

引用格式: 万静, 郭雅志. 基于多段落排序的机器阅读理解研究[J]. 北京化工大学学报(自然科学版), 2019, 46(3): 93–98.

WAN Jing, GUO YaZhi. Machine reading comprehension based on multi-passage ranking[J]. Journal of Beijing University of Chemical Technology (Natural Science), 2019, 46(3): 93–98.

基于多段落排序的机器阅读理解研究

万 静 郭雅志

(北京化工大学 信息科学与技术学院, 北京 100029)

摘 要: 针对多段落的机器阅读理解问题, 在双向注意力流 (BiDAF) 模型的基础上, 结合双向长短期记忆网络 (BiLSTM) 和 self-attention 机制构建了多段落排序 BiDAF (PR-BiDAF) 模型, 利用该模型定位答案所在的段落, 然后在预测段落中寻找最终答案的始末位置。实验结果表明, 相较于 BiDAF 模型, 本文提出的 PR-BiDAF 模型的段落选择正确率、BLEU4 指标及 ROUGE-L 指标分别提高了约 13%、6% 和 4%。

关键词: 机器阅读理解; 双向注意力流 (BiDAF) 模型; self-attention 机制

中图分类号: TP391 **DOI:** 10.13543/j.bhxbzr.2019.03.014

引 言

近年来, 自然语言处理技术发展迅速, 在各个领域都取得了引人注目的成绩。机器阅读理解作为自然语言处理的核心任务之一, 也越来越受到各方关注。机器阅读理解是指通过让机器阅读文本和问题, 从而得到一段简洁、有效的文字作为该问题的答案。机器阅读理解可以充分体现机器的智能水平, 因此提升机器阅读理解水平对于推动人工智能发展具有重要意义。

Rajpurkar 等^[1]于 2016 年推出英文的机器阅读理解 SQuAD 数据集, 该数据集通过众包方式构建, 包含 10 万个问题和 536 篇维基百科文章, 因其数量、质量等方面的优势吸引了众多研究者的目光。但在中文阅读理解领域, 权威的大型数据集较少, 中文阅读理解研究进展缓慢; 此外大多数研究热点都集中在单段落抽取式的阅读理解类型上, 即从文本中抽出一个片段作为相应的答案。但在真实场景中往往需要从多个段落中找到最相关的答案, 因此针对多段落机器阅读理解问题有待进一步研究。

与此同时, 深度学习技术在自然语言处理 (NLP) 方面取得的良好效果为解决阅读理解问题提

供了新的思路, 其中被广泛应用的模型是模型 Seq2Seq^[2]。相比基于语法句法分析和特征工程的传统方式, 基于 Seq2Seq 的深度学习模型在机器阅读理解方面可以取得更好的效果, 并且减轻了人工选取和构造特征方面的压力。

由 Seo 等^[3]提出的双向注意力流模型 (bi-directional attention flow, BiDAF) 属于 Seq2Seq 框架, 它在 SQuAD、CNN/DailyMail^[4]等面向单段落机器阅读理解的数据集上获得了较好的表现。BiDAF 采用双向的注意力机制, 提取到的蕴含向量包含更多的关联性信息, 但在多段落的机器阅读理解问题方面, 还需要考虑不同段落与问题之间的整体相似度以及段落与段落之间的信息关联, 使得到的预测结果更精准。

本文针对 BiDAF 模型对多段落信息考虑不充分的问题, 提出基于多段落排序机制的 BiADF (passage rank-BiDAF, PR-BiDAF) 模型。PR-BiDAF 使用双向长短期记忆网络 (bi-directional long short-term memory network, BiLSTM) 对问题和段落分别编码, 通过 self-attention 机制提取问题和段落中具有代表性的重要信息, 再将段落和问题的整体向量输入 softmax 层得到段落的关联度排序, 并输出得分最高的段落; 在预测段落中, 通过 BiDAF 模型的双向注意力机制融合段落和问题的编码信息, 最终预测和生成答案片段。最后通过在中文多段落阅读理解数据集 DuReade 上的多组实验对比证明了本文方法的有效性。

收稿日期: 2018-08-14

基金项目: 国家自然科学基金 (51577006)

