

基于语义联合的知识库问答方法

苟苗苗, 徐涛, 吴瑜, 戴玉刚

西北民族大学中国民族语言文字信息技术教育部重点实验室, 甘肃 兰州

收稿日期: 2023年9月28日; 录用日期: 2023年11月2日; 发布日期: 2023年11月13日

摘要

基于知识图谱的问答(KBQA)是问答系统的重要组成部分。然而, 现有大多数知识图谱问答系统侧重于回答单个三元组查询的简单问题, 对回答涉及多个实体和关系的复杂问题的正确率较低。为了提高正确率, 本文采用对比学习方法来计算语义相似度, 并提出了语义联合模型框架, 将实体消歧和关系匹配任务进行联合建模, 以避免误差传递问题。最终使用本文的方法在CCKS2019KBQA数据集上进行实验, 实验结果表明, 与BERT模型相比, 对比学习模型在计算语义相似度方面更具优势, 并且语义联合建模的效果也优于先进行实体消歧再进行关系匹配的方法。

关键词

知识库问答, 对比学习, 语义联合, 相似度

Knowledge Base Question Answering Method Based on Semantic Fusion

Miaomiao Gou, Tao Xu, Yu Wu, Yugang Dai

Key Laboratory of Information Technology and Education Ministry of Chinese Ethnic Languages and Writings, Northwest Minzu University, Lanzhou Gansu

Received: Sep. 28th, 2023; accepted: Nov. 2nd, 2023; published: Nov. 13th, 2023

Abstract

Translation: Knowledge-based question answering (KBQA) is an important component of question answering systems. However, most existing KBQA systems primarily focus on answering simple questions that involve single triple pattern queries, resulting in lower accuracy when it comes to answering complex questions involving multiple entities and relationships. To improve the accuracy, this paper adopts contrastive learning to calculate semantic similarity and proposes a semantic joint modeling framework that combines entity disambiguation and relationship matching

tasks to avoid error propagation. The proposed method is evaluated on the CCKS2019KBQA dataset, and the experimental results demonstrate that compared to the BERT model, the contrastive learning model has a greater advantage in computing semantic similarity, and the effectiveness of semantic joint modeling is superior to the approach of first performing entity disambiguation and then relationship matching.

Keywords

QA Based on Knowledge Base, Contrastive Learning, Semantic Fusion, Similarity

Copyright © 2023 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

知识库是存储知识的数据库，存储包含规则联系的事实的数据，知识库有多种表现的方式，知识图谱就是使用最广泛的一种知识库。问答是自然语言处理中的重要任务，它的目的是回答由人提出的自然语言问题，知识库问答[1]就是指通过检索储存在知识库中的三元组来获取问题的答案实体，来返回问题的真实答案的一种问答方法。而在基于知识库的问答中，如何精准检索到和问题匹配的正确答案是问答任务的关键。知识库问答按照检索信息的难易程度进行分类，可以分为简单问题和复杂问题。

简单问题是结构简单的问题，比如“叔本华信仰什么宗教？”这种只需要在知识库中只进行一次检索就可以的到答案的问题。复杂问题是指多跳和多约束问题，多跳问题在知识库中搜索时通过多步检索得到答案。例如，“请回答《西风颂》作者的主要成就”，这个问题首先要找出问题中的实体“西风颂”，然后再知识库中搜索与问题相关的包含《西风颂》的三元组(珀西·比希·雪莱, 代表作品, 西风颂)，确定与“《西风颂》作者”对应的实体“珀西·比希·雪莱”，再在知识库中寻找与“珀西·比希·雪莱”相关的三元组(珀西·比希·雪莱, 主要成就, 创作大量浪漫主义诗歌)从而获取到《西风颂》作者的主要成就为“创作大量浪漫主义诗歌”。复杂问题的示例图如图1所示。

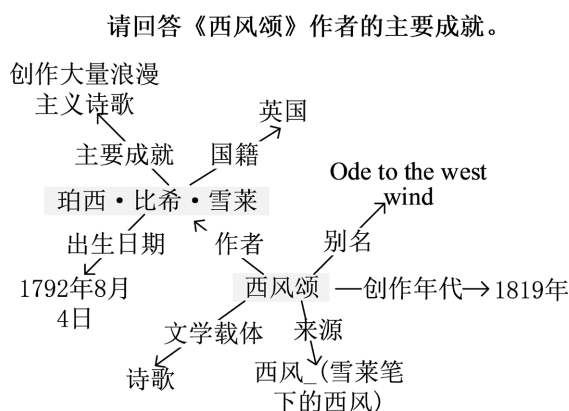


Figure 1. Example of complex problem query

图1. 复杂问题查询实例

现在知识库问答的解决方法主要的有两种，Lan 等人[2]的综述中总结了复杂知识库问答方法，重点

归纳了基于语义解析和基于信息检索的方法。前者将问句解析成为由实体和关系构成的语法树，然后用语法树生成语句查询知识库，然而这种严重依赖于中间逻辑形式(如 SPARQL)的昂贵注释。后一种方法不解析问题，而是直接表示实体，并根据它们与输入问题的相关性对其进行排序。其中，首先检索与问题相关的路径，然后对其进行缩小。路径检索对于问答系统的性能至关重要，因为短的路径极有可能排除答案，而过长的路径则可能引入影响问答系统性能的噪声[3]。

