

机器阅读理解的双问答网络

韩晓, 冯峰, 冯延建, 郑静尧

摘要

阅读理解设置有三种形式：问题，答案和背景。问答或问题生成的任务旨在在给对方时推断答案或问题

基于背景。我们提出了一种新颖的双向神经序列转导模型，它连接三种模态，允许它同时和相互学习两个任务

相互受益。在训练期间，模型接收问题 - 上下文 - 答案三元组作为输入，并通过分层注意过程捕获跨模态交互。与以前利用数据级问题生成和问答的二元性的联合学习范式不同，我们通过镜像网络结构和部分共享不同层的组件来在架构级解决这样的双重任务。这使得知识可以从一个任务转移到另一个任务，帮助模型找到每种模态的一般表示。对四个公共数据集的评估表明，我们的双学习模型优于单一学习模型以及问答和问题生成任务的最新联合模型。

介绍

最近机器阅读理解的任务 (Rajpurkar等人2016; Nguyen等人2016) 受到了NLP和AI研究团体越来越多的关注。该任务尝试使机器在阅读段落落后回答问题。这项任务有三种形式：问题，答案和背景。已经提出了许多成功的问答模型来通过模拟问题和上下文模态之间的相互作用来填充答案模态 (Seo等人2016; Yu等人2018)。直观地说，问题和答案非常相似，因为它们既是短文本又与给定的上下文密切相关。很少有最近的作品已经认识到这种关系并从不同的角度进行探索 (Song, Wang和Hamza 2017; Wang, Yuan和Trischler 2017; Tang et al. 1188)。他们有希望的结果表明：(i) 问题和答案的作用可以在给定背景的情况下切换，因为相同的网络结构可以用于问答 (QA) 和问题生成 (QG)；(ii) 可利用这种可逆性或二元性来改善质量保证和质量控制。

在这项工作中，我们将机器阅读理解视为一个双重学习问题，即学习回答与学习要求。利用其中的共性并不奇怪

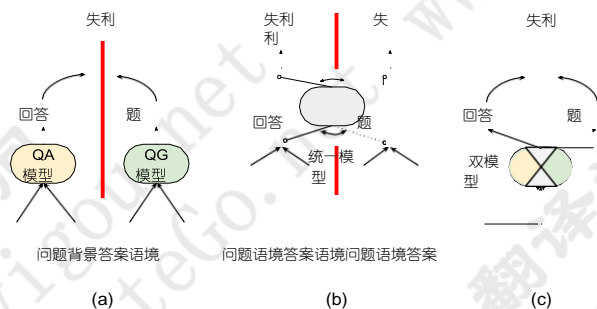


图1：三种学习范例，用于利用问答和问题生成的任务关联。红线代表数据/模型级别的分离。(a) 两个具有关节损失功能的分离模型 (Song, Wang和Hamza 2017)。(b) 具有交替训练输入和两个独立损失函数的统一模型 (Wang, Yuan和Trischler 2017)。(c) 这项工作：具有本地共享结构的统一架构，可同时学习两项任务。与 (a) 和 (b) 相反，在我们的学习范例中没有数据级别或模型级别的分离。

任务，学习范式中需要一些共享方案。图1显示了以前的双重学习范例和这项工作中的范式。具体来说，我们提出了一种新的双向神经序列转导模型，它共同解决了这两个任务。该模型在结构上是对称的，其组件在不同级别的两个任务之间共享，如图2所示。在训练期间，模型接收问题 - 上下文 - 答案三元组作为输入，并通过分层注意过程捕获跨模态交互。在测试期间，模型在给定基于上下文的对应物时生成答案或问题作为序列。我们的贡献总结如下：

- 我们提出了一种新颖的双向神经序列转导模型，可以同时学习问答和问题生成。我们将在两个任务中扮演类似角色的网络组件联系起来，以便在培训期间传递跨任务知识。问题，上下文和答案的跨模态交互由一对对称的分层关注过程捕获。据我们所知，我们是第一个在机器阅读理解设置下以这种粒度利用二元性的人。
- 我们证明了我们的模型在四个公共数据集上的有效性，并取得了可喜的成果。我们的模型优于单一学习版本和最先进的模型

