

# 大连理工大学

## 专业学位硕士研究生中期考核报告

论文题目：基于图卷积网络的文档级别关系抽取研究

姓名：徐奕斌

学号：31909185

专业/领域：计算机技术

培养类型：☒全日制 ☐在职学习 ☐非全日制

指导教师：姚卫红

实践导师：周惠巍

入学日期：2019.09

报告日期：2021.1.23

报告地点：线上

研究生院制表

## 说 明

中期检查是保证学位论文质量、工作进度和研究生培养质量的重要措施。

一、考核内容：首先，考查学位论文内容是否符合专业学位人才培养要求，开题时确定的校内外实习实践环节完成情况；其次，考查学位论文内容完成情况、阶段性成果是否正确，开题时方案是否需调整或已做了哪些调整，后续工作思路是否正确、工作进度是否有保障、预期目标能否实现、论文质量是否能够保证以及论文工作存在的问题等；第三，考查研究生对本专业/领域的专业知识、技术标准规范等掌握程度，考查学生的工作态度等。

二、报告撰写：中期报告正文字数不少于 4000 字，正文及参考文献等撰写要求参见《大连理工大学博士学位论文格式规范》。

三、报告保存：研究生根据专家的意见修改完善中期报告，经导师签字后生成 PDF 文档，上传到研究生信息管理系统，导师审核，由学部（学院）进行网上确认。

## 中期考核报告正文

撰写大纲：

- 1) 校内外实习实践工作总结；
- 2) 开题时拟定的研究方案、进度计划；若开题时的研究方案已经调整，应说明调整的原因、调整后该领域的国内外研究状况分析、研究内容、研究方法、进度计划等；
- 3) 学位论文的研究进展完成情况、阶段性成果和创新点论述；
- 4) 后续工作的设想、可能遇到的困难和问题及条件保障措施；
- 5) 已发表、录用的论文和已投稿的论文情况。
- 6) 参考文献（不占字数）。

# 基于图卷积网络的文档级别关系抽取研究

## 1 校内外实习实践工作总结

校内实习是在实践导师的指导下，根据实践导师指定的任务展开的。我的实习实践的主要研究任务是生物医学领域的问题蕴含识别。

生物医学问题蕴含识别是指识别给定两个生物医学问题之间是否存在蕴含关系，给定两个生物医学问题 Q1 和 Q2，如果 Q2 的答案是 Q1 的完整答案或者部分答案，那么就可以说 Q1 蕴含了 Q2。

生物医学问题蕴含识别是自然语言处理中的一个重要任务，在问答系统、信息抽取等多个领域中有重要的应用，这一概念是 Asma B 等人<sup>[1]</sup>在 2016 年提出的。经历了早期基于词语级别的分类方法和基于句子编码的深度学习方法<sup>[2]</sup>，到逐词记忆力模型的方法<sup>[3]</sup>，再到匹配的长短期记忆模型和一种新的注意力机制模型<sup>[4]</sup>。由于任务是针对问题对的，那么如何衡量两个问题之间的相似度，是第一个需要解决的问题；第二个需要解决的问题就是存在一些问题对在识别中有一定的困难，这种问题对的特征是问题对之间的主要内容相似，但问题类型却完全不同；再者就是训练语料不充分，且与测试语料差异较大。这三个问题，使得生物医学问题蕴含识别成为一个棘手的任务。

在实践中，我搭建的模型主要由 Embedding 层、交互的 Transformer<sup>[5]</sup>模型以及分类层组成。其中 Embedding 层采用 BioBERT<sup>[6]</sup>，BioBERT 是在 BERT<sup>[7]</sup>的基础上又增加了大量的生物医学相关预料进行预训练得到的，其中会包含大量的生物医学相关的先验知识，以此来得到相应的问题 1、问题 2 各自的上下文表示，其中蕴含丰富生物医学知识信息；然后利用交互的 Transformer 模型对问题 1 和问题 2 的句子表示进行处理，以获取问题 1 的交互的上下文关系表示为例，将问题 1 的句子表示作为 Key、Value 矩阵，将问题 2 的句子表示作为 Query 矩阵，通过这种方式可以在获得问题 1 本身上下文长距离依赖的同时，还能建立问题 1 和问题 2 之间的联系。最后利用线性层进行分类，用交叉熵损失函数进行模型的训练，取得了较为先进的结果。