基于信息检索的知识库问答方法首先要通过命名实体识别得到问句中的实体与候选实体，然后在知识库中进行若干次检索获得与包含答案的候选路径，再计算自然语言问题与候选路径之间的语义相似度、Jaccard 距离等元素，然后使用逻辑回归模型对候选路径进行排序，找出候选路径中最优的路径，输出答案。因为基于信息检索的方法更方便寻找答案，在逻辑理解与构造数据集方面也较为简单，所以获得了更多的探寻。然而基于信息检索的方法在根据实体和候选实体检索知识库的这个过程中，路径数目会随跳数增加而呈指数级增长，引入更多噪声，使问答系统变得繁冗，也导致了正确率的降低。

为了避免引入更多噪声，提高正确率，本文提出了一个基于对比学习的语义联合模型框架，将实体消歧和关系匹配整合到一个统一的框架中，利用实体所连接的关系信息，同时完成实体消歧和关系匹配任务，防止误差传播，最大可能减少错误路径和无关路径，解决复杂问题，提高系统整体性能。

2. 基于语义联合的知识库问答方法

开放域知识图谱问答一般有三个模块，实体识别、实体消歧和关系匹配。在实体识别模块中，计算机识别出问题中的实体提及。但是自然语言疑问句中的实体提及通常可能含有多个意思，实体存在的缩写、别名、嵌套以及问句与知识库中结构化语义之间的差距，例如“苹果”可能是水果，也可能是手机，在这种情况下就需要通过实体消歧来找到知识库中准确对应的实体。用户问句的意图通常具有不同的表现形式，关系匹配任务用于匹配问句意图和知识图谱中相近的关系。但是大多数研究都将实体消歧和关系匹配视为独立子任务导致误差传递，使得整体系统准确率不佳。

本文构建基于对比学习的实体消歧关系匹配任务联合框架如图 2 所示，将实体消歧与关系匹配统一建模联合，减少传递误差，提高相似度计算准确率。在得到问句中识别的实体提及，在 Neo4j 图数据库中进行模糊匹配，得到候选集合，用 mention2id 词典对其进行过滤，然后使用语义联合模型控制路径扩展，最大可能减少错误路径和无关路径，降低噪声，提高系统整体性能。

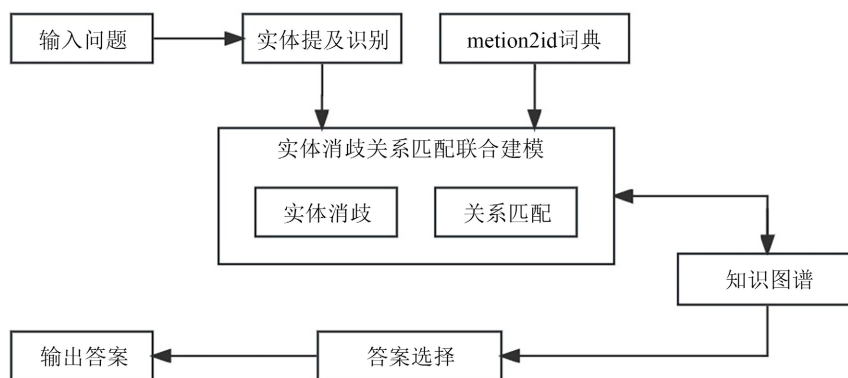


Figure 2. Semantic-based knowledge base question answering approach

图 2. 基于语义联合的知识库问答方法

2.1. 实体命名识别

命名实体识别[4]是为了找到自然语言问题 q 中实体，常用的实体识别方法是基于模型的方法。常用

的实体识别模型为序列标注模型，可以给自然语言中的每个词映射最优的标签序列。

BERT 是一种双向的 Transformers 模型，通过掩码机制进行预训练，取代了传统的单向语言模型或对单向语言模型进行浅层拼接的方法。在文本建模方面，BERT 利用基于注意力机制的方法，使输入序列中的每个词都能够关注到其他词，从而更好地实现对语义层面的理解。通过双向编码的机制，BERT 能够更好地将词的上下文信息融合到词向量中，使得词向量的表示更加准确。