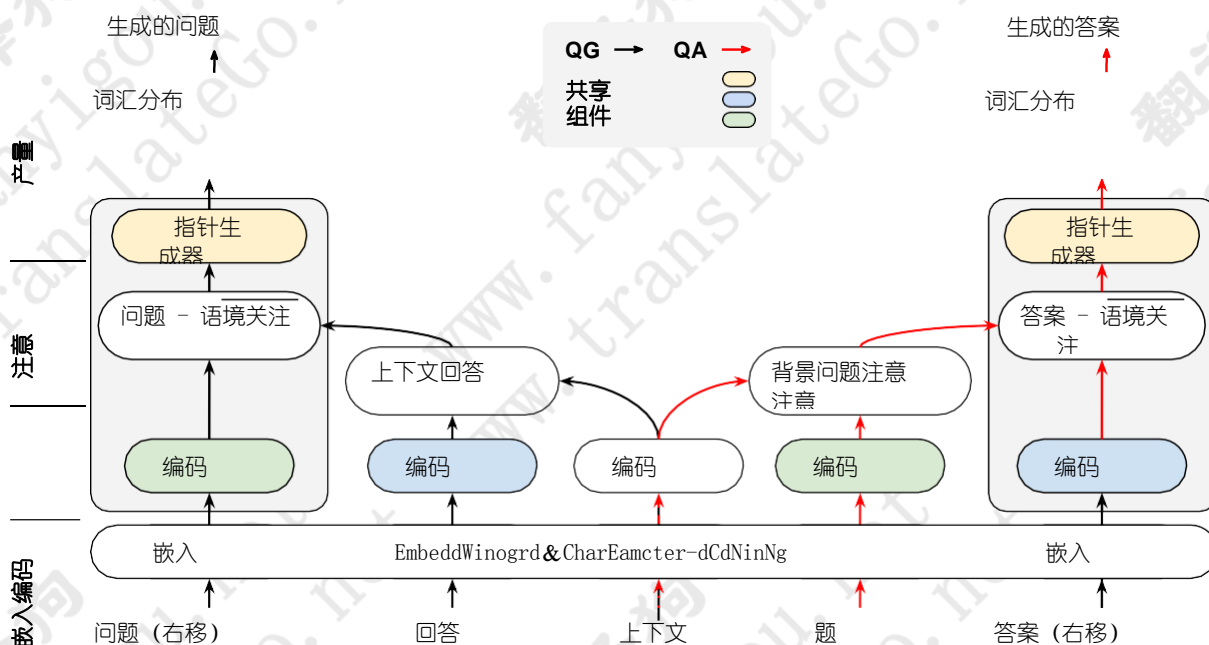


图2: 双问答网络 (DAANET) 的模型架构。最好看的颜色。它是一个由四层组成的分层过程: 嵌入, 编码, 注意和输出。侧面的矩形超级块可以分别被视为QG和QA的解码器。QG和QA任务的计算流程分别用黑色箭头和红色箭头绘制。共享 (包括部分共享) 组件使用相同的颜色填充, 即答案编码器为蓝色, 问题编码器为绿色, 指针生成器为黄色。注意, QA和QG也共享上下文编码器。为清楚起见, 我们只绘制一个上下文编码器块。在测试期间, 移位的输入将由模型自己生成的前一步骤中的单词替换。

联合问答模型。我们提供了强有力的证据证明将二元性纳入模型结构以提高阅读理解能力。本文中使用的代码和数据可在线获取, 以便将来进行比较¹。

相关工作

我们的工作涉及机器阅读理解 (MRC), 问答 (QA) 和问题生成 (QG) 方面的现有工作。在本节中, 我们将从这些角度简要回顾以前的工作。

在过去几年中已经开发了大量的MRC模型, 例如BiDAF (Seo等人2016), S-Net (Tan等人2017), R-Net (Wang等人2017), match-LSTM (Wang and Jiang 2016), ReasonNet (Shen et al. 2017), Document Reader (Chen et al. 2017), Reinforced Mnemonic Reader (Hu et al. 2018), Fu-sionNet (Huang et al. 2017) 和QANet (Yu等人, 2018年)。大多数现有模型依赖于假设答案是给定段落的连续跨度。在这种假设下, 答案可以简化为一对两个整数, 分别代表其在段落中的起点和终点位置。在这项工作中, 我们没有做出这样的假设。我们的模型生成答案和问题作为序列, 因此相同的模型体系结构可用于问答和问题生成。总的来说, 我们相信这种生成模型能够为复杂的语义提供更好的表达能力, 因此更有可能在机器阅读理解中实现真正的突破。

关于问题生成的大部分先前工作依赖于特征工程, 手工模板和语言规则 (Aldabe等人2006; Heilman 2011; Liu, Calvo和Rus 2010; Bordes, Weston和Usunier 2014; Dhingra, Pruthi和Rajagopal 2018)。最近, 受到机器翻译中序列到序列模型的显著成功的鼓舞 (Sutskever, Vinyals和Le 2014; Luong, Pham和Manning 2015), 解析 (Vinyals等人2015) 和文本摘要 (Nallapati等) 使用深度神经网络来解决QG的兴趣迅速增加 (Du, Shao和Cardie 2017; Yuan等人2017; Duan等人2017)。这些研究的主要动机之一是使用QG来丰富QA的训练数据。例如, 通过使用深度神经网络将事实从知识库转换为自然语言问题, 生成具有30M QA对的语料库 (Serban等人, 2016)。

最后, 一些最近的MRC工作已经认识到QA和QG之间的关系并且以不同的方式利用它。尽管本研究分享了类似的目标, 但我们的工作从以下观点来看是独一无二的:

- 与 (Yang et al. 2017) 不同, 它依赖于答案的连续跨度假设和 (Sachan和Xing 2018; Tang等人2017; 2018) 考虑回答作为句子选择任务, 我们将问题回答和问题生成都视为序列转换使得可以对两个任务使用相同的模型架构。与 (Tan et al. 2017) 中提出的提取 - 合成框架不同, 我们的模型是完全端到端的训练。
- 与培训QA的 (Song, Wang和Hamza 2017) 不同

¹<https://github.com/hanxiao/daanet>

模型和QG模型独立使用相同的架构和 (Wang, Yuan和 Trischler 2017) 在同一模型的QA和QG示例之间交替训练数据, 我们的双向模型在训练期间直接消耗问题-上下文-答案三元组。给定三元组, 参数将在从QA和QG方向接收梯度的意义上“更新”两次。因此, 与其他两个相关工作相比, 我们模型中的参数训练得更充分。

- 与 (Tang et al. 2018) 使用协作检测器连接两个任务的训练不同, 我们提出了一种新的共享方案, 以利用架构级别的二元性。在两个任务中扮演类似角色的组件和参数绑定在一起。通过耦合两个对称的分层关注过程来捕获问题, 答案和上下文之间的相互作用。该共享方案极大地减少了两个任务所需的参数总数。

双问答网络

我们首先在MRC设置中制定询问和回答的双重学习问题, 然后呈现我们的模型: 双问答网络 (DAANET)。最后, 我们总结了模型中使用的注意力和二元性。

问题制定

与仅关注QA的传统MRC问题不同, 这项工作中考虑的问题是二分法: QA和QG。具体而言, 模型应该能够在根据上下文给出对应方时推断出答案或问题。

在形式上, 我们将上下文段落表示为 $C: c_1, \dots, c_n$, 问题句子为 $Q: q_1, \dots, q_m$ 和答案句子为 $A: a_1, \dots, a_k$ 。在续集中, 我们遵循这种表示法并使用 n, m, k 来表示上下文的长度, 问题和

