

基于宏观语义表示的宏观篇章关系识别方法*

周懿, 褚晓敏, 蒋峰, 李培峰, 朱巧明

(苏州大学 计算机科学与技术学院, 江苏 苏州 215006)

摘要: 宏观篇章分析旨在分析相邻段落或段落群之间的语义联系, 是自然语言处理领域其他任务的基础工作。本文研究了宏观篇章分析中的关系识别问题, 提出了一个宏观篇章关系识别模型。该模型利用基于词向量的宏观篇章语义表示方法和适用于宏观篇章关系识别的结构特征, 从两个层面提高了模型分辨宏观篇章关系的能力。在汉语宏观篇章树库(MCDTB)上的实验表明, 该模型在大类分类中 F1 值达到了 68.22%, 比基准系统提升了 4.17%。

关键词: 宏观篇章关系识别; 宏观篇章结构特征; 宏观篇章语义表示

Macro Discourse-level Relation Classification Based on Macro Semantics Representation

ZHOU Yi, CHU Xiaomin, JIANG Feng, LI Peifeng, ZHU Qiaoming

(School of Computer Sciences and Technology, Soochow University, Jiangsu, Suzhou, 215006)

Abstract: The macro discourse-level analysis aims to analyze the semantic relations between adjacent paragraphs or paragraph groups and is a basic work of other tasks in the field of natural language processing. This paper mainly proposes a classification model to resolve the problem of relational classification in macro discourse-level analysis. This model introduces a distribute representation of macro discourse semantics on word vectors and a set of structure features to improve the performance of classifying macro discourse-level relations in two dimensions. The experimental results on the Macro Chinese Discourse Tree Bank (MCDTB) show that the F1 value of our model reaches 68.22%, with 4.17% improvement.

Key words: macro discourse-level relation classification; structure features of macro discourse; representation of marco discourse

1 引言

随着自然语言处理的发展,其处理的信息的粒度呈现出由细到粗的变化趋势。具体而言,其处理的对象已经从字词等细粒度单元拓展到句子等较粗粒度的单元上。篇章作为一个比句子更大的一种文本分析粒度也愈发受到人们的重视。

篇章分析的主要任务是挖掘篇章单元之间的内在结构和语义关系,此处的篇章单元可以是句子、复句、句群或段落等。篇章分析分为宏观和微观两个层面,微观层面主要研究段落内的句子和连续两个句子间的关系,而宏观篇章分析主要研究段落及更高层次的段落群和章节之间的关系。和词法、句法分析一样,篇章分析作为篇章级的基础研究能够对更高层次的自然语言处理问题,如问答系统^[1]、情感分析^[2]、信息抽取^[3]等提供更有效的支撑。

在篇章分析的任务中,篇章关系识别,尤其是隐式篇章关系识别始终是一个重点和难点。在宏观篇章关系识别的任务中,由于汉语文章的写作方法,段落与段落间很少出现标识宏观语义联系连接词。即使出现连接词,也很难将它们和标识段内关系的微观连接词区分开来。因此,中文文本的宏观篇章关系的识别均是隐式关系的识别,而且相对于一般的隐式篇章识

* 收稿日期:

定稿日期:

基金项目: 国家自然科学基金(No. 61772354, No.61773276, No.61472265); 国防科技先导计划(No. 17-ZLXDXX-02-06-02-04); 江苏省科技计划(No. BK20151222)

作者简介: 周懿(1995——),男,硕士研究生,主要研究方向:自然语言处理;李培峰(1971——),男,教授,硕士生导师,主要研究方向:中文信息处理、自然语言处理;朱巧明(1963——),男,教授,博士生导师,主要研究方向:中文信息处理, Web 信息处理。

别，它有着论元长度长，论元间关系复杂的特点，因而难度更大。本文以 CTB8.0^[4]中的一个篇章(ghtb_0010.nw.raw)来说明宏观篇章之间的关系，如例 1 所示。