本文尝试利用 BERT-BiLSTM-CRF 模型[5]对问题中实体指称的起点位置及终点位置进行预测，寻找问句中的实体。BERT-BiLSTM-CRF 模型结构如图 3 所示，主要由 BERT 层、BiLSTM 双向长短期记忆网络，CRF 条件随机场三层组成。

在 BERT-BiLSTM-CRF 模型中输入问句后，通过 BERT 模型获得问句相应的字向量表示，然后将字向量输入到双向长短期记忆网络进行对字向量的双向编码；最后用 CRF 条件随机场进行特征解码，将自然语言问句映射得到最优的标签序列，从而完成对整个问题命名实体识别的流程。

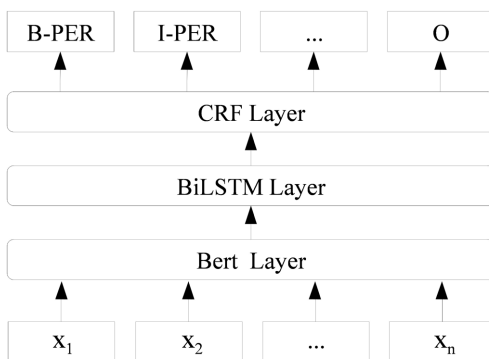


Figure 3. BERT-BiLSTM-CRF model architecture
图 3. BERT-BiLSTM-CRF 模型结构

2.2. 对比学习模型

2.2.1. Sentence-BERT 模型

Sentence-BERT [6]使用基于 BERT 模型[7]孪生网络结构来获得两个句子的词向量表示，然后对其进行预训练以建立相似度模型。使用 BERT 模型进行相似度计算会消耗很长的时间进行语义相似性的搜索，而且句子表示不适合无监督的任务，Sentence-BERT 减小 Bert 语义搜索的巨大耗时，使其适用于句子相似度计算。

在 Sentence-BERT 模型的预训练阶段，首先使用 BERT 的孪生网络获取句子的向量表示。这意味着将两个句子输入到两个共享参数的 BERT 模型中，得到这两个句子的词向量表示。然后，在句子长度的维度上对所有词向量求平均值。这一步骤是将 BERT 输出的词向量输入到池化层进行平均池化操作，以获得两个句子的句向量表示。最后，使用余弦相似度公式计算词向量 U 和 V 之间的相似度，从而计算出两个句子之间的相似度。Sentence-BERT 模型的结构如图 4 所示。

我们使用 Sentence A 和 Sentence B 的特征向量作为孪生网络模型的输入，输出两个句子的词向量进入 Pooling 层进行平均池化，得到句子的向量表示 \$x_i, y_i\$，然后采用余弦相似度来计算问题 Sentence A 和 Sentence B 的相关度得分如式(1)。

$$\cos(\theta) = \frac{\sum_{i=1}^n (x_i + y_i)}{\sqrt{\sum_{i=1}^n (x_i)^2} + \sqrt{\sum_{i=1}^n (y_i)^2}} = \frac{a \cdot b}{\|a\| \times \|b\|} \quad (1)$$

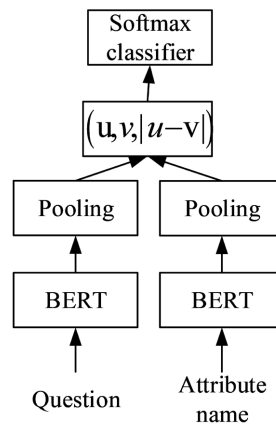


Figure 4. Sentence-BERT model architecture
图 4. Sentence-BERT 模型结构

评价本文问题和实体属性相关度的准确标准使用的指标是估计距离 $\hat{\theta}$ 与真是距离 θ 的距离的函数，最常用的函数是距离的平方，但是因为估计距离 $\hat{\theta}$ 具有随机性，所以求该函数的期望值，公式如(2)给出了均方误差公式。在训练模型的过程中，我们使用最小化均方误差作为损失函数来训练相似度模型。

$$\text{MES}(\hat{\theta}) = E \left\{ \left[\hat{\theta} - E(\hat{\theta}) \right] + \left[E(\hat{\theta}) - \theta \right] \right\}^2 \quad (2)$$