分别是一个答案。上下文 C 由两个任务共享。给定 C , QA任务被定义为基于问题 Q 找到答案 A ; QG任务被定义为基于答案 A 找到问题 Q 。与以前主要假设答案为连续跨度的MRC模型相比, 我们将QA和QG任务视为生成问题并在神经序列中联合解决它们转导模型。

型号说明

我们提出的双问答网络的高级架构如图2所示。该神经网络转换模型接收字符串序列作为输入, 并通过嵌入层, 编码层, 关注层处理它们, 最后处理到输出层生成序列。

- 1. 嵌入图层。** 嵌入层将每个单词映射到高维向量空间。向量表示包括字级和字符级信息。该层的参数由上下文, 问题和答案共享。对于单词嵌入, 我们使用预先训练的256维GloVe (Pennington, Socher和Manning 2014) 单词向量, 在训练期间固定。一切外的单词词汇被映射到一个<UNK>令牌。除此之外, 还有三个特殊标记: <PAD>,

<START>和<END>。嵌入<START>, <END>和<UNK>是可训练和随机初始化, 而<PAD>被固定为零向量。

对于字符嵌入, 每个字符表示为200维可训练矢量。因此, 每个单词可以表示为一系列字符向量, 其中序列长度被截断或填充为16。接下来, 我们进行内核宽度为3的1D CNN, 然后沿时间轴进行最大池化。这为每个单词提供了固定大小的200维向量。最后, 字符嵌入和字嵌入向量的串联被线性投影到300维空间, 然后传递到高速公路网络 (Srivastava, Greff和Schmidhuber 2015), 如下所示:

$$e: "re_x, e_{\text{char}} sH_1 \cdot v_1, g: " \sigma \\ p e H_2 \cdot v_2 q, e_{\text{out}}: "g d e \cdot p l' g q d p e H_3 \\ \cdot v_3 q,$$

其中 $H_1 \in \mathbb{R}^{456 \times 300}$, $H_2, H_3 \in \mathbb{R}^{300 \times 300}$ 和 $v_1, v_2, v_3 \in \mathbb{R}$ 是学习参数; 你 σ 表示S形函数。输出的维数是300。

- 2. 编码层。** 编码层分别包含三个用于上下文, 问题和答案的编码器。它们由QA和QG共享, 如图2所示。也就是说, 给定QA和QG双任务, 原始任务的编码器和双任务的解码器被强制相同。参数共享方案用作影响两个任务的训练的正则化。它还有助于模型为每种模态找到更通用和稳定的表示。

每个编码器由以下基本构建块组成: 元素完全连接的前馈块, 堆叠的LSTM (Gers, Schmidhuber和Cummins 2000) 和自我关注块 (Vaswani等人2017)。自注意块允许编码器中的每个位置基于由其点积测量的相似性来参与所有位置。注意层中描述了自注意块中使用的注意功能。每个块之后是层标准化 (Ba, Kiros和Hinton 2016)。编码器的最终输出是所有块输出的串联, 如图3所示。

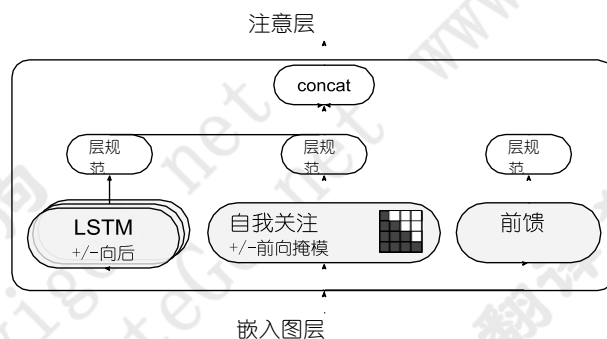


图3: 上下文/问题/答案编码器的体系结构。在问答编码器中, LSTM是单向的, 并且在计算自我关注时使用前向掩码。

上下文编码器由前馈网络, 3层双向LSTM和4头自我关注组成。问答编码器遵循相同的组成, 但需要进行必要的修改以防止“泄漏”

向左的信息。具体地，在问答中，编码器LSTM是单向的，并且在计算自我关注时使用前向掩模。正向掩模仅保留后期位置到早期位置的注意力，同时将剩余部分设置为0，如下所示，

$$\text{SelfAttn}_{ij} = \frac{1}{d} \exp(\mathbf{R}^T \mathbf{p}_{ij} \mathbf{Q}^T \mathbf{M}) \mathbf{e}_{ij}$$

\mathbf{M}_{ij} : “如果 $i \leq j$, 则为0 别的’8,

其中 \mathbf{e}_{ij} 是嵌入层的输出， $\mathbf{RR}^{300 \times d}$ 是学习参数。在实验部分中对编码器的不同成分进行基准测试。

3. 注意层。注意层将迄今为止观察到的所有信息与上下文编码C，应答编码A和问题编码Q以分层方式融合。在QG的时间t，我们首先将答案编码折叠到上下文编码中，即再次折叠成先前生成的编码Q的问题。质量保证部分遵循类似的注意程序。在这项工作中，我们将第一个折入步骤称为“上下文”注意，将第二个折入步骤称为“上下文”注意，其中上下文是第一个折入步骤的输出。图4显示了这个两步注意程序。