分析例 1 的篇章可知，整个篇章的主题就是标题所示的“中国进出口银行在日本获债券信用高等级”这一事件，因而(1)是本文的主题段落。(2)一方面重述了(1)所陈述的事实，同时附加了信息“与日本评级机构内部对中国主权信用等级的评级一致”，并未细化本文的主题内容，只是对(1)起到补充说明的作用。(3)、(4)、(5)段描述(1)中所述事件的详细过程，对段落(1)进行了解说。而在(3)、(4)、(5)段内部，(4)说明了(3)中采取的“向日本评级机构提出评级申请”这一行为的目的，(5)段说明提出申请后进出口银行所采取的一系列行为，与(3)中的内容有着明确的时间上的先后关系。由以上分析，我们可以得到如图 1 所示的篇章间的关系。

例 1 chtb_0010. nw. raw 内容

中国进出口银行在日本获债券信用高等级

(1)新华社北京二月十六日电中国进出口银行最近在日本取得债券信用等级 A A 一，这是日本金融市场当前对中国银行的最高债券评级。

(2)日本公社债研究所确定中国进出口银行债券信用等级为 A A 一，与日本评级机构内部对中国主权信用等级的评级一致。

(3)去年十月，中国进出口银行聘请日本野村证券公司作顾问，向日本著名的评级机构日本公社债研究所提出正式评级申请。

(4)进出口银行决定先在日本取得信用评级是为进入国际资本市场融资创造作准备，以便扩大资金来源，支持中国机电产品和成套设备出口。

(5)进出口银行通过书面介绍、实地考察等形式向日本公社债研究所全面介绍了今年来中国金融体制改革的情况、银行成立的背景、银行管理和运营机制、银行业务发展现状以及以后的发展目标，使之对中国进出口银行有了较深的了解。

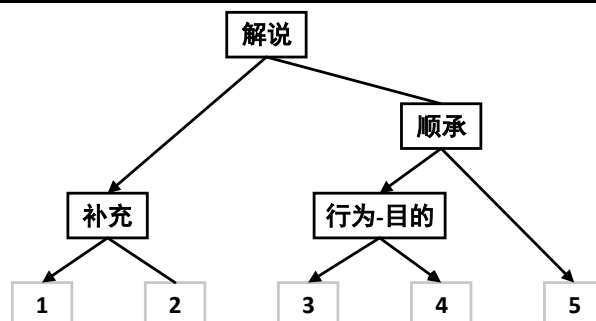


图 1 chtb_0010. nw. raw 宏观结构

本文提出了一种基于词向量的宏观篇章语义表示和一组在宏观篇章关系识别中适用的特征并给出了一个基于该特征的宏观篇章关系识别方法。

2 相关工作

目前篇章关系识别的任务在宏观层面上的研究尚属空白，但在微观的层面上已经有了比较广泛的研究，研究主要涉及基于修辞结构的篇章树库和基于连接依存树的篇章树库这两类语料资源。

2.1 基于修辞结构理论的篇章树库

修辞结构篇章树库(RST-DT)^[5]是以 Mann 和 Thompson^[6,7]提出的修辞结构理论(RST)为理论支撑的篇章树库。RST 提出了“命题-证据”的关系模式，其中，命题是涵盖了作者陈述的观点的篇章，这个观点读者不一定认同，而证据是为了为命题提供支撑的篇章。据此，RST-DT 树库标注了 16 种关系大类和 78 种小类。同时，还标注了篇章单元，“核-卫星”模式的主次类型，篇章结构等，将文本组织成了层次化的篇章结构树。

在 RST-DT 树库上, Hernault^[8]等提出了 HILDA 分析器, HILDA 分析器使用两个支持向量机分别进行篇章单元识别和主次-关系标签标注, 实现了一个自底向上构建自动篇章树的框架, 在篇章关系识别的任务上得到了 50.9%的 F1 值。Joty^[9,10]等认识到句内和句间的关系分布上有差异, 使用了两个动态条件随机场模型针对句内和句间关系分别建模, 并使用动态规划算法对篇章树的构建进行优化, 在篇章关系识别的任务上得到了 55.73%的 F1 值。Feng 和 Hirst^[11,12]认识到篇章结构对于关系识别的重要性, 提出了先识别篇章结构再识别篇章关系的两步走策略, 使用每组两个的两组线性条件随机场模型, 在篇章关系识别上获得了 58.2%的正确率。Wang^[13]等使用基于转移的方法将篇章树构建转化成 shift-reduce 序列, 提出了先标注结构-主次, 再进行标签标注的两步模型, 在篇章关系识别上获得了 59.7%的正确率。