这次的实践的圆满完成，离不开实践导师的认真指导，也离不开实践导师学生们的热情帮助，让我能快速地熟悉生物医学领域地自然语言处理，准确地掌握相关研究的研究方法，同时不断地开拓自己的知识面，在交流中总结出适合自己的研究方法并且进行深入的研究，取得了一定的成果。

在实践的过程中，我不仅仅学会了使用自然语言处理常用的模型，还掌握了做研究的方法，学会了如何正确的看待问题、解决问题，这对于我以后的研究、工作和生活有着非常重大的影响。

## 2 开题时拟定的研究方案、进度计划

### 2.1 研究方案

实体关系抽取任务旨在从给定的文本中找到多个实体，并判断它们之间是否存在关系、存在何种关系。文本中的实体一般可以分为两种：一种需要自行训练分类器模型进行识别，另一种在文本中已标出。本文的主要研究目标是自动判定实体间的关系，即在给定标准实体的基础上进行实体关系抽取研究。

以往的研究大多以句子作为输入样例，自动抽取句子内实体间的关系，该任务被称为句子级别关系抽取或提及关系抽取。句子级别关系抽取往往需要手动抽取、过滤样例，这会导致部分样例丢失，对实验结果产生影响。以文章作为输入，对文章中所有出现的实体对（包含跨句的）进行关系抽取的任务被称为文档级关系抽取或全局关系抽取。文档级别关系抽取相比于句子级别关系抽取任务难度更大，因为更多的实体数量和更长的上下文序列使得对实体之间的关系建模变得更加困难。以往的句子级别关系抽取模型很难捕获长距离尤其是跨句实体间的复杂语义交互，因而不适用于文档级别关系抽取。如何以整篇文章作为输入，并高效地编码实体、对实体间复杂的交互进行建模，即实现文档级别的关系抽取是本文的研究重点。

本文研究基于图卷积网络的文档级别关系抽取研究，该研究首先通过文本编码器编码文本获得实体的向量表示，然后通过图神经网络对实体间的复杂关系进行建模，最终获得文档级别实体关系抽取模型。

#### （1）基于图卷积网络的文档级别关系抽取模型

首先通过 Glove<sup>[8]</sup>模型编码文本，并通过 LSTM 对上下文信息进行交互，或直接使用 BERT 预训练模型对文档进行编码，通过对 token 和 mention 的平均操作获得实体的向量表示，接下来进行初始无向图的构建，主要是提出合适的规则去进行实体之间的连接，获得合适的邻接矩阵。通过使用图卷积网络更新结点表示，使实体的表示中含有部分与他相邻实体的信息，通过多次图卷积操作，使实体间的本地信息和全局信息充分交互，最终使用头尾实体的向量表示，通过全连接神经网络进行分类。

#### （2）融合了知识表示的用于文档级别关系抽取任务的图卷积网络模型

根据文本信息，抽取知识库中相应的结构化信息，得到知识三元组信息。利用知识表示学习方法，获得低维、稠密的知识表示，一般得到的是实体向量和关系向量表示。通过注意力机制，将知识表示引入，融合知识表示和文本信息，进行图卷积操作并更新实体表示，构建高性能的文档级别关系抽取模型。