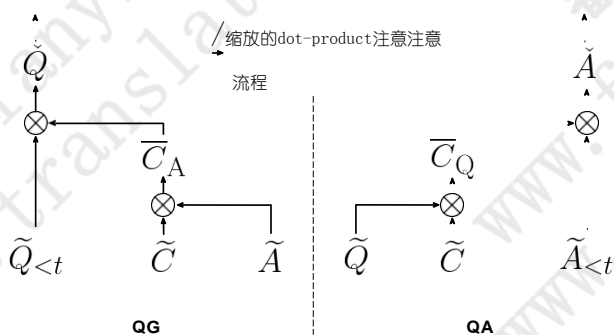


图4：注意层中实现的两步注意流程，其中输入 $\mathbf{Q}_r, \mathbf{C}_r, \mathbf{A}_r, \mathbf{Q}_{r:t}, \mathbf{A}_{r:t}$ 来自以前的编码层。

更具体地说，上下文答案注意力计算如下：我们首先计算每对上下文和答案词之间的归一化相似度为 $\text{softmax}(\frac{1}{d} \mathbf{C}_r^T \mathbf{U}_A \mathbf{R}^T \mathbf{V}_Q \mathbf{Q}_t)$ ，其中得分函数

$\mathbf{R}^T \mathbf{V}_Q \mathbf{Q}_t$ 是比例因子的倍增因子，这是Vaswani等人(2017)中所提出的。上下文回答

然后将注意力计算为答案的加权和，即 \mathbf{C}_A ：“spC, Aq”
A. 上下文问题注意遵循相同的过程，即 \mathbf{C}_Q ：s C, Q Q. 输出第一个折叠步骤在RRR，与上下文序列的大小相同。

第二个折入步骤涉及计算问题上下文和回答 - 上下文关注，这对于分别生成有意义的问答序列是必不可少的。它允许生成序列中的每个位置在第一个折叠步骤中以上下文顺序参与所有位置。在时间t，问题 - 上下文关注的输入由 \mathbf{Q}_r 来自

前面的编码层和 \mathbf{C}_A 从上下文回答注意，输出是 \mathbf{Q} ：“spQ_在, $\mathbf{C}_A \mathbf{Q} - \mathbf{C}_A$ 。答案 - 背景注意遵循相同的程序，即A：“spA_在, $\mathbf{C}_Q \mathbf{Q} - \mathbf{C}_Q$ 。

4. 输出层。输出层一次生成一个字的输出序列。在每个步骤中，模型都是自回归的 (Graves 2013)，在生成下一个时消耗先前作为输入生成的单词。在这项工作中，我们

使用指针生成器作为输出层的核心组件 (参见Liu和Manning 2017)。它允许通过指向从上下文复制单词，并从固定词汇表中采样单词。这有助于准确地再现信息，尤其是在QA中，同时保留生成新单词的能力。

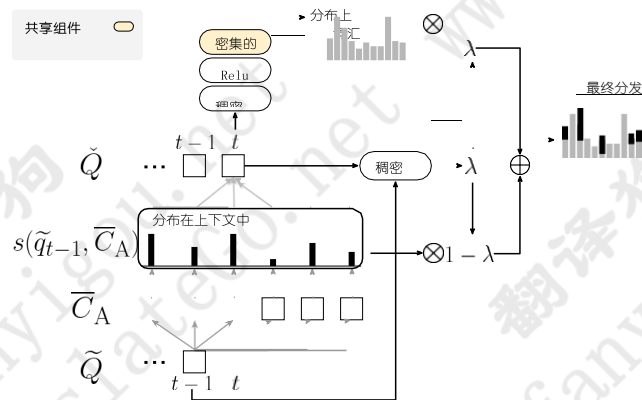


图5：QG输出层的体系结构。每个方块代表一个单词。最终分布是两个离散分布的混合。质量保证遵循相同的架构，几乎没有替代品。黄色密集层由QA和QG共享。

输出层的结构如图5所示。具体地说，在时间t，生成的单词w的概率被定义为两个离散分布 P_{pwq} 的混合：

$$\lambda P_{pwq}^{t-1} + (1-\lambda) P_{pwq}^t$$

其中 P_{pwq}^t

是预定义词汇表 $P_{w|C}^{t-1}$ 的分布是基于“上下文”注意分数 (复制的单词被合并) 在当前上下文C中对单词的分布，并且 λ 是在生成单词和从上下文复制之间进行选择的软开关。形式上，最后的词汇QG的分发遵循以下表格：

$$P_{w|C}^t = \text{softmax}(\frac{1}{d} \mathbf{C}_A^T \mathbf{U}_A \mathbf{R}^T \mathbf{V}_Q \mathbf{Q}_t)$$

$\mathbf{R}^T \mathbf{V}_Q \mathbf{Q}_t$

其中 softmax ；“q是注意层中定义的得分函数： q_i 和 $q_{i'}$ 分别对应于Q的 t 行和Q的 t' 行。可学习的变量是 $\mathbf{W}_1 \mathbf{R}^{1024 \times d}$, $\mathbf{W}_2 \mathbf{R}^{1024 \times d}$, $\mathbf{b}_1 \mathbf{R}^{1024}$, $\mathbf{b}_2 \mathbf{R}^{1024}$, \mathbf{P} 和 \mathbf{b}_2 ，其中V表示数字中的单词数和标量 b_2 。注意， \mathbf{W}_1 和 \mathbf{b}_2 共享由QA和QG共享，而 \mathbf{W}_2 , \mathbf{b}_1 , \mathbf{b}_2 是任务特定参数。这种共享方案有两个优点。首先，共享 \mathbf{W}_1 显著减少了参数的数量，

这有助于模型更好地捕捉注意层中的跨模态交互。其次，它强制模型在投射到词汇空间之前首先将解码后的信息从自己的空间投影到共同的潜在空间。这个共同的潜在空间可以作为连接两个任务的纽带，提供另一个传递知识的渠道。人们也可以将其视为我们在网络的上层设置的正规化，以影响培训过程。在实验部分，我们表明这种共享方案确实有助于提高性能。

