目h,要不就是

这个结果显示,如果固定一个比较大的key\_size(比如128),那么我们可以调整模型的hidden\_size

型的例子是高斯混合模型:我们知道高斯分布只是一个常见的简单分布,但多个高斯分布叠加而成的 高斯混合分布(也叫高斯混合模型,GMM)就是一个更强的分布,理论上来说,只要叠加的高斯分布 足够多,高斯混合分布能逼近任意概率分布。这个例子告诉我们,想要增加Attention中分布的表达能 力,又不想增加key\_size,那么可以考虑叠加多个低秩分布。

合分布,然后再用另一个参数矩阵降维,但这并不是特别重要的改进,所以不在主要篇幅介绍。 是不是真的有效,当然还是得靠实验结果来说话。这篇论文的实验阵容可谓空前强大,它同时包含了 BERT、ALBERT、T5为baseline的实验结果! 众所周知, BERT、ALBERT、T5均是某个时间段的 NLP最优模型,尤其是T5还是处在superglue的榜首,并且远超出第二名很多,而这个Talking-Heads Attention则几乎是把它们的辉煌战绩又刷到了一个新高度!

上述表格只是原论文实验结果的冰山一角,这里再放出一个实验表格,让大家感受感受它的实验阵 Table 12: T5 on SuperGLUE Language-Understanding Benchmark [Wang et al., 2019] (dev). Experiments described in Section 7.1 and appendix A. ReCoRDReCoRD RTE WiC BoolQ CB CBCoPA MultiRC MultiRC Score

代码都算不上,感觉更像是直接把实验中的Python代码复制到了论文中,还是复制到论文主体上!笔 者印象里,只有那些不入流的水论文才会这样做,所以笔者看到的第一想法就是水文一篇。也就 Google的大佬们才能这么任性,要不是耐着心多扫了几眼,要不是不小心扫到了T5等字眼,要不是回 去看作者居然清一色是Google的,这篇强大的论文就被笔者当作垃圾论文放到回收站了 不过,任性

 $O[n, d_v] = einsum(W[n, m], V[m, d_v])$ 

 $Y[n, d_Y] = einsum(O[n, d_V], P_o[d_Y, d_V])$ 

W[n, m] = softmax(L[n, m], reduced\_dim=m) # Attention weights

本文介绍了两个关于Multi-Head Attention的后续改进工作,虽然改进细节不一致,但可以说它们都是 针对"低秩瓶颈"这个问题而提出的,有种殊途同归之感。两个工作都来自Google,实验内容都很丰 富,所以结果都比较有说服力,正在做模型结构改进工作的读者可以参考参考。