图 1 是融合了知识表示的用于文档级别关系抽取任务的图卷积网络模型框

架图。

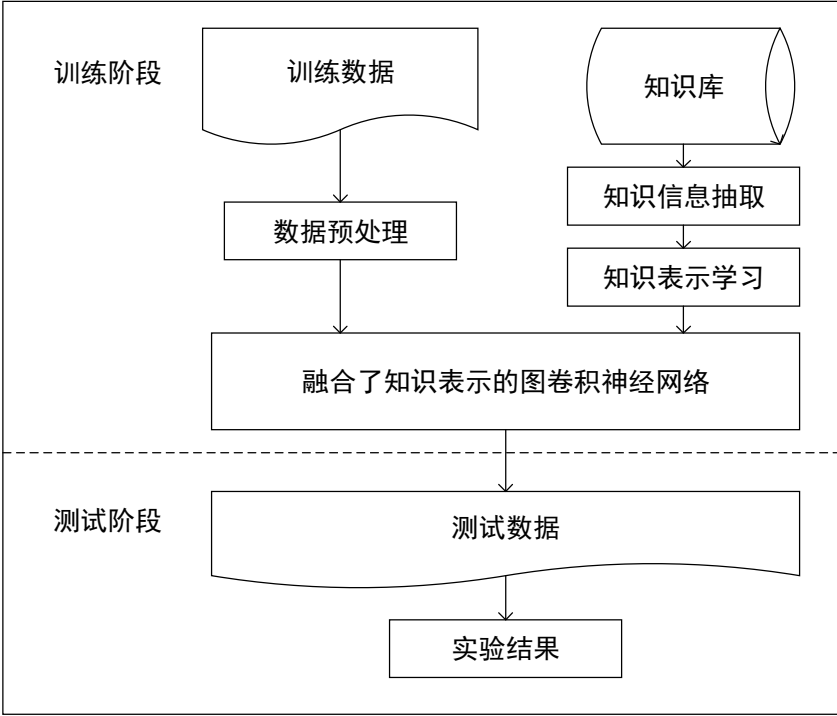


图 1 融合知识表示的用于文档级别关系抽取任务的图卷积网络模型框架图

## 2.2 进度计划

（1）2020 年 7 月到 2020 年 8 月完成基本语料的预处理工作。并收集大量无标注的语料作为深度学习学习词向量的原始数据。

（2）2020 年 9 月到 2020 年 10 月，学习深度学习的理论知识，了解基本深度神经网络模型，如卷积神经网络（CNN）、循环神经网络（RNN）、长短期记忆网络（LSTM）、注意力机制（Attention）。研究如何深入使用神经网络抽取相关特征以及编码文本信息。

（3）2020 年 10 月到 2020 年 11 月，熟悉预训练模型 Bert 的原理和使用方法，尝试用 Bert 模型编码文本并获得实体表示。

（4）2020 年 12 月到 2021 年 1 月寻找相关的知识库，并完成知识库结构化信息的抽取。熟悉知识表示学习工具如 TransE 的使用方法，利用知识表示学习构建实体与关系的知识表示。

（5）2021 年 2 月，学习图卷积网络的理论知识，研究利用图卷积网络对实体之间的复杂关系进行建模。探索如何解决图卷积网络的过渡平滑问题。

（6）2021 年 3 月，学习并深入研究多源信息融合的方法，通过引入注意力机制等各类引入外部信息的方法，有效地融合知识表示信息和文本信息。尝试构

建不同结构的图卷积神经网络模型。

(7) 2021 年 4 月到毕业整理实验数据、图表与看过的相关论文，撰写和修改毕业论文，并参加毕业答辩。

### 3 研究进展完成情况、阶段性成果和创新点论述

#### 3.1 研究进展完成情况

本文的研究分为三部分：数据预处理、基于图卷积网络的文档级别关系抽取模型的构建、融合知识表示的文档级别关系抽取模型的构建。当前研究进展较为顺利，已完成前两部分的工作。

##### (1) 数据预处理阶段

本文采用大规模文档级别关系抽取数据集 DocRED<sup>[9]</sup>训练以及验证提出的模型。DocRED 的训练集有 3053 篇文章，验证集和测试集各有 1000 篇文章。DocRED 数据集标注了 132375 个实体、56354 个关系事实以及 96 种常见的关系类型。数据集中，约 40.7% 的关系事实只能从多个句子中提取出，61.1% 的关系事实需要多种推理方法。除了标注数据集之外，DocRED 也提供了一个包含 101875 篇文章的远程监督数据集。

训练以及测试时，我们对于每一篇输入文章手动构建一张初始图，用作后续的图卷积操作的输入。对于一篇文章，假设其中存在  $N$  个实体，记为  $\mathbf{E} = \{e_v\}_{v=1}^N$ ，每一个实体  $e_v$ ，他们的若干 mention 记为  $\{m_i\}_{i=1}^M$ ，我们按照如下规则构造图  $G(\mathbf{A}, \mathbf{E})$ ，其中  $\mathbf{A}$  是邻接矩阵。

① 我们将实体集合  $\mathbf{E}$  中的每个实体作为图中的一个结点，即  $\mathbf{E}$  也是图像  $G$  的结点集合。

② 如果两个实体在同一个句子中出现，我们就将图中表示这两个实体的结点连接，并且把这个句子添加到边上。由于实体对可以出现在不同的句子中，因此一条边上可以有多个句子。