第一作者: 女, 1975 年生, 副教授

E-mail: wanj@mail.buct.edu.cn

1 基于多段落排序的阅读理解模型

1.1 模型结构

如图 1 所示, PR - BiDAF 模型由输入层、BiDAF、passage rank 三部分组成。输入层负责对段落和问题进行词向量嵌入和编码处理,并作为后续操作的输入。输入层中的多段落阅读理解数据存储形式为 $\langle P_1, P_2, \dots, P_n, Q, A \rangle$,其中 P_i 代表 n 个待阅读的段落, Q 代表问题, A 为答案。BiDAF 模型用于抽取单段落答案,按照顺序可以分为双向注意力计算层(bi-attention)、解码层(decode)和答案预测层(answer predict)。BiDAF 模型主要负责处理单个段落内的答案预测,通过段落和问题的输入,得到 start scores vector 和 end scores vector 分别代表的答案在段落内的始末位置,并选取乘积值最大区域作为预测的答案片段。passage rank 分为 self-attention 提取信息层 (self-attention)^[5]、全连接层 (fully connect)

和 softmax 层^[6]。self-attention 层通过提取句子内部重要信息来表征该语句,再经过全连接层得到第 i 个段落 P_i 与问题 Q 的匹配度得分 g_i ,然后把 N 个段落的得分输入 softmax 层,输出得分最高的段落作为模型预测的正确段落,并将该段落的答案片段作为输出结果。

1.2 输入层

模型的输入数据包括一个问题和多个相关段落。目前在自然语言处理中通常使用词向量技术来表示文字,然后通过编码层处理后生成蕴含上下文语义信息的表示向量。本文选取 Glove 预训练词向量^[7],将未出现在该词向量中的文字进行随机初始化处理。

为了解决长距离依赖以及上下文信息融合问题,编码部分采用 BiLSTM 网络结构,将语句相反的两种顺序分别输入长短期记忆网络 (LSTM) 中编码,然后进行拼接,这样得到的表示向量就能蕴含上