对于QA, P_{tpwq} 的参数形式类似地定义

首先用 t_i 替换 q_i , r 用 C_q 替换 C_A , 用 q 替换 q_i ; q q_t
然后装备任务专用参数 W_1, W_2, b_1, b_2 ; 最后与QG共享 W 和 b 共享。

损失函数

在训练期间，模型被馈送问题 - 上下文 - 应答三元组 Q, C, A ，并且来自输出层的解码的 Q 和 A 被训练为分别类似于 Q 和 A 。至

实现这一点，我们的损失函数由两部分组成：序列转导模型中广泛使用的负对数似然丢失和惩罚重复生成文本的覆盖损失，这类似于（见，刘和曼宁）2017年）。

我们在培训中聘用教师强制策略模型从时间 $t-1$ 接收地面令牌作为输入并在时间 t 预测令牌。 具体而言, 给定三元组 Q, C, A 的目的是最小化关于所有模型参数的以下损失函数:

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad \log p(p, q) \text{ 在 } A, \\ / Q, \quad Cq' \kappa$$

QG损失

```

10000010000000000000
0000000000000000000m
0000000000000000000i
0000tO1O
“OO100000000000n

```

QG损失

QA的覆盖率下降

$\log p_{pppa_t} "a_t / A_{\text{在}}, Q,$

质量保

哪里 QC 和 s 自动控制 对应于 t 一排问题 -

上下文和分别从关注层的第二折入步骤获得的答案 \hat{r} 上
 下文关注分数： $\frac{\hat{r}}{\kappa}$ 表示两个向量的元素最小值之和；
 κ 是用于加权覆盖损失的超参数。

- 在关注层，我们开发了两步注意力，折叠到目前为止所观察到的所有信息中，以产生最终序列。第一个折叠步骤捕获问题/答案和上下文之间的交互，并将其表示为新的上下文序列。第二个折叠步骤模拟了序列到序列模型中的典型编码器 - 解码器注意机制 (Bahdanau, Cho 和 Bengio, 2014; Gehring 等, 2017)。
- 最后，在输出层，我们回收从第二个折叠步骤获得的注意力词汇，并将其用作指针生成器中的复制分布。最终的词汇因此被定义为对分布的内插

当前上下文单词和预定义大词汇表的分布。这有助于从上下文中复制单词，同时保留生成新单词的能力。

模型中的二元性

我们的模型在两个地方利用了QA和QG的二元性。

- 由于我们认为QA和QG都是序列生成问题，我们的架构是反射对称的，如图2所示。左QG部分是具有相同结构的右QA部分的镜像。这种对称性也可以在注意力计算和损失函数中找到。因此，答案模态和问题模态通过上下文模式在双向过程中连接，允许模型根据上下文推断给定对方的答案或问题。
- 我们的模型包含不同级别的QA和QG之间的共享组件。从底部开始，嵌入

任务。此外, QG中的答案编码器在QA中重用

层和上下文编码器总是在两者之间共享任务。此外，QG中的答案编码器在QA中重用

用于生成答案序列，反之亦然。在之上
在指针生成器中，QA和QG在最终投影到词汇空间之前
共享相同的潜在空间。利用问题和答案之间的循环一
致性

在不同层面规范培训过程，帮助模型，以找到每种模态的更一般的表示。

实验结果

实施细节

由于所提出的模型中的每个组件都是可区分的

模型中的注意事项

我们的模型在三个地方利用了注意机制。

- 在编码层中，我们将自我关注作为捕获远程依赖关系的一种方式。与LSTM相比，信号在网络中必须经过的最大路径长度在自我关注方面要短得多，这使得学习远程依赖性变得更加容易。在编码上下文时，我们允许前向和后向信号在自注意网络中遍历。当自我注意编码问答时，我们使用掩码防止反向信号保留自回归属性。

参数可以通过反向传播进行训练。我们使用“fan-avg”策略随机初始化参数，即从均匀分布中采样，其宽度是输入和输出连接的平均数量 (Glorot和Bengio 2010)。我们筛选出来，出现在数据小于5倍的话，并与<UNK>令牌，产生一个具有90000个字的词汇替换它们。我们修改了嵌入这个词预训练的256维GloVe单词向量 (Pennington, Socher和Manning 2014)。字符嵌入的大小为66 256，并在训练期间学习。Dropout主要应用于编码层，保持率为

0.9。覆盖损失权重 κ 为1.0。我们使用Adam优化器 (Kingma和Ba 2014) 来最小化损失函数。通过限制其 f_2 - 范围小于或等于5.0来剪切梯度。利用逆指数函数将学习率从零增加到0.001

然后修复剩下的训练。批量大小根据经验设置为16。在测试期间，我们分别对QA和QG进行自回归解码。时间t的输入是它

自己生成的单词来自前面的步骤。解码当模型遇到第一个<END>或序列包含超过100个单词时终止。

数据集

我们的实验在四个数据集上进行：SQuAD (Ra-jpurkar等人2016)，MSMARCO (Nguyen等人2016)，Wik-iQA (Yang, Yih和Meek 2015) 和TriviaQA (Joshi等人2017)。我们对MSMARCO进行二次采样以加速实验。所有数据都预处理为SQuAD格式，可在线获取。表1总结了实验中使用的训练和测试数据的统计数据。

数据集	#培养	#测试	n	m	k
队	86,821	5,928	117.1	10.1	3.1
MSMARCO	120,000	24,000	405.8	4.1	12.3
TriviaQA	120,000	24,000	631.9	12.2	1.5
WikiQA	873	126	471.2	6.4	24.1