③ 通过分析训练数据，我们发现在两个邻接的句子中，后一个句子中的代词往往指代的是前一个句子中出现的实体。为了应对这一问题，我们标注训练数据中每个词的词性，选择最常用的代词构建代词此表。如果某个句子中出现词表中的代词，我们就将它拼接到前一个句子上作为一个句子。

##### (2) 基于图卷积网络的文档级别关系抽取模型的构建

基于图卷积网络的文档级别关系抽取模型由四个部分组成：一个编码层、一个上下文感知注意力引导的图卷积模块（CAGGC）、一个多头注意力引导的图卷积模块（MAGGC）和一个分类层。图 2 是模型结构的示意图。

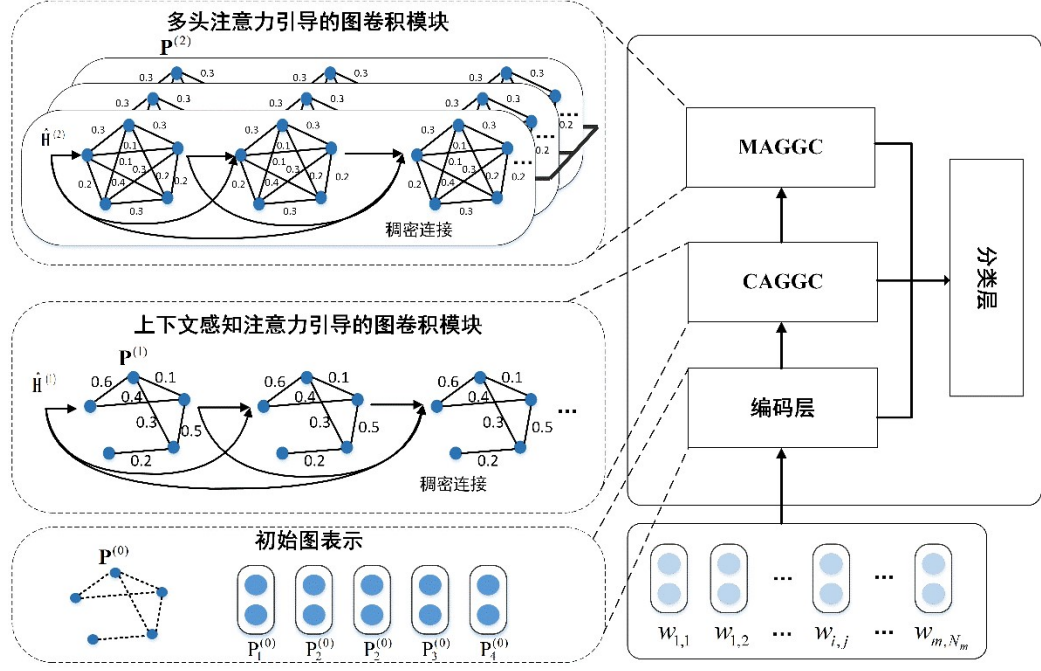


图 2 基于图卷积网络的文档级别关系抽取模型结构图

### ① 编码层

利用编码器将给定的一个文档矩阵  $\mathbf{D} = (w_{1,1}, w_{1,2}, \dots, w_{i,j}, \dots, w_{m,N_m})$  编码到隐藏的向量空间：

$$\mathbf{D}' = \text{Encoder}(w_{1,1}, w_{1,2}, \dots, w_{i,j}, \dots, w_{m,N_m}) = (h_{1,1}, h_{1,2}, \dots, h_{i,j}, \dots, h_{m,N_m}) \quad (1)$$

编码器可以是 BiLSTM 或者 BERT。  $h_{i,j} \in \mathbb{R}^d$  是第  $i$  个句子中第  $j$  个词的表示， $d$  是隐层表示的维度。

接下来要计算构建图的结点表示，也就是文档中的实体表示。假设提及  $m_Q$  是实体  $e_v$  的提及，出现在第  $i$  个句子中，从第  $s$  个词开始到第  $t$  个词结束，我们计算它的表示  $m_Q = \frac{1}{t-s+1} \sum_{j=s}^t h_{i,j}$ ，然后计算出实体表示  $P_v = \frac{1}{J} \sum_{Q=1}^J m_Q$ ，其中  $J$  是实体  $e_v$  的提及数量。