2.2.2. ConSERT 模型

ConSERT [8] (Contrastive Framework for Self-Supervised Sentence Representation Transfer)是一种基于对比学习的句子表示迁移方法，旨在通过使用共享参数的 BERT 模型对两个句子进行向量表示来获得更具区分性的语义向量表示。该模型通过引入对比学习的思想，在训练过程中通过比较同一句子的不同增强样本之间的相似性以及不同句子之间的差异性，进一步提升语义表示的质量。与 Sentence-BERT 相似，ConSERT 采用了共享参数的 BERT 模型架构，但利用对比学习的思想对模型进行了进一步训练，以获得更加区分度的语义向量表示。这种对比学习的方法能够生成更具区分性的句子表示，从而在各种自然语言处理任务中取得更好的性能。通过对比学习，ConSERT 能够在无监督的自我训练过程中学习到更具语义区分度的句子表示，从而提高了句子表示的迁移能力和表达能力。这使得 ConSERT 成为一个有效的句子表示迁移方法，可在多种自然语言处理任务中发挥作用，如文本分类、语义匹配等。

ConSERT 模型的结构如图 5 所示，由三个主要部分组成：共享参数的 BERT 模型层、对比损失层和数据增强模块。在每个批次中，有 N 个输入文本，数据增强模块采用不同的方法生成两个增强样本。这两个增强样本经过共享参数的 BERT 模型层进行向量表示，然后通过平均池化映射到相同的维度上，并在对比损失层计算对比损失。损失函数是通过计算两个句子之间的余弦相似度来度量它们的相似性。对比损失层的目标是最大化正样本的相似度，并将负样本的相似度最小化，以增强语义表示的区分度。通过这种对比学习的训练方式，ConSERT 能够生成更具区分性的语义向量表示，从而在各种自然语言处理任务中获得更好的性能。损失函数如式(3)所示：

$$L_{i,j} = -\log \frac{\exp(\text{sim}(r_i, r_j)/\tau)}{\sum_{K=1}^{2N} \mathbf{1}_{K \neq i} \exp(\text{sim}(r_i, r_j)/\tau)} \quad (3)$$

其中 $\text{sim}(\cdot)$ 函数为余弦相似度函数， $\mathbf{1}$ 为指示函数， r 为句子对的向量表示， τ 为温度超参，设置为 0.1，最终的对比损失是通过平均批次内的所有损失计算得到的。对比损失函数的目标是拉近相似句子，拉远

不相似句子。通过找到每个样本的增强样本，并计算它们之间的相似度，以及将其他样本作为负样本，可以提高语义表示的区分度。最终的对比损失是批次内所有样本损失的平均值。

上述为在无监督情况下的对比损失，ConSERT 在有监督的情况下，采取损失融合的方式进行，将有监督数据视为分类任务数据，并结合对比损失进行训练，式(4)和式(5)为分类损失函数：