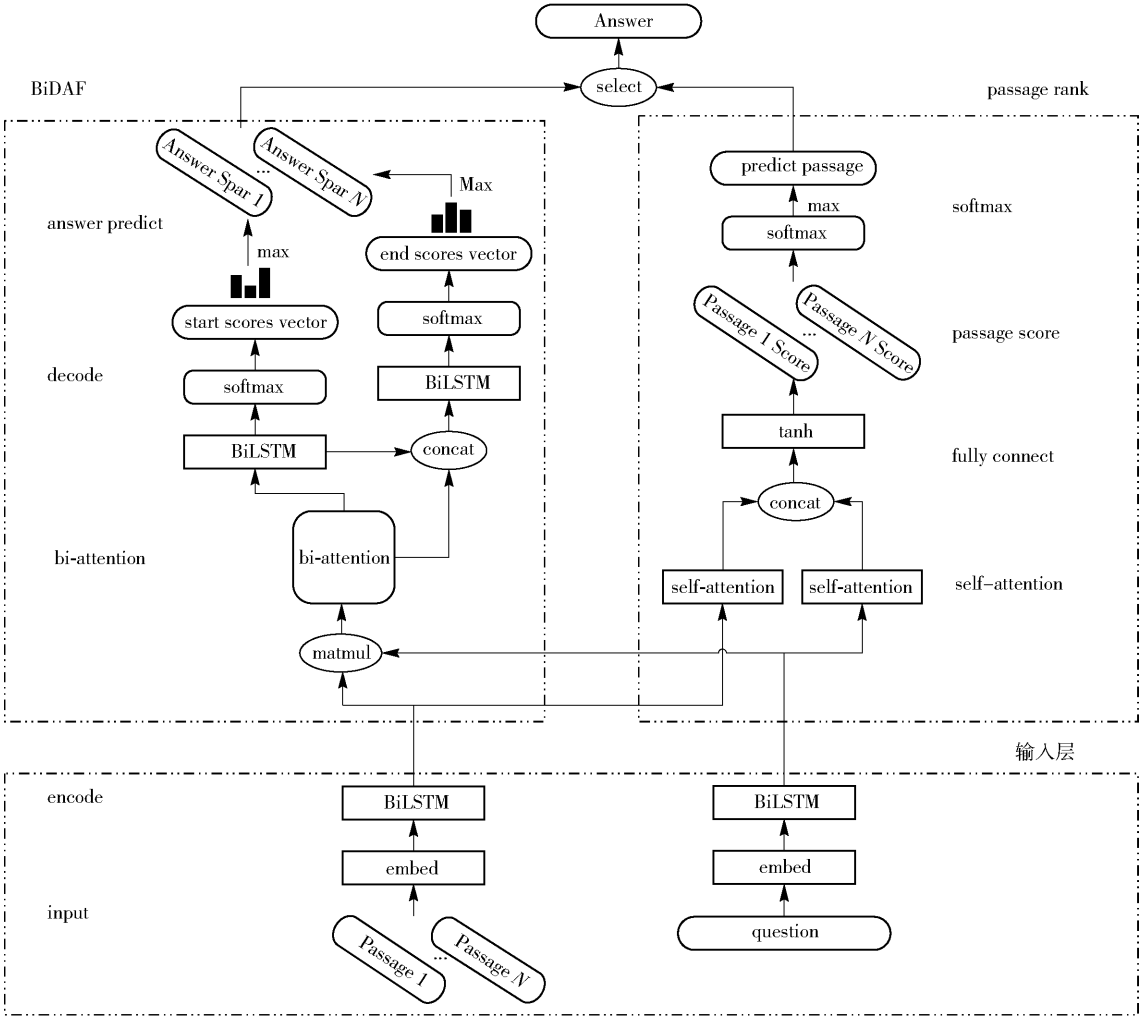


图 1 PR-BiDAF 模型

Fig. 1 PR-BiDAF model

下文信息。BiLSTM 的公式定义如下^[8]

$$\begin{cases} \mathbf{h}_t^+ = \text{LSTM}(\mathbf{h}_{t-1}, \mathbf{e}_t) \\ \mathbf{h}_t^- = \text{LSTM}(\mathbf{h}_{t+1}, \mathbf{e}_t) \\ \mathbf{h}_t = [\mathbf{h}_t^+, \mathbf{h}_t^-] \end{cases} \quad (1)$$

式中, \mathbf{e}_t 代表 t 时刻的输入向量, \mathbf{h}_t^+ 、 \mathbf{h}_t^- 分别代表 t 时刻正向和反向的隐含层向量, \mathbf{h}_t 为 t 时刻正反隐含层向量拼接后的结果。

1.3 基于 BiDAF 的单段落答案抽取

对于单段落的答案抽取, 定位答案范围对生成的答案质量有关键性影响。BiDAF 模型的主要特点是在特征提取阶段引入文章和问题双向注意力机制, 得到原文关于问题的表示向量后用双向 LSTM 进行语义信息聚合, 最后通过解码层的输出向量来分别预测答案起始位置和结束位置。

模型单轮训练的输入包括单个问题和多篇文章两部分。经过分词等预处理后, 以文章和问题的词向量表示为输入, 经过 BiLSTM 的编码分别得到文章的编码 \mathbf{H} 和问题的编码 \mathbf{U} 。由式(2)计算 \mathbf{H} 和 \mathbf{U} 的相似度矩阵 \mathbf{S}_{ij} , 再基于该矩阵计算双向的注意力机制。

$$\begin{cases} \mathbf{S}_{ij} = \alpha(\mathbf{H}_{:,i}, \mathbf{U}_{:,j}) \\ \alpha(\mathbf{h}, \mathbf{u}) = \mathbf{W}_s^T[\mathbf{h}; \mathbf{u}; \mathbf{h} \odot \mathbf{u}] \end{cases} \quad (2)$$

式中, $\mathbf{H}_{:,i}$ 表示 \mathbf{H} 的第 i 个列向量, $\mathbf{U}_{:,j}$ 表示 \mathbf{U} 的第 j 个列向量, \mathbf{W}_s^T 表示可训练的权重向量, \odot 表示逐个元素相乘, $[\cdot; \cdot]$ 表示按行连接向量。