表1: 数据集摘要。统计数据包括培训/测试Q, A, C三胞胎的数量; 上下支 (n), q 问题 (m) 和答案 (k) 中的平均单词数。

评估指标

由于问题和答案都是在我们的模型中生成的，我们采用BLEU-1, 2, 3, 4 (Papineni et al. 2002) 和 Meteor (Denkowski和Lavie 2014) 得分来自机器翻译, ROUGE-L来自文本摘要 (林2004) 评估一代的质量。这些度量是通过将机器生成的文本与基于单词之间的精确, 词干, 同义词和释义匹配的一个或多个人类生成的引用对齐来计算的。较高的分数是优选的, 因为它表明更好地与groundtruth对齐。

双重学习的有效性

表2总结了我们的模型与单学习对应 (单声道), JointQA (Wang, Yuan和Trischler 2017) (jqa), 没有强化学习的多视角QG (Song, Wang和Hamza 2017) 的表现 (mpqg) 和一个简单的序列到序列模型与注意机制 (s2s)。jqa和mpqg都与我们的工作有着共同的目标, 即利用QA和QG的双重性来相互改进。我们通过屏蔽DAANET中所有与任务无关的结构和参数来实现mono。可以观察到, DAANET明显优于单一学习对手, 平均利润率为5.7pp在Rouge-L和5.3u的Bleu-4 QA。它也胜过在大多数情况下, 最先进的联合模型。这个结果为将二元性纳入模型结构以提高阅读理解能力提供了有力的证据。人们还可能注意到, DAANET的改进在QA中比在QG中更为显着。这是因为有许多方法可以在QG中传达相同的问题, 例如通过替换

因此, 基于字面匹配的度量可能低估了生成的问题的质量。我们将调查更好的QG绩效指标作为未来的工作。

		QA				QG					
		dAAN _{mono}	jqa	mpqg	s2s	DAAN _{mono}	jqa	mpqg	s2s		
隊	B1	37.28	29.62	26.66	14.65	15.47	33.95	31.28	34.48	29.92	28.73
	B2	32.66	25.37	22.14	11.17	11.44	18.99	16.23	18.84	14.82	13.01
	B3	29.31	22.23	18.82	8.74	8.74	12.35	9.92	12.09	8.64	6.73
	B4	26.38	19.56	16.11	6.93	6.75	8.71	6.64	8.32	5.49	3.74
	RL	43.85	37.38	41.43	26.96	24.21	35.58	33.32	34.31	32.24	30.92
公	23.21	19.73	21.70	12.64	8.95	14.18	12.52	13.89	11.70	11.41	
吨											
MSMARCO	B1	44.98	41.20	41.26	35.51	29.41	58.45	59.04	56.28	49.42	53.20
	B2	38.09	34.70	34.53	28.51	22.84	45.65	46.14	43.64	37.25	40.01
	B3	34.83	31.62	31.37	25.43	19.66	35.98	37.47	34.14	28.40	30.31
	B4	32.82	29.70	29.42	23.60	17.67	28.79	29.18	26.99	22.10	22.86
	RL	44.25	42.23	41.30	35.23	32.75	57.90	58.19	56.13	50.90	53.57
公	22.60	22.28	22.06	16.93	15.33	29.52	29.78	27.94	24.14	25.73	
吨											
TriviaQA	B1	62.13	48.47	47.67	41.41	50.14	45.61	43.26	37.57	34.07	31.93
	B2	59.40	45.13	44.54	37.91	45.07	33.25	29.91	23.49	19.91	17.43
	B3	53.41	39.68	40.90	31.45	32.88	26.60	22.81	16.91	13.23	10.45
	B4	40.93	29.10	34.54	20.10	19.18	22.32	18.33	11.48	9.46	6.72
	RL	61.97	47.89	46.86	41.94	51.87	45.31	42.69	38.26	34.57	32.81
公	38.80	29.64	28.45	24.79	31.38	21.35	19.67	16.64	14.54	13.43	
吨											
WikiQA	B1	37.31	36.47	31.74	33.53	7.49	14.13	13.49	14.45	15.53	11.06
	B2	31.76	30.89	26.31	28.81	2.82	7.67	7.28	7.76	7.78	6.37
	B3	29.74	28.78	24.33	27.02	1.21	2.89	0.00	0.00	0.00	0.00
	B4	28.61	27.56	23.25	25.98	0.56	1.71	0.00	0.00	0.00	0.00
	RL	37.32	35.50	29.91	36.41	10.74	18.86	19.61	19.24	20.03	18.45
公	19.50	18.90	14.95	20.08	3.47	5.71	5.71	5.63	5.47	4.27	
吨											

表2: 四个数据集上的问答和问题生成结果。B代表Bleu-1, 2, 3, 4; RL是Rouge-L; Mt是流星。越高越好; 最好的大胆。

消融研究

表3总结了DAANET的性能及其在我们的SQuAD测试集上的消融, 其中我们比较了编码器的不同组成, 有无注意机制和不同的共享方案。首先, 人们可以观察到LSTM是编码中最重要的部分, 没有它QA时性能下降约10pp, QG下降2.5pp。自我关注也有助于改善质量保证3pp和QG 0.6pp。接下来, 删除上下文关注会导致灾难性的表现。这符合我们的直觉因为没有上下文关注意味着答案生成独立于给定问题并且仅取决于给定的上下文。最后, 表3的最后部分证明了所提出的共享方案的有效性。