$$L_{ce} = CrossEntropy(W(r_1, r_2, |r_1 - r_2|) + b, y) \tag{4}$$

$$L_{joint} = L_{ce} + \alpha L_{ci} \tag{5}$$

其分类损失使用交叉熵损失函数， r_1, r_2 为句子向量表示，并且和 $|r_1 - r_2|$ 拼接起来， α 作为对比损失的平衡参数。

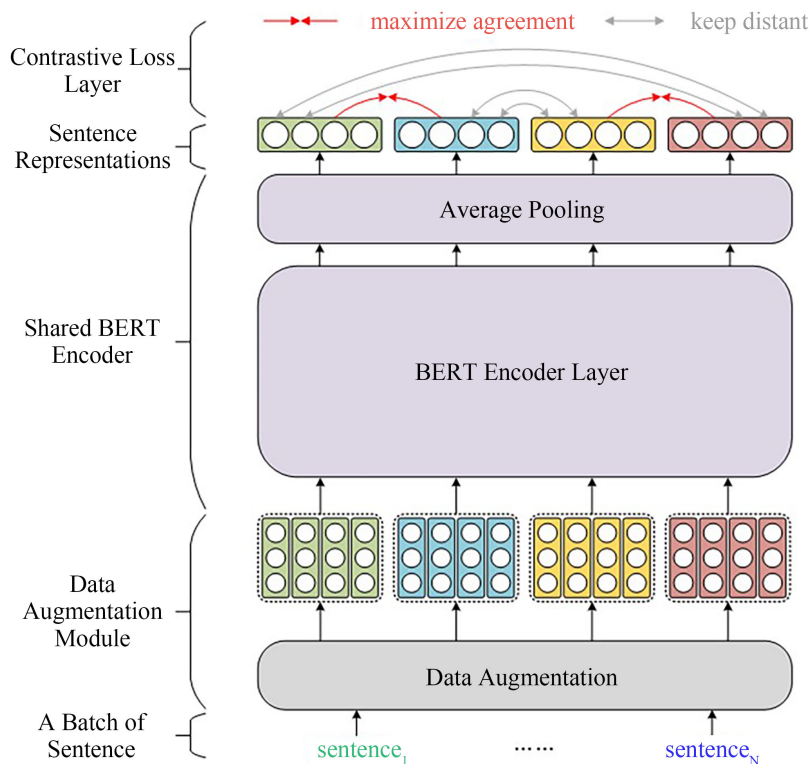


Figure 5. ConSERT model architecture
图 5. ConSERT 模型结构

2.2.3. CoSENT 模型

CoSENT [9] (Cosine Sentence)模型使用两个参数共享的 BERT 构成孪生网络，和 Sentence-BERT 模型结构相同，对于输入句子 U 和 V ，输出它们各自的语义向量，然后进行池化操作，最后使用余弦相似度函数进行相似度计算。CoSENT 使用对比学习的方式是针对一对句子而不是一个句子对比。具体来说，CoSENT 模型在训练阶段，记 h^+ 为所有的正样本对集合， h^- 为所有的负样本对集合，对于任意的正样本对 $(h_i, h_j) \in h^+$ 和负样本对 $(h_k, h_l) \in h^-$ ，都有：

$$\cos(u_i, u_j) > \cos(u_k, u_l) \tag{6}$$

其中 u_i, u_j, u_k, u_l 分别为 h_i, h_j, h_k, h_l 的句向量表示。CoSENT 模型使用对比损失优化句子对的余弦相似度替换 Sentence-BERT 的训练任务，获得了句子更有区分度的语义向量，其损失函数如公式(7)所示：

$$\log\left(1 + \sum (h_i, h_j) \in h^+, (h_k, h_l) \in h^- - e^{\lambda(\cos\cos(u_i, u_j) - \cos\cos(u_k, u_l))}\right) \quad (7)$$

其中 λ 是一个大于 0 的超参数，后续实验取为 20。该损失函数可以再训练过程中，将语义相似的句子对在向量空间中的表示拉近，不相似的句子对远离，来得到更有区分度的句子向量表示。

2.3. 基于对比学习的语义联合模型

在先执行实体消歧任务再执行关系匹配任务的过程中，会导致误差的传递，如果实体消歧模型选出的实体就已经偏离了问句，那么关系匹配模型将无法正确找到关系，进而无法在知识图谱中找到正确的答案。同时，在这种情况下，实体消歧过程中无法利用关系匹配阶段的信息，如有些候选实体根本没有正确的关系，而实体消歧任务中仍然可能会选择这个候选实体，最终导致结果错误。因此本文用出了语义联合建模框架来完成实体消歧关系匹配联合任务，其直接计算候选集合路径和问句之间的语义相似度。

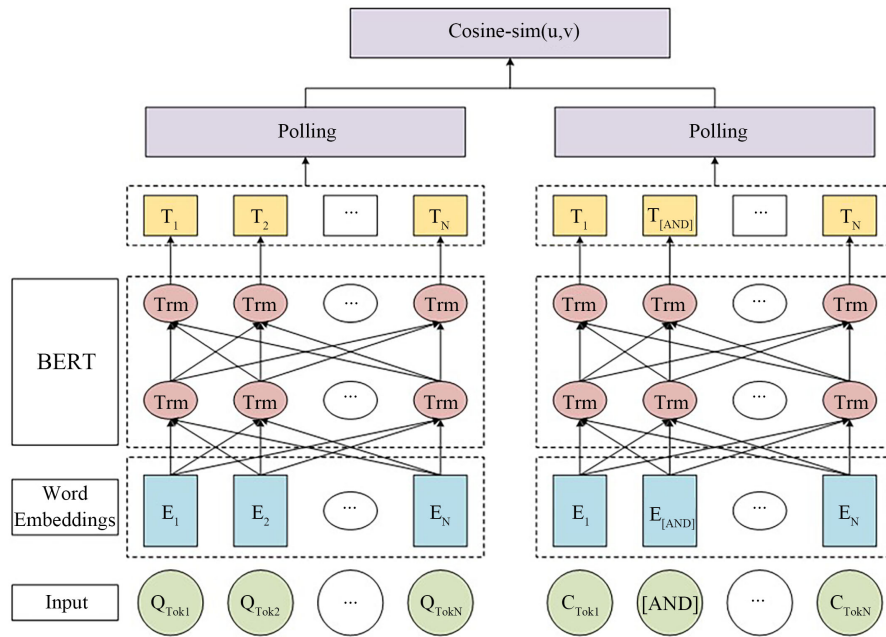


Figure 6. Framework diagram of a contrastive learning-based semantic fusion model
图 6. 基于对比学习的语义联合模型的框架图