利用文章到问题方向(P2Q)的注意力机制为每个段落中的词找出问题中与其最相关的词。对 \mathbf{S} 按列计算 softmax 得到注意力向量 \mathbf{a}_i , 然后将 \mathbf{a}_i 与 \mathbf{U} 相乘得到矩阵 $\tilde{\mathbf{U}}$, 公式如下

$$\begin{cases} \mathbf{a}_i = \text{softmax}(\mathbf{S}_{:,i}) \\ \tilde{\mathbf{U}}_{:,i} = \sum_j \mathbf{a}_{ij} \mathbf{U}_{:,j} \end{cases} \quad (3)$$

式中, $\mathbf{S}_{:,i}$ 表示 \mathbf{S} 的第 i 个列向量, $\tilde{\mathbf{U}}_{:,i}$ 表示输出矩阵的第 i 列向量。

利用问题到文章方向(Q2P)的注意力机制对每个问题中的词找出段落中与其最相关的词, 再对相似度矩阵 \mathbf{S} 按列提取最大值, 将所有最大值经过 softmax 处理后, 由式(4)计算得到 $\tilde{\mathbf{h}}$

$$\begin{cases} \mathbf{b} = \text{softmax}(\max_{\text{col}}(\mathbf{S})) \\ \tilde{\mathbf{h}} = \sum_i \mathbf{b}_i \mathbf{H}_{:,i} \end{cases} \quad (4)$$

式中, \max_{col} 表示按列取最大值, \mathbf{b} 为注意力向量, $\tilde{\mathbf{h}}$ 为特征向量。最后将 $\tilde{\mathbf{h}}$ 按列重复 T 次得到特征矩阵 $\tilde{\mathbf{H}}$ 。

将 \mathbf{H} 、 $\tilde{\mathbf{U}}$ 、 $\tilde{\mathbf{H}}$ 3 个矩阵拼接起来, 并分别经过两个 BiLSTM 解码层后得到特征向量 \mathbf{G} 、 \mathbf{M} 、 \mathbf{M}^2 , 再通过公式(5)得到代表答案初始位置和结束位置可能性大小的向量 $\mathbf{p}_{\text{start}}$ 和 \mathbf{p}_{end}

$$\begin{cases} \mathbf{p}_{\text{start}} = \text{softmax}(\mathbf{W}_{\text{start}}^T[\mathbf{G}; \mathbf{M}]) \\ \mathbf{p}_{\text{end}} = \text{softmax}(\mathbf{W}_{\text{end}}^T[\mathbf{G}; \mathbf{M}^2]) \end{cases} \quad (5)$$

式中 $\mathbf{W}_{\text{start}}^T$ 和 $\mathbf{W}_{\text{end}}^T$ 分别表示 $\mathbf{p}_{\text{start}}$ 和 \mathbf{p}_{end} 的权重矩阵。最后通过循环比较所有位置可能性的大小, 找到乘积最大的两个坐标点, 从而得出预测的起始和结束范围。

1.4 基于多段落排序机制的预测模型

在多段落阅读理解问题中, 最终的答案来自某一最相关的段落, 所以计算段落与问题的关联度并选取正确段落是提高答案匹配度的关键所在。

本文首先进行段落与问题的整体关联度匹配计算, 进而转化成正确段落和非正确段落间的二分类问题, 通过多段落排序模型(passage rank)训练得到预测结果, 直接在该段落内寻找最终的答案片段。在提取语句内部特征时引入 self-attention 机制, 目的是提取句子内部不同重要程度的特征, 从而得到句子的语义蕴含向量。self-attention 通过比较语句内部所有词语间的重要程度, 生成该语句的注意力向量。其机制定义如下

$$\begin{cases} \mathbf{h} = (\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3, \dots, \mathbf{h}_n) \\ \mathbf{A} = \text{softmax}(\mathbf{W}_1 \tanh(\mathbf{W}_2 \mathbf{h}^T)) \\ \mathbf{R} = \mathbf{A} \mathbf{h} \end{cases} \quad (6)$$

式中, \mathbf{h}_i 为隐含层向量矩阵, \mathbf{A} 为注意力向量, \mathbf{W}_1 和 \mathbf{W}_2 为 \mathbf{A} 的权重矩阵, \mathbf{R} 为经过注意力权重计算后的表示向量。