### ② 上下文感知注意力引导的图卷积模块（CAGGC）

我们利用 CAGGC 来构建一个部分连通图，不同于传统图卷积的是，GCGCN 不仅考虑结点的表示，还考虑了构造图中的边的表示。

GCGCN 首先利用实体感知的词级别注意力机制计算句子表示：

$$\alpha_{i,j}^c = \text{softmax}(z^T \tanh(\mathbf{W}_1 h_{i,j} + \mathbf{W}_2 \mathbf{x}_{pos(i,j)}^c + b_1)) \quad (2)$$

$$h_i^c = \sum_{j=1}^m \alpha_{i,j}^c h_{i,j} \quad (3)$$



$$h_i = \mathbf{W}_{wa} [h_i^u; h_i^v] + b_{wa} \quad (4)$$

其中  $c \in \{u, v\}$  表示两个实体中的任意一个,  $\mathbf{x}_{pos(i,j)}^c$  表示当前词和实体的相对位置,  $\mathbf{W}$  和  $b$  是可训练参数。

然后利用实体感知的句子级别注意力和门控机制计算实体间的边的表示, 得到边的表示矩阵  $\hat{\mathbf{H}}^{(1)}$ :

$$\beta_i^c = \sigma(\mathbf{W}_3^T \tanh(\mathbf{W}_4 h_i + \mathbf{W}_5 \mathbf{P}_c^{(0)} + b_2)) \quad (5)$$

$$\hat{h}_{u,v}^{c(1)} = \frac{1}{S} \sum_{i=1}^S \beta_i^c h_i \quad (6)$$

$$\hat{h}_{u,v}^{(1)} = \mathbf{W}_{sg} [\hat{h}_{u,v}^u; \hat{h}_{u,v}^v] + b_{sg} \quad (7)$$

传统图卷积中使用的邻接矩阵由 0 和 1 组成, 表示结点之间是否存在边缘连接, 不能有效控制实体之间的信息传播。因此本文提出了一种综合考虑结点信息和边信息的加权邻接矩阵计算方法, 结点  $u$  和  $v$  之间的权重为:

$$A_{u,v}^{(1)} = \frac{\exp(\tanh(\mathbf{W}^T (\mathbf{W}_u \mathbf{P}_u^{(0)} + \mathbf{W}_v \mathbf{P}_v^{(0)} + \mathbf{W}_e \hat{h}_{u,v}^{(1)})))}{\sum_{u \in \text{neighbour}(v)} \exp(\tanh(\mathbf{W}^T (\mathbf{W}_u \mathbf{P}_u^{(0)} + \mathbf{W}_v \mathbf{P}_v^{(0)} + \mathbf{W}_e \hat{h}_{u,v}^{(1)})))} \quad (8)$$

本文的图卷积操作同样考虑了边的表示:

$$\mathbf{P}_v^k = \text{ReLU} \left( \sum_{u \in \text{neighbour}(v)} A_{u,v}^{(1)} (\mathbf{W}_{node}^k \tilde{\mathbf{P}}_u^{k-1} + \mathbf{W}_{edge}^k \hat{h}_{u,v}^{(1)} + b^k) \right) \quad (9)$$

本文采用稠密连接<sup>[10]</sup>去融合前  $k-1$  个子层的输出, 并记融合结果为  $\tilde{\mathbf{P}}_u^{k-1}$ :

$$\tilde{\mathbf{P}}_u^{k-1} = [\mathbf{P}_u^{(0)}; \mathbf{P}_u^1; \dots; \mathbf{P}_u^{k-1}] \quad (10)$$

为了在不改变输出结点维度的情况下融合这些表示, 本文使用线性层来减少它们的维度, 即将这些子层输出的维度变为原先的  $1/K$ 。最终, 我们利用稠密连接融合初始结点表示和 CAGGC 所有子层的输出表示, 将融合的结点表示作为一个模块的输入。

### ③ 注意力引导的图卷积模块 (MAGGC)

基于传统图卷积的关系抽取模型只能在直接连接或紧密连接的实体之间建立交互。为了解决这个问题, 我们在 MAGGC 模块中引入注意引导的图卷积<sup>[11]</sup>。它利用多头注意力收集所有结点之间的交互信息, 特别是通过多跳路径连接的结点。MAGGC 对边的计算和图卷积操作和 CAGGC 一致, 不同点在于利用多头注意力机制计算邻接矩阵, 使局部联通图变为全联通图:

$$\mathbf{A}^{(2)} = \text{softmax} \left( \frac{(\mathbf{W}_Q \mathbf{P}^{(1)})^T (\mathbf{W}_K \mathbf{P}^{(1)})}{\sqrt{d}} \right) \quad (11)$$