图 6 为基于对比学习的语义联合模型的框架图，其使用对比学习的方法来学习问题以及候选集合的语义向量表示，使其更容易区分开来，并将实体消歧和关系匹配统一建模联合完成，以避免误差传递。首先将候选实体集 $E = \{e_1, e_2, \dots, e_n\}$ 中的每一个候选实体 e_i 和其所连关系集 $R_i = \{r_1^i, r_2^i, \dots, r_n^i\}$ 通过一个特殊的 [AND] 标识符进行连接，构成候选集合集 $C_i = \{e_1 r_1^1, \dots, e_1 r_n^1, e_2 r_1^2, \dots, e_i r_n^i\}$ 。其次将问句 Q 和候选集合集 C 输入进共享参数的 BERT 层，得到它们的向量表示。然后将这些向量分别输入进池化层以获得固定大小的句子嵌入，其表示为：

$$H^q = \text{polling}(\text{Bert}(Q)) = \{h_1^q, h_2^q, \dots, h_n^q\} \quad (8)$$

$$H^c = \text{polling}(\text{Bert}(C)) = \{h_1^c, h_2^c, \dots, h_n^c\} \quad (9)$$

其中池化层默认使用平均池化策略。最后使用余弦相似度函数计算他们的相似度：

$$sim_s = \cos(H^q, H^c) \tag{10}$$

其中 sim_s 为问句与候选集合的相似度得分集合。直观地说，一些候选关系可以提供一些语义信息给实体消歧。如果知道问句中的关系，就可以通过其提供的语义信息排除一些候选实体。例如，问句“水浒传有多少集？”包含词语“多少集”对应的关系“集数”。为了实体消歧，有理由将注意力集中在连接有“集数”的候选实体上，例如“水浒传(电视剧)”而不是“水浒传(小说)”。因此，本章构建了一个语义联合模型来同时进行实体消歧和关系匹配。

2.4. 知识库问答的实现

在知识库查以基于信息检索的方式查询一个问题，返回包含问题答案的候选路径。在多跳问答中，若是跳数多于或者少于可以找到问题答案的真实跳数，则会引入影响问答系统性能的噪声，并且考虑到大多数复杂问题可以在两条内解决，为了减少候选路径的规模，去除大量与问题无关的路径，因此本文只选择两跳的扩展路径，并在这个过程中加入语义联合模型，通过计算问题与路径的相似度对路径每一跳质量评估，保留高质量路径。

对于问句中的主题实体集合 E_{topic} 中的实体 e ，如问句“白鹿原的作者出生地在哪儿”，检索问句的实体“白鹿原”，得到实体属性为“作者”、“出版社”、“总策划”和“导演”等属性，然后计算问题与每个候选属性的语义相似度，自然排除了与问题无关的实体属性与包含这些属性的路径，得到(白鹿原, 作者, 陈忠实)，然后根据第一跳得到的属性“陈忠实”检索，得到与“陈忠实”有关的属性“国籍”、“民族”、“出生地”和“出生日期”等，再次计算实体属性与问题的相关度，得到(陈忠实, 出生地, 陕西)，最终得到问题的答案“白鹿原的作者出生地在陕西”。

因为在根据“白鹿原”扩展路径时，第二跳可能会出现(俯仰关中, 作者, 陈忠实)这种情况，也有可能第一跳中出现实体关系与正确无关的情况，所以在路径扩展中方式如图7。

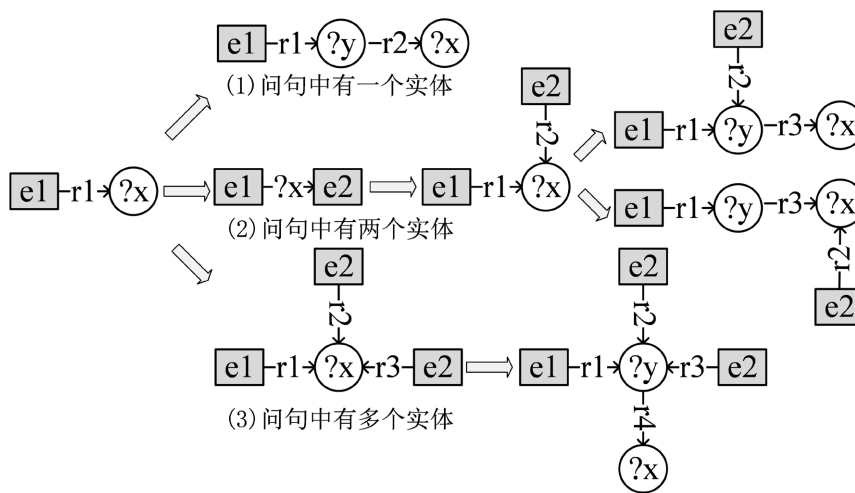


Figure 7. Path expansion methods
图7. 路径扩展方法

(1) 对于单实体问题的答案路径：在第一跳中，首先通过问句中识别的实体提及，在 Neo4j 图数据库中进行模糊匹配，得到候选集合，之后使用 mention2id 词典对其进行过滤，只保留词典实体提及对应的候选实体及其关系，使用联合语义建模框架计算问句 Q 和实体 - 关系对的语义相似度，选取排名前 n 的候选路径集合 C ；得到第一跳的候选集合中实体提及对应属性以后，根据属性在 Neo4j 图数据库中进行