所有段落和问题经过 self-attention 层提取特征后, 分别得到第 i 个段落的表示向量 \mathbf{r}_i^p 和问题表示向量 \mathbf{r}^q ; 将每个段落表示向量分别与问题表示向量同时输入全连接层, 得到段落相关度评分; 最后选取评分最高的段落作为段落的预测结果, 即从该段落中寻找答案。passage rank 部分的定义如下

$$\begin{cases} \mathbf{g}_i = \mathbf{V}_g^T(\tanh(\mathbf{W}_g[\mathbf{r}_i^p, \mathbf{r}^q])) \\ P_{\text{max}} = \max(\text{softmax}(\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_n)) \end{cases} \quad (7)$$

式中, $[\mathbf{r}_i^p, \mathbf{r}^q]$ 为 \mathbf{r}_i^p 和 \mathbf{r}^q 的拼接向量; \mathbf{V}_g 和 \mathbf{W}_g 为 \mathbf{g}_i 的权重矩阵; \tanh 表示以 \tanh 为激活函数的全连接

神经网络; g_i 表示第 i 个段落的评分($i = 1, 2, \dots, n$), 其值越大, 与问题的关联度越高, 反之越低。最后经过 \max 选出匹配分数最高的段落。

在模型训练阶段, 将每个问题的答案所在的正确段落作为模型的训练标签。由于该问题实际上等价于二分类问题, 因此模型选择交叉熵作为损失函数, 同时引入 Dropout 机制^[9] 和 L2 正则项^[10] 来避免过拟合。交叉熵损失函数定义如下

$$l = - \sum_{i=1}^k [y_i \lg g_i + (1 - y_i) \lg (1 - g_i)] \quad (8)$$

式中, y_i 取 1 或 0, 分别代表是否为正确的段落。

2 实验部分

2.1 实验数据集

为了验证本文提出的 PR-BiDAF 模型的有效性, 从百度发布的大规模中文阅读理解数据集 DuReader^[11] 中选取共 30 万条数据, 其中包括训练集 27 万条, 开发集 1 万条, 测试集 2 万条。

数据分为 Search 和 Zhidao 两大类, 分别来自百度搜索和百度知道两个真实应用场景。内容为百度用户查询的真实问题, 每个问题对应 5 个候选文档及其人工整理的优质答案。为了提供更丰富的问题种类, DuReader 数据集中的问题按答案类型被分为 3 类: Entity (实体)、Description (描述) 和 YesNo (是非)。对于实体类问题, 其答案一般包含具体的一个或多个词语; 描述类问题的答案包含一段连续性的文字描述; 是非类问题的答案为“是”或者“否”。数据集样例见表 1。

表 1 DuReader 数据集样例

Table 1 The sample of the DuReader data set

项目	描述
Question	板蓝根颗粒的功效与作用
Question Type	DESCRIPTION-FACT
Answer1	清热解毒、凉血; 用于温热发热、发斑、风热感冒、咽喉肿烂、流行性乙型脑炎、肝炎、腮腺
Answer2	板蓝根的用途不仅是治疗感冒, 板蓝根的功效与作用很多, 对多种细菌性、病毒性疾病都有较好的预防与治疗作用。
Document1	板蓝根对感冒、流感、流脑、腮腺炎、肺炎等疾病都有良好的预防和治疗效果……
:	:
Document5	当出现流鼻涕, 鼻子堵塞或头痛的症状时, 很多人都会习惯性地去药店买板蓝根, 喝两包板蓝根冲剂, 我记得小时候……

2.2 评价方法

本文选取 ROUGE-L 和 BLEU4 作为评价答案优劣的评价指标^[12], 同时利用段落选择准确率来衡量 PR-BiDAF 模型预测段落的效果。

ROUGE-L 指标参考最长公共子序列 (LCS) 来比较两个字符串 X 、 Y 的相似度, 其值越高, 两个字符串越相似。ROUGE-L 指标 (F_{lcs}) 计算公式为

$$F_{\text{lcs}} = \frac{(1 + \beta^2) R_{\text{lcs}} P_{\text{lcs}}}{R_{\text{lcs}} + \beta^2 P_{\text{lcs}}} \quad (9)$$

式中, $R_{\text{lcs}} = \frac{L(X, Y)}{m}$, $P_{\text{lcs}} = \frac{L(X, Y)}{n}$, β 为常数, R_{lcs} 和 P_{lcs} 分别表示召回率和准确率, $L(X, Y)$ 是 X 、 Y 的 LCS 长度, m 、 n 分别为 X 、 Y 的长度。