#### ④ 分类层

本文将编码器层、CAGGC 和 MAGGC 的输出拼接起来，用全连接层处理，获得最终的结点表示：

$$\mathbf{P} = \tanh(\mathbf{W}_p[\mathbf{P}^{(0)}; \mathbf{P}^{(1)}; \mathbf{P}^{(2)}] + b_p) \quad (12)$$

在分类前，将最终的结点表示与实体类型嵌入和相对距离嵌入拼接起来，然后将其输入双线性函数和全连接层，得到关系特征，用于关系预测，公式如下：

$$\mathbf{P}'_u = [\mathbf{P}_u; t_u; d_{u,v}] \quad (13)$$

$$\mathbf{P}'_v = [\mathbf{P}_v; t_v; d_{v,u}] \quad (14)$$

$$P(r|u, v) = \text{sigmoid}(\mathbf{P}'_u^T \mathbf{W}_r \mathbf{P}'_v + \mathbf{W}_t[\mathbf{P}'_u; \mathbf{P}'_v] + b_r) \quad (15)$$

DocRED 是一个多关系抽取任务，我们使用二元交叉熵函数作为损失函数进行模型的训练，公式如下：

$$Loss = - \sum_{D \in S} \sum_{u \neq v} \sum_{r_i \in R} \mathbb{I}(r_i = 1) \log P(r_i|u, v) + \mathbb{I}(r_i = 0) \log(1 - P(r_i|u, v)) \quad (16)$$

### 3.2 阶段性成果

#### 3.2.1 系统设置

本文使用 100 维的 GloVe 词向量和 128 维隐层大小的 BiLSTM。使用的 BERT-Base 层数为 12，隐层维度为 768。两种词向量最终都映射到 128 维，位置嵌入和实体类别嵌入维度为 20。MAGGC 和 CAGGC 各有 4 个子层，MAGGC 所用的多头注意力机制头数为 4。采用 Adam 作为优化器，权重衰减为 0.0001，学习率设置为 5e-6。Dropout 概率为 0.2。

#### 3.2.2 实验结果及分析

实验采用广泛使用的指标  $F_1$  值作为评测标准。由于一些关系事实在训练集中出现过，模型可能在训练期间记住了这些关系事实，导致在测试集上的结果不能很好的展示出模型的泛化性能，所以采用了忽视了这些关系事实的  $F_1$  值（记为  $Ign F_1$ ）作为额外的评价指标。

我们比较了所提出的模型和目前的一些先进的文档级关系抽取模型的性能，结果如表 1 所示，在基于 GloVe 和基于 BERT 的模型中，我们基于图卷积网络的模型总体上优于其他模型，这验证了层次推理方法能够区分关键的实体级、句子级和文档级推理信息，实现全面的文档级关系推理，而简单地对文档进行编码不能有效地建模实体之间的复杂关系。基于 BERT 的模型比基于 GloVe 的模型有了很大的性能提升，这表明 BERT 是一种强大的上下文关系编码器。表 1 中的结果证明了本文的模型可以用两个基于图卷积网络的文档级关系推理模块增强

全局上下文表示。

表 1 与其他先进方法的结果的比较

模型		验证集		测试集	
		Ign $F_1$ (%)	$F_1$ (%)	Ign $F_1$ (%)	$F_1$ (%)
1	CNN	41.58	43.45	40.33	42.26
	LSTM	48.44	50.68	47.71	50.07
	BiLSTM	48.87	50.94	48.78	51.06
	ContextAware	48.94	51.09	48.40	50.70
2	GREG-Context	-	-	-	52.88
	HIN-GloVe	51.06	52.95	51.15	53.30
	本文方法-GloVe	51.14	53.05	50.87	53.13
3	BERT-RE	-	54.16	-	53.20
4	BERT-Two-Step	-	54.42	-	53.92
	HIN-BERT	54.29	56.31	53.70	55.60
	本文方法-BERT	55.43	57.35	54.53	56.67