精确匹配，计算问句 Q 和实体 - 关系 - 关系对的语义相似度，选取排名前 n 的实体 - 关系 - 关系对加入候选路径集合 C 。

(2) 对于两个实体问题的答案路径：首先查询多实体之间相似的关系，构成候选路径集合 C ；若无相似关系，每个实体提及进行第一跳，如果有相同属性则合并，加入候选路径集合 C ；若第一跳有相同属性，则将属性作为实体在 Neo4j 图数据库精确匹配，计算语义相似度，排名前 n 则加入候选理解集合 C ；若第一跳无相同属性，则联合语义建模框架计算问句 Q 与两对实体 - 关系对的语义相似度，选取排名前 n 关系对一一进行第二跳，合并两个实体 - 关系 - 关系对和两个实体 - 关系，不能合并则舍弃。

(3) 对于多实体问题的答案路径：首先查询多实体之间相似的关系，构成候选路径集合 C ；然后第一跳，对多实体相似的关系进行扩展，合并想同的属性，计算多实体 - 关系对于问题 Q 的相似度，排名前 n 加入候选路径集合 C 。

在得到所有的候选路径集合 C 后，使用语义联合建模框架计算问题 Q 与候选路径集合 C 的语义相似度，选取排名最高的路径最为最终路径，在 Neo4j 图数据库中匹配最终答案。

3. 实验

3.1. 数据集分析

本文使用了 PKUBASE 大规模开放域知识库，在 CCKS2019KBQA 数据集上进行测试，其中，训练集有 2298 条，验证集有 766 条，测试集有 766 条，简单问题和复杂问题的比例为 1:1。一条完整的问答数据包括问题、SPARQL 查询语句和问题的答案，如表 1 为一条问答数据的示例。

本文使用的 PKUBASE 知识库包含了三个文件，分别是知识库的三元组，有 41,009,141 条数据；实体类型，有 25,182,627 条数据；实体别称，有 13,930,117 条数据。知识库使用 Neo4j 图数据进行储存，并提供查询，mention2id 词典为数据集所提供，用于辅助开放域问题中主题实体提及映射到知识图谱中的实体上，然后在知识图谱的三元组中根据问题寻找若干个实体或者属性名作为答案。

Table 1. Example of a question-answer data

表 1. 一条问答数据的示例

问句	《论衡》的作者是哪个民族的人民？
SPARQL 查询语句	<code>select ?y where { ?x <代表作品><论衡>. ?x <民族> ?y. }</code>
答案	<汉族>

3.2. 实验设置

本实验使用 BERT、Sentence-BERT、Con SERT、Co SENT 四个模型，实验相关软件及训练环境为 Pycharm Community Edition 2020.2.5 x64、Python3.7 版本、Tensorflow-gpu 1.14、Pytorch 1.7。

实验使用的所有 BERT 预训练模型为 bert-base-chinese 的权重，训练了 30 个 epoch，最大序列长度为 64，批次大小为 16，学习率为 $2e-5$ 。

3.3. 实体识别模型效果

本文针对 BERT-CRF、BiLSTM-CRF、BERT-BiLSTM-CRF 三个实体识别模型训练，实验结果如表 2 所示，可以看出用 BERT-CRF 模型在识别问题中的实体提及任务上的效果要优于 BiLSTM-CRF 模型，BERT-BiLSTM-CRF 实现实体命名识别的效果高于另外两个模型。

Table 2. Entity named recognition performance
表 2. 实体命名识别效果

模型	P	R	F1
BiLSTM-CRF	0.7431	0.8108	0.7806
BERT-CRF	0.8367	0.8512	0.8358
BERT-BiLSTM-CRF	0.8437	0.8612	0.8523

3.4. 语义联合模型任务效果

实体消歧关系匹配联合任务中,对各个模型的准确率进行了实验,使用准确率(Accuracy)作为评价指标。联合匹配任务模型效果如表 3 所示,实体消歧关系匹配任务的联合模型中候选实体 - 关系数量较多,实验结果表明和 BERT 模型相比,对比学习模型能够学习到更深层次的语义信息,在语义联合模型框架中 CoSENT 模型表现效果最好,但是 conSERT 模型也有不错的速率。