案例分析

表4列出了DAANET和单一学习模型的一些生成问题和答案。根据上下文给出黄金对应物生成问题或答案。在前两个样本中, DAANET完美地工作, 而

切除	QA				QG			
	布鲁-1	布鲁-4	胭脂-L	流星	布鲁-1	布鲁-4	胭脂-L	流星
没有LSTM的编码器	20.80	11.22	30.84	14.84	31.60	6.30	32.75	12.38
编码器没有自我关注	33.24	22.69	39.90	21.31	33.56	8.14	34.10	13.65
没有背景 - 关注	4.96	0.38	4.61	1.81	25.74	2.05	26.99	8.47
没有复制机制	8.24	0.94	11.23	3.99	26.32	2.22	27.29	8.83
非共享问答编码器	29.82	19.43	37.75	19.47	33.23	8.51	35.26	13.86
非共享上下文编码器	22.90	13.45	29.33	14.67	33.37	8.17	34.94	13.75
非共享输出词汇表投影	32.64	21.46	40.37	20.77	33.71	8.15	34.77	13.85
DAANET	37.28	26.38	43.85	23.21	33.95	8.71	35.58	14.18

表3: DAANET的性能及其对SQuAD的消融。 越高越好; 最好的大胆

单一学习模型无法提供所需的输出。 在第三个示例中, DAANET生成的问题与groundtruth相比更具可读性。 根据经验, 我们发现DAANET提出的问题在语义上与引用的问题类似, 但措辞不同。 尽管人类可能认为这些是好的案例, 但DAANET在所有基于对齐的评估指标下仍然在QG上得分很低。 这解释了QG的得分总体上低于QA, QG的改善不如QA。 我们在补充文件中附上了更多生成的样本。

结论

我们提出了双问答网络, 这是一种双向神经序列转导模型, 它解决了机器阅读理解的问题和问题生成。 我们通过在多层面共享本地结构来利用QA和QG任务的二元性。 注意机制用于多层以捕获上下文, 问题和答案的相互作用。 我们证明了我们的模型在四个公共数据集上的有效性, 为将二元性带入模型结构以提高阅读理解能力提供了有力的证据。 一个有趣的未来方向是研究多个DAANET的协作。 例如, 可以多次堆叠DAANET, 使得从较低网络生成的问题和答案被馈送到上层网络。 这种“自举”策略可以被认为是隐含的数据增加方式, 可以减轻标记的MRC数据的不足。

参考

阿尔达部, 我。 De Lacalle, ML; Maritxalar, M. ; Martinez, E. ; 和Uria, L. 2006. Arikurri: 基于语料库和nlp技术的自动问题生成器。 在智能辅导系统国际会议上, 584-594。 斯普林格。

Ba, JL; 基洛斯, JR; 和Hinton, GE 2016. 图层规范化。 的arXiv: 1607. 06450。

Bahdanau, D. ; Cho, K. ; 和Bengio, Y. 2014. 通过联合学习对齐和翻译的神经机器翻译。 的arXiv: 1409. 0473。

Bordes, A. ; 韦斯顿, J. ; 和Usunier, N. 2014. 用弱监督嵌入模型打开问题回答。 在ECML-PKDD中, 165-180。 斯普林格。

陈, D. ; Fisch, A. ; 韦斯顿, J. ; 和Bordes, A. 2017。

背景在10世纪的过程中, 挪威战争带最初破坏性地侵入法国河流, 演变成更为永久的营地, 包括当地妇女和个人财产。 诺曼底公国, 始于911年作为一个封地, 是由国王之间的圣克萊尔河畔埃普特条约建立的.....

问题Nor诺曼底公国什么时候成立?

回答 911

DAANET: 诺曼底公国何时开放?

911

单 z 诺曼底公国什么时候开始?

10世纪

背景..... 墨尔本有许多博物馆, 艺术画廊和剧院的所在地, 也被称为“澳大利亚体育之都”。 墨尔本板球场是澳大利亚最大的体育场, 也是1956年夏季奥运会和2006年英联邦运动会的东道主.....

问题Australia澳大利亚最大的体育场是什么?

答案 墨尔本板球场

DAANET: 澳大利亚最大的体育场是什么?

墨尔本板球场

单 z 什么是最大的国家生产总值 (AFL)?

墨尔本板球场

Context Sky UK Limited (前身为英国天空广播公司或BSkyB) 是一家服务于英国的英国电信公司。 Sky为英国的消费者和企业提供电视和宽带互联网服务以及固定电话服务。 它是英国最大的付费电视广播公司, 截至2015年拥有1100万客户.....

问题BSkyUK Limited以前叫什么名字?

答案 英国天空广播

DAANET sky什么是天空的名字有限?

英国天空广播或BSkyB

单 z 加拿大最大的付费电视服务是什么?