双语评价替补 (bilingual evaluation understudy, BLEU) 是一种基于 N-gram 的匹配规则, 当 $N = 4$ 时即为 BLEU4, 其值越大两个字符串的语义联系越紧密。BLEU (F_{bleu}) 的计算公式为

$$\begin{cases} F_{\text{bleu}} = f \exp \left(\sum_{n=1}^N W_n \lg p_n \right) \\ f = \begin{cases} 1, & l_c > l_r \\ e^{(1 - \frac{l_c}{l_r})}, & l_c \leq l_r \end{cases} \end{cases} \quad (10)$$

式中, $N = 4$, W_n 为各阶权重值, f 为惩罚因子, p_n 为 n 阶精度, l_c 和 l_r 分别代表两个字符串 c 和 r 的长度。

段落选择准确率计算公式如下

$$P_{\text{Accuracy}} = \frac{N_{\text{acr}}}{N_{\text{tot}}} \quad (11)$$

式中, N_{acr} 为正确预测段落的问题数量, N_{tot} 为总问题数量。

2.3 参数设置

本文考察了词向量维度和批量处理大小 (batch_size) 两个参数对实验结果的影响。词向量维度取 200、250、300、350, 均从维基百科中文语料上训练得到; batch_size 取 8、16、32。通过多组对比实验得到最佳的 batch_size 为 32, 词向量维度为 300; 对于未出现的词语进行随机初始化后也加入以上训练过程。

设置 BiLSTM 网络结构的隐含层节点为 150, 初始化学率速率为 0.001; 使用 Adam 优化算法^[13] 对权重进行更新; 为了防止过拟合化, Dropout 值设为 0.2。

2.4 实验结果及分析

分别在 Search 数据集、Zhidao 数据集以及总数

数据集 (All) 上进行实验, 并选择 BiDAF (模型 1)、型 3) 与本文方法 (模型 4) 进行对比实验, 实验结果 Match-LSTM^[14] (模型 2) 以及人类的表现 Human (模如表 2 所示。

表 2 不同模型的实验结果
Table 2 Experimental results with different models

数据集	BLEU4 指标/%				Rouge-L 指标/%				准确率/%			
	模型 1	模型 2	模型 3	模型 4	模型 1	模型 2	模型 3	模型 4	模型 1	模型 2	模型 3	模型 4
Search	23.1	23.1	55.1	26.0	31.1	31.2	54.4	32.9	36.2	36.6		50.6
Zhidao	42.5	42.5	57.1	44.9	47.5	48.0	60.7	49.2	39.3	40.0		52.3
All	31.9	31.8	56.1	37.7	39.0	39.2	57.4	43.0	38.7	39.1		51.5

从表 2 中可以看出, 相较于模型 1 和模型 2, 本文模型在最好情况下 BLEU4 指标提升 6% 左右, Rouge-L 指标提升 4% 左右, 段落选择准确率最高提升 13%; 但是与 Human 模型的结果相比, 3 个指标在各数据集上的表现仍有较大差距, 最高相差 11.9%, 原因是本文模型仍基于字符层级抽取特征, 不能像人类一样真正理解段落和问题。

另外, 为了验证段落选择准确率对答案选取质量的影响, 利用正确段落标签替代 passage-rank 模型, 并将其作为预测段落进行另一组对照实验 (命名为 Label-BiDAF), 实验结果如表 3 所示。

表 3 对照实验结果
Table 3 Contrast experimental results

模型	BLEU4 指标/%	Rouge-L 指标/%	准确率/%
BiDAF	31.8	39.0	38.7
PR-BiDAF	37.7	43.0	51.5
Label-BiDAF	48.5	50.8	100.0

表 3 表明, 段落预测准确率达到 100% 时, 本文模型的两项指标相比基线模型均有较大提升。这一结果同样证明, 对段落先打分排序再选出正确段落的方法对得到正确答案有较大帮助。