2.2 基于连接依存树的篇章树库

相比修辞结构理论, 基于连接依存树的体系借鉴了“谓词—论元”的模式, 凸显了连接词的作用, 以连接词为核心标注与之相关的篇章单元, 依据有无连接词将篇章关系分为显式篇章关系和隐式篇章关系, 代表性的语料资源有宾州篇章树库(PDTB)^[14], 在汉语上有汉语篇章树库(CDTB)^[15]等。

在 PDTB 上, Lin^[16]等探索了各种上下文特征、词对特征、句法特征等, 对隐式篇章关系进行识别, 得到了 40.2%的正确率。Park^[17]等通过特征集优化算法对特征进行选择, 分类的性能有所提高。Qin^[18]等尝试将对抗生成网络用于篇章分析, 构建了对抗生成模型来从隐式篇章关系中获得包含隐藏连接词的段落表示, 获得了 44.61%的正确率。

在 CDTB 上, 李艳翠^[19]等构建了基于连接依存树的汉语篇章分析平台, 在微观篇章关系识别上, 同时考虑显式和隐式篇章关系的情况下在句内和句间的分类任务上分别达到了 78.4%和 69.6%的正确率。Kong^[20]等建立了一个端到端的篇章结构分析器, 在“解说”, “并列”, “因果”三大类的分类上分别取得了 51.8%, 85.8%, 57.1%的 F1 值。

3 宏观篇章树库 (MCDTB)

宏观篇章分析的任务在国内的研究目前还处于起步阶段。Jiang^[21], Chu^[22]等在宾州汉语树库(CTB)的基础上标注了语料的宏观结构, 形成了汉语宏观篇章树库(MCDTB)。

表 1 MCDTB 语料库中宏观关系类型分布

关系大类类型	关系小类类型	关系小类数目	关系小类占比	关系大类数目	关系大类占比
解说类	解说关系	994	30.44%	1203	36.85%
	陈述举例关系	39	1.19%		
	举例陈述关系	20	0.61%		
	评价关系	127	3.89%		
	总结关系	23	0.70%		
并列类	并列关系	1003	30.72%	1648	50.47%
	顺承关系	125	3.83%		
	递进关系	10	0.31%		
	对比关系	17	0.52%		
	补充关系	493	15.10%		
因果类	因果关系	49	1.50%	414	12.68%
	果因关系	103	3.15%		
	行为目的关系	14	0.43%		
	目的行为关系	11	0.34%		
	背景关系	237	7.26%		

MCDTB 以段落为基本篇章单元, 使用自底向上的方式对段落及段落以上的篇章从篇章

主题，段落主题，篇章摘要，篇章结构，篇章主次，篇章关系等方面进行了标注。总共标注了 720 篇文章的 2870 个关系。在类别上，MCDTB 将这些关系分为“解说类”、“并列类”、“因果类”，然后进一步细分为 15 个小类，具体分布如表 1 所示。

4 宏观篇章关系识别方法

本文进行的任务是在 MCDTB 上进行三个大类的划分。依据 MCDTB 的标注方式，本文将篇章结构树的结构视为已知条件进行关系分类。由于 MCDTB 中的关系分为二元关系和多元关系，在本文中，二元关系表示为一个元组([Arg1,Arg2],Label)，而多元关系表示为([Arg1,Arg2,...,Argn],Label)。参照 RST-DT 上的一些研究，本文将多元关系以右连接的方式转化为二元关系，例如对于多元关系([a,b,c],Label)，转化之后为([a,c],Label)，([b,c],Label)两个元组。最终问题转化为对二元关系进行“解说类”、“并列类”、“因果类”的三分类问题。

由于宏观篇章分析分析的是段落及更高层次的篇章单元之间的关系，在进行微观篇章分析时常用的语法、句法信息很难被有效地利用。单个词和词性相对宏观篇章单元而言粒度过小，难以表示篇章本身的语义和篇章之间的语义关系。本文认为在宏观篇章分析的时候应当考虑粒度更大的特征，提出了一种基于词向量的宏观篇章单元表示方法和一组用于宏观篇章关系识别的特征。

4.1 基于词向量的宏观篇章单元表示方法

基于词向量的宏观篇章单元表示方法通过词向量训练算法训练得到词向量模型 \mathbf{WV} ，从中获取篇章单元中每个词的词向量 \mathbf{wv}_i ，再通过式(1)计算得到整个篇章单元表示。

$$\mathbf{Discourse}_{\mathbf{WV}} = \frac{\sum_{i \in \mathcal{W}} \mathbf{wv}_i}{N} \quad (1)$$