BSkyB

表4: DAANET和mono的选定输出。 黄色文字与问题有关; 绿色文本与答案有关。

阅读维基百科以回答开放域问题。 在ACL中, 1870-1879。

- Denkowski, M. 和Lavie, A. 2014. Meteor universal: 针对任何目标语言的语言特定翻译评估。在统计机器翻译@ EACL。
- Dhingra, B. ; Pruthi, D. ; 和Rajagopal, D. 2018. 简单有效的半监督问答。在NAACL-HLT中。杜, X. 邵, J. ; 和Cardie, C. 2017. 学习问: 用于阅读理解的神经问题生成。在ACL中。Duan, N. ; 唐, D. ; 陈, P. ; 和周, M. 2017年。问答的问题生成。在EMNLP, 866-874。
- Gehring, J. ; Auli, M. ; Grangier, D. ; Yarats, D. ; 和Dauphin, Y. 2017. 卷积序列到序列学习。在ICML中。Gers, FA; Schmidhuber, J. ; 和康明斯, FA 2000. 学会忘记: 用lstm进行持续预测。神经计算12: 2451-2471。
- Glorot, X. 和Bengio, Y. 2010. 了解训练深度前馈神经网络的难度。在AISTATS, 249-256。
- Graves, A. 2013. 使用递归神经网络生成序列。的arXiv: 1308.0850。
- Heilman, M. 2011. 自动事实问题生成文本。胡, M. ; 彭, Y. ; 黄, Z. ; 邱, X. Wei, F. ; 和周, M. 2018. 用于机器读数的增强型记忆读取器理解。在IJCAI, 4099-4106。
- Huang, H.-Y. ; 朱, C. ; 沉, Y. ; 和Chen, W. 2017. Fusionnet: 通过充分意识到的注意力与机器理解的应用融合。的arXiv: 1711.07341。
- Joshi, M. ; Choi, E. ; 焊接, DS; 和Zettlemoyer, L. 2017. Triviaqa: 一个用于阅读理解的大规模远程监督挑战数据集。在ACL中。加拿大温哥华: 计算语言学协会。
- Kingma, DP和Ba, J. 2014. Adam: 一种随机优化的方法。CoRR abs / 1412.6980。
- 林, C.-Y. 2004. Rouge: 自动评估摘要的软件包。文本摘要分支出来。
- 刘, M. ; 卡尔沃, RA; 和Rus, V. 2010. 文献综述写作支持的自动问题生成。在智能辅导系统国际会议上, 45-54。斯普林格。Luong, T. ; Pham, H. ; 和Manning, CD 2015. 有效基于注意力的神经机器翻译的方法。在EMNLP, 1412-1421。
- Nallapati, R. ; 周, B. ; dos Santos, CN; aglarGuðlehnre; 和Xiang, B. 2016. 使用序列到序列rnns及其后的抽象文本摘要。在CoNLL。
- Nguyen, T. ; 罗森伯格, M. ; 宋, X. ; 高, J. ; Tiwary, S. ; Majumder, R. ; 和邓, L. 2016年。马克女士: 人类生成的机器阅读理解数据集。
- Papineni, K. ; Roukos, S. ; 沃德, T. ; 和朱, W.-J. Bleu: 一种自动评估机器翻译的方法。在ACL中, 311-318。计算语言学协会。
- Pennington, J. ; Socher, R. ; 和Manning, C. 2014. Glove: 用于单词表示的全球向量。在EMNLP, 1532-1543。
- Rajpurkar, P. ; 张, J. ; Lopyrev, K. ; 和Liang, P. 2016. Squad: 10万多个关于机器理解文本的问题。在EMNLP。
- Sachan, M. 和Xing, E. 2018. 自我训练, 共同学习提问和回答问题。在NAACL, 第1卷, 629-640。

看到。; 刘, PJ; 和Manning, CD 2017. 重点: 使用指针生成器网络进行汇总。 在ACL中。

Seo, M.; Kambhavi, A.; Farhadi, A.; 和Hajishirzi, H. 2016. 机器理解的双向注意力流动。 的arXiv: 1611.01603。

塞尔班, 第四; Garcia-Dura'n, A.; Gulcehre, C.; 安, S.; Chandar, S.; Courville, AC; 和Bengio, Y. 2016. 使用递归神经网络生成仿真问题: 30m仿真问题 - 答案语料库。 在ACL中。

沉, Y.; Huang, P.-S.; 高, J.; 和陈, W. 2017. 理由: 学会停止阅读机器理解。 在ACM SIGKDD, 1047-1055。 ACM。

宋, L.; 王, Z.; 和Hamza, W. 2017. 一个统一的基于查询的生成和问答的生成模型。 的arXiv: 1709.01058。

Srivastava, RK; Greff, K.; 和Schmidhuber, J. 2015. 公路网。 的arXiv: 1505.00387。

Sutskever, 我。 Vinyals, O.; 和Le, QV 2014. 用神经网络进行序列学习的序列。 在NIPS, 3104-3112。

Tan, C.; Wei, F.; 杨, N.; 杜, B.; Lv, W.; 和周, M. 2017. S-net: 从答案提取到机器阅读理解的答案生成。 的arXiv: 1706.04815。

唐, D.; Duan, N.; 秦, T.; Yan, Z.; 和周, M. 2017. 问题回答和问题生成是双重任务。 的arXiv: 1706.02027。

唐, D.; Duan, N.; Yan, Z.; 张, Z.; 太阳, Y.; 刘, S. Lv, Y.; 和周, 硕士, 2018年。 学习合作回答问题。 在NAACL, 第1卷, 1564-1574。

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; 琼斯, L.; Gomez, AN; 凯撒, Ł.; 和Polosukhin, I. 2017年。 注意力就是你所需要的。 在NIPS, 5998-6008。

Vinyals, O.; 凯撒, Ł.; Koo, T.; Petrov, S.; Sutskever, 我。 和Hinton, G. 2015. 语法作为外语。 在NIPS, 2773-2781。

Wang, S. 和Jiang, J. 2016. 使用match-lstm和答案指针的机器理解。 的arXiv: 1608.07905。

王, W.; 杨, N.; Wei, F.; Chang, B.; 和周, M. 2017. 用于阅读理解和问答的门控自匹配网络。 在ACL中, 第1卷, 189-198。

王, T.; 袁, X. 和Trischler, A. 2017年。 问答和问题生成的联合模型。 arXiv preprint arXiv: 1706.01450。

杨, Z.; 胡, J.; Salakhutdinov, R.; 和Cohen, WW 2017. 生成域自适应网络中的半监督qa。 在ACL中。 杨, Y. Yih, W.-t.; 和Meek, C. 2015. Wikiqa: 一个挑战开放域问答的数据集。 在EMNLP, 2013-2018。

Yu, AW; Dohan, D.; Luong, M.-T.; 赵,

R.; 陈, K.; Norouzi, M.; Qa 2018, Le, QV 2018. Qanet: 将局部卷积与全球自我关注相结合, 以便阅读理解。 的arXiv: 1804.09541。

袁, X. 王, T.; aglarGulcehre; Sordani, A.; Bachman, P.; Subramanian, S.; 张, S.; 和Trischler, A. 2017. 通过文本到文本神经问题生成的机器理解。 在RepANP @ ACL中。