3 结束语

本文在 BiDAF 模型的基础上, 通过对段落与问题关联匹配度进行排序, 选取最有可能产生答案的段落, 同时引入 self-attention 机制提取语句内部重要特征作为语句表示向量, 最后在多段落中文数据集 DuReader 上进行了实验验证。实验结果表明, 相比 BiDAF 基线模型, 本文方法的段落选择准确率, BLEU4 评价指标及 ROUGE-L 指标均有不同程度提高, 表明本文模型在机器阅读理解问题中取得了较好的效果; 但与人类的表现仍有较大差距, 因此本

文模型还有较大提升空间。下一步工作将考虑引入多种注意力机制算法提取特征, 提升段落选择准确率, 从而更好地实现多段落机器阅读理解。

参考文献:

[1] RAJPURKAR P, ZHANG J, LOPYREV K, et al. SQuAD : 100, 000 + questions for machine comprehension of text [C] // Conference on Empirical Methods in Natural Language Processing. Austin, 2016: 2383 – 2392.

[2] SERBAN I V, SORDONI A, LOWE R, et al. A hierarchical latent variable encoder-decoder model for generating dialogues [C] // Association for the Advance of Artificial Intelligence. San Francisco, 2017: 3295–3301.

[3] SEO M, KEMBHAVI A, FARHADI A, et al. Bidirectional attention flow for machine comprehension [C] // The 5th International Conference on Learning Representations. Toulon, 2017.

[4] CHEN D, BOLTON J, MANNING C D. A thorough examination of the CNN/daily mail reading comprehension task [C] // Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics. Berlin, 2016: 2358–2367.

[5] LIN Z H, FENG M W, SANTOS C N D, et al. A structured self-attentive sentence embedding [C] // The 5th International Conference on Learning Representations. Toulon, 2017.

[6] MEMISEVIC R, ZACH C, HINTON G, et al. Gated softmax classification [C] // International Conference on Neural Information Processing Systems. Vancouver, 2010: 1603–1611.

[7] PENNINGTON J, SOCHER R, MANNING C D. GloVe: global vectors for word representation [C] // Conference on Empirical Methods in Natural Language Processing. Doha, 2014: 1532–1543.

- [8] GREFF K, SRIVASTAVA R K, KOUTNIK K J, et al. LSTM: a search space odyssey[J]. IEEE Transactions on Neural Networks and Learning Systems, 2015, 28(10): 2222–2232.
- [9] SRIVASTAVA N, HINTON G, KRIZHEVSKY A, et al. Dropout: a simple way to prevent neural networks from overfitting [J]. Journal of Machine Learning Research, 2014, 15(1): 1929–1958.
- [10] PHAISANGITTISAGUL E. An analysis of the regularization between L2 and dropout in single hidden layer neural network[C]//International Conference on Intelligent Systems, Modelling and Simulation. Rehovot, 2017: 174–179.
- [11] HE W, LIU K, LYU Y J, et al. DuReader: a Chinese machine reading comprehension dataset from real-world applications[J/OL]. (2017–11–15) [2018–07–20]. <https://arxiv.org/abs/1711.05073>.
- [12] LIU C W, LOWE R, SERBAN I V, et al. How not to evaluate your dialogue system: an empirical study of unsupervised evaluation metrics for dialogue response generation[C]//Conference on Empirical Methods in Natural Language Processing. Austin, 2016: 2122–2132.
- [13] KINGA D, ADAM J B. A method for stochastic optimization[C]//International Conference on Learning Representations. San Diego, 2015.
- [14] WANG S H, JIANG J. Machine comprehension using match-LSTM and answer pointer[C]//The 5th International Conference on Learning Representations. Toulon, 2017.

Machine reading comprehension based on multi-passage ranking

WAN Jing GUO YaZhi

(College of Information Science and Technology, Beijing University of Chemical Technology, Beijing 100029, China)

Abstract: To solve the problem of machine reading comprehension of multi-paragraphs, we propose a model named PR-BiDAF which uses BiLSTM and self-attention based on the bi-directional attention flow (BiDAF) model. The model is used to locate the paragraph in which the answer is located, and then to find the beginning and end of the final answer in the prediction paragraph. Experiments show that, compared with the BiDAF model, the paragraph selection accuracy, BLEU4 index and ROUGE-L index of the PR-BiDAF model proposed in this paper are increased by about 13%, 6% and 4% respectively.

Key words: machine reading comprehension; bi-directional attention flow (BiDAF) model; self-attention

(责任编辑:汪 琴)