Table 3. Accuracy of the entity disambiguation and relation matching joint task model
表 3. 实体消歧关系匹配联合任务模型准确率

模型	P
BERT	0.8331
Sentence-BERT	0.8512
conSERT	0.8725
CoSENT	0.8765

3.5. 问答系统性能评估

本文最终在 CCKS2019KBQA 数据集上进行测评,评测最终结果采用的评价指标为 AverageF1 值。本文在识别问题中的实体之后,做以下两个实验:(1) 先做实体消歧然后进行关系匹配的方法作为 Pipeline 方法;(2) 将两个子任务联合进行的方式即使用语义联合匹配框架。最终问答结果如表 4 所示。实验结果表明了将实体消歧关系匹配任务进行联合匹配框架均优于 Pipeline 方法,证明了该框架在开放域知识图谱问答中的优势性。

Table 4. Evaluation results of the CCKS2019KBQA dataset
表 4. CCKS2019KBQA 数据集测评任务结果

模型	F1 值
Sentence-BERT (Pipeline)	0.7158
conSERT (Pipeline)	0.7311
CoSENT (Pipeline)	0.7332
Sentence-BERT	0.7349
conSERT	0.7375
CoSENT	0.7397

4. 结论

本文介绍了一种基于语义联合的知识库问答方法,实现了对简单问题和复杂问题的统一处理。系统

使用实体消歧关系匹配联合任务模型框架, 充分利用实体和关系之间的交互信息, 获取更加具有区分度的语义向量表示, 提升了问题与路径之间的相关能力, 在任务数据集上取得了比先做实体消歧然后进行关系匹配的方法更好的效果, 证明了基于对比学习的联合语义建模框架在开放域知识图谱问答中的有效性。

虽然本文的知识库问答方法提高了系统的正确率, 但仍然有欠缺之处。为了解决复杂问题, 在识别到问句中的实体提及后, 在 Neo4j 图数据库中进行模糊匹配和二次匹配, 减少了噪声, 但是最终获得的路径还有很多, 整个系统依然十分繁冗。因此后续我们会将解决复杂问题的研究重点放在相关子图的推理中, 减少答案的候选集合, 结合本文的方法计算问题与候选路径集合的相似度, 达到简化系统的目的。

基金项目

1) 国家档案局科技计划项目《基于大数据智能驱动的档案信息资源挖掘与共享利用服务研究》的阶段性成果, 项目批准编号: 2021-X56; 2) 甘肃省委办公厅(甘肃省档案局), 甘肃省档案科技项目《甘肃省档案信息资源大数据分析及其数据可视化研究与应用》研究成果之一, 甘档发[2020]48号 GS-2020-X-07, 2020-09至2021-09, 在研, 主持。

参考文献

- [1] Chen, W., Zha, H., Chen, Z., *et al.* (2020) HybridQA: A Dataset of Multi-Hop Question Answering over Tabular and Textual Data. *Findings of the Association for Computational Linguistics: EMNLP 2020*, Online, 1026-1036. <https://doi.org/10.18653/v1/2020.findings-emnlp.91>
- [2] Lan, Y., He, G., Jiang, J., *et al.* (2021) A Survey on Complex Knowledge Base Question Answering: Methods, Challenges and Solutions. *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence Survey Track*, 4483-4491. <https://doi.org/10.24963/ijcai.2021/611>
- [3] Bao, J.W., Duan, N.Y., *et al.* (2016) Constraint-Based Question Answering with Knowledge Graph. *Proceedings of the 26th International Conference on Computational Linguistics Technical Papers*, 2503-2514.
- [4] 徐增林, 盛泳潘, 贺丽荣, 王雅芳. 知识图谱技术综述[J]. 电子科技大学学报, 2016, 45(4): 589606.
- [5] 谢腾, 杨俊安, 刘辉. 基于 BERT-BiLSTM-CRF 模型的中文实体识别[J]. 计算机系统应用, 2020, 29(7): 48-55. <http://www.c-s-a.org.cn/1003-3254/7525.html>
- [6] Reimers, N. and Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings Using Siamese BERT-Networks. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, 3982-3992. <https://doi.org/10.18653/v1/D19-1410>
- [7] Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019) BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Vol. 1, Minneapolis, Association for Computational Linguistics, 4171-4186.
- [8] Yan, Y., Li, R., Wang, S., *et al.* (2021) Consert: A Contrastive Framework for Self-Supervised Sentence Representation Transfer. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, 5065-5075. <https://doi.org/10.18653/v1/2021.acl-long.393>
- [9] 苏剑林. CoSENT(一): 比 Sentence-BERT 更有效的句向量方案[EB/OL]. <https://spaces.ac.cn/archives/8847>, 2022-01-06.