3.3 创新点论述

- 我们提出了新颖的全局上下文增强图卷积网络(GCGCN)，该网络以实体为结点，实体对的上下文作为结点之间的边，以捕获丰富的全局上下文信息。
- 我们提出两个层次推理模块，分别是用于部分联通图的上下文感知注意力引导的图卷积模块（CAGGC）和用于完全连接图的多头注意力引导的图卷积模块（MAGGC），这两个模块可以逐步考虑更多的全局上下文。
- 我们提出的端到端的文档级别关系抽取模型，能够自动识别输入文档内实体对间的复杂关系。在 DocRED 数据集上取得了最优的结果，实验结果证明了方法的有效性

4 后续工作的设想、可能遇到的问题

4.1 下一步研究计划

- （1）为了验证模型的泛化性能，尝试在更多的文档级别关系抽取数据集上进行实验。
- （2）当前在人工标注数据上用有监督的方法训练的模型大多取得了较好的性能。但是手工标注数据所花费的代价过高，这在一定程度上阻碍了文档级别关系抽取的研究，构建远程监督训练语料是此问题的很好的解决办法。DocRED 数据集提供了大规模的远程监督语料，因此考虑利用远程监督语料进一步提升现有模型性能。
- （3）在自然语言处理任务中融合外部知识往往可以产生令人满意的效果。实体关系抽取任务中，尤其是生物医学领域的关系抽取研究，有很多的结构化知识可以利用，很多知识表示模型也被提出，本文考虑将学习到的知识表示融入到

已经提出的模型中，进一步提升文档级别关系抽取的性能。

## 4.2 可能遇到的问题

(1) 由于大规模远程监督语料是基于知识库自动标注的，其中包含了大量的噪声，即：未收录进知识库的关系可能未被标注；在知识库中标注为某种关系的实体对在特定的文本中可能并未体现该种关系。这些噪音对于模型的训练存在一定的误导作用。如何设计合理的降噪机制是当前利用大规模远程监督语料所面临的首要问题。

(2) 预训练语言模型含有丰富的外部知识，并且可以很好的编码文本的语义表示，因此被广泛用于当前各类自然语言处理任务中。但是这些模型参数量十分巨大，在原始语料中直接进行参数微调需要消耗大量的时间和空间资源。我们又考虑利用大规模远程监督数据，这会大大增加模型训练的时间和空间复杂度，造成模型调参困难、效率较低。对模型进行优化是需要考虑的问题。

(3) 将知识表示融合到当前的模型是否可行，当前的模型是否存在缺陷，是否需要引入新的模型方法，还需要进一步的实验研究。

## 5 已发表、录用和已投稿的论文情况

### 5.1 已发表论文情况

[1] Zhou H W, **Xu Y B**, Liu Z, Yao W H, Lang C K, Jiang H B. Global Context-enhanced Graph Convolutional Networks for Document-level Relation Extraction[C]. Proceedings of the 28th International Conference on Computational Linguistics (COLING), pages 5259-5270, 2020.

## 参 考 文 献

- [1] Asma B. Abacha and Dina D. Fushman. Recognizing question entailment for medical question answering. AMIA [J]. Proceeding of the AMIA Symposium, pages 310-318, 2016.
- [2] Bowman S R, Angeli G, Potts C, et al. A large annotated corpus for learning natural language inference [J]. Computer Science, 2015.
- [3] Rocktäschel T, Grefenstette E, Hermann K M, et al. Reasoning about Entailment with Neural Attention [J]. ICLR, 2016.
- [4] Wang S, Jiang J. Learning Natural Language Inference with LSTM [J]. arXiv preprint arXiv:1512.08849, 2015.
- [5] Vaswani A, Shazeer N, et al. Attention Is All You Need [C]. Conference and Workshop on Neural Information Processing Systems (NeurIPS), 2017.
- [6] Jinhyuk L, Wonjin Y, Sungdong K, et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining [J]. arXiv preprint arXiv: 1901.08746, 2019.
- [7] Jacob D, Ming-Wei C, Kenton L, Kristina T. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding [J]. arXiv preprint arXiv: 1801.04805, 2018.
- [8] Jeffrey P, Richard S, and Christopher M. Glove: Global vectors for word representation [C].

- Conference on Empirical Methods in Natural Language Processing, pages 1532–1543, 2014.
- [9] Yao Y, Ye D M, et al. DocRED: A Large-Scale Document-Level Relation Extraction Dataset [C]. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 764–777, 2019.
  - [10] Huang G, Liu Z, et al. Densely Connected Convolutional Networks [C]. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2261-2269, 2017.
  - [11] Guo Z J, Zhang Y, Lu W. Attention Guided Graph Convolutional Networks for Relation Extraction[C]. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics(ACL), pages 241–251, 2019.

**大连理工大学专业学位硕士研究生学位论文中期考核报告评审意见表**

学 号	31909185	学生姓名	徐奕斌	专业/领域	计算机技术
第一次开题 <input checked="" type="checkbox"/>			第二次开题 <input type="checkbox"/>		
实 践 导 师 信 息					
姓 名	周惠巍	性 别	女	职称/职务	副教授
专 业	计算机科学与技术	所在单位	大连理工大学	联系电话	15542679503
通讯地址	辽宁省大连市甘井子区凌工路 2 号 创新园大厦			E_mail	zhouhuiwei@dlut.edu.cn
<p>校内导师考核意见（对校内外实践实习情况、参加学术活动情况、学位论文工作及中期报告撰写情况、企业实践实习情况及计划、学生的学习和工作态度等进行考查，给出考核成绩和具体改进意见和建议）：</p> <p>1) 考核成绩：<input checked="" type="checkbox"/> 优秀，<input type="checkbox"/> 良好，<input type="checkbox"/> 中等，<input type="checkbox"/> 及格，<input type="checkbox"/> 不及格</p> <p>2) 是否通过：<input checked="" type="checkbox"/> 通过，<input type="checkbox"/> 不通过</p> <p>3) 关于中期考核报告撰写质量及学位论文工作的具体意见（可加页）：</p> <p>徐奕斌同学在过去一年半积极参与实验室的科研活动，学习认真努力，科研态度端正，具有良好的工作习惯，目前徐奕斌同学主要研究文档级别关系抽取任务，使用当前效果较好的图卷积网络对文档中的实体之间的交互进行建模，具有一定的研究价值，当前，徐奕斌同学已经完成了数据预处理和基本模型搭建，并在相关评测任务上完成了初步实验，取得了一定的成果，发表一篇会议论文。徐奕斌同学已完成中期报告的撰写，内容充实，同意进行中期答辩。</p>					
导师签字：姚卫红 2020 年 1 月 22 日					

实践导师考核意见（对学位论文工作及中期报告撰写情况、企业实践实习情况及计划、学习和工作态度等进行考查，给出考核成绩和具体改进意见和建议）：

- 1) 考核成绩：☒ 优秀，☐ 良好，☐ 中等，☐ 及格，☐ 不及格
- 2) 是否通过：☒ 通过，☐ 不通过
- 3) 关于中期考核报告撰写质量及学位论文工作的具体意见（可加页）：

徐奕斌同学在过去一年半积极参实习实践活动，学习认真努力，科研态度端正，具有良好的工作习惯，该同学毕业论文研究内容为文档级别关系抽取研究，当前已完成部分科研计划，后续研究正在进行。当前已完成中期报告的书写，内容充实，书写规范，同意进行中期答辩。

导师签字：周惠巍  
2020年 1月 22日

评 议 专 家 组		姓名	职称	学科专业	是否博导	签字
	组长	姚念民	教授	计算机	是	姚念民
	成员	孙亮	副教授	计算机	否	孙亮
		任健康	副教授	计算机	否	任健康
		刘昊	高级工程师	计算机	否	刘昊

专家组评审意见（对课程学习情况、校内实践实习情况、参加学术活动情况、学位论文工作及中期报告撰写情况、企业实践实习情况及计划、学生的学习和工作态度等进行考查，给出考核成绩和具体改进意见和建议）：

- 1) 考核成绩：☒优秀，☐良好，☐中等，☐及格，☐不及格
- 2) 是否通过：☒ 通过，☐ 不通过
- 3) 关于中期考核报告撰写质量及学位论文工作的具体意见（可加页）：

徐奕斌同学在前期的工作中充分分析了课题任务需求，熟练掌握了理论基础，完成了数据预处理和基本模型的搭建工作，所做的研究工作具有重要的理论意义和实际价值。期间该生工作安排合理，中期报告条理清晰，撰写规范，课题进展符合预期计划，可以按照既定计划开展下一步论文工作。

组长签字：姚念民

2021 年 1 月 23 日