其中， \mathcal{W} 是篇章单元中所有词的集合， \mathbf{wv}_i 是 \mathcal{W} 中第 i 个词在词向量模型 \mathbf{WV} 中的表示。 N 表示篇章单元中词的数量。

考虑到目前常用的词向量算法中，Word2Vec^[23]能很好地表示词语的局部信息，而 GloVe^[24]同时考虑了局部信息和全局词共现信息，正与宏观篇章关系识别任务中既要考虑两段落间内容的关系，同时统筹考虑全文主题的要求相一致，本文通过式(2)得到两种词向量间的差异，意图表示与该词相关的全局信息，最后通过结合篇章的局部表示和全局信息得到式(3)为篇章单元最终的向量表示。

$$\mathbf{Global}_i = \mathbf{GloVe}_i - \mathbf{Word2Vec}_i \quad (2)$$

$$\mathbf{Discourse} = \frac{\sum_{i \in \mathcal{W}} (\mathbf{Word2Vec}_i + \lambda \times \mathbf{Global}_i)}{N} \quad (3)$$

其中， \mathbf{GloVe}_i 和 $\mathbf{Word2Vec}_i$ 分别表示第 i 个词在 GloVe 模型和 Word2Vec 模型下的表示， λ 是全局信息的权重参数。

4.2 宏观篇章关系识别的结构特征

在宏观篇章的关系识别任务上，目前还没有可供参考的研究，在结构特征上，本文整合了 Jiang^[21]等在宏观篇章主次识别时使用的和 Hernault 的 HILDA^[8]中做微观篇章分析时使用的特征中在宏观篇章识别关系的任务上最有效的特征集，并把 Feng^[12]等在后剪辑时使用的节点所处的深度信息也作为结构特征来使用。

上述三人的工作在将原来的篇章结构树转化为二叉树后不再考虑转化前的多叉树的结构，本文认为树本来的结构对于篇章关系的分类，尤其是对并列类和其他两类的区分有着至关重要的作用。因此将二叉化前树的结构也作为结构特征来使用。

基于上述讨论，本文最终使用了表 2 所示的 5 组特征，其中，基础组织结构特征是

Jiang,Hernault,Feng 等先前的研究中使用的特征, *originalStructure*是二叉化前树的结构特征, *Vec_{w2v}*是使用 Word2Vec 训练的仅考虑局部信息的宏观篇章语义表示, *Vec_{global}*是仅考虑全局信息的宏观篇章语义表示, *Vec_{w2v+global}*是加上了全局信息补正之后的宏观篇章语义表示。

表 2 本文使用的特征

基础组织结构特征 (9 个)
Arg1 的开始位置和结束位置
Arg2 的开始位置和结束位置
关系的开始段落距离文章开头的段落数
关系的结束段落距离文章结尾的段落数
Arg1 和 Arg2 包含句子数的大小关系
关系节点所处深度
<i>originalStructure</i> (2 个)
二叉化之前该节点拥有的子节点数目
二叉化之前该节点拥有的非叶子节点数目
<i>Vec_{w2v}</i> (2 个)
Arg1 的宏观篇章表示
Arg2 的宏观篇章表示
<i>Vec_{global}</i> (2 个)
Arg1 的全局信息表示
Arg2 的全局信息表示
<i>Vec_{w2v+global}</i> (2 个)
加上全局信息补正的 Arg1 的宏观篇章表示
加上全局信息补正的 Arg2 的宏观篇章表示

5 实验

5.1 实验设置

本文使用 python 的 sklearn 包提供的 SVC 分类器¹, 参数均使用默认值, 篇章单元的词向量表示使用中文维基语料, 经由 Word2Vec 和 GloVe 训练成 50 维词向量, 训练时窗口大小为 5。数据集大小为 MCDTB 的全部 720 篇文章, 二叉化后共 3265 条关系。

考虑到样本集相对较小, 实验采用五倍交叉验证的方式, 将 720 篇文章按段落数平分为五份, 如有 58 篇 7 段的文章, 则给每个样本集 11 篇, 再将剩下三篇随机分派给三个样本集。然后将五个样本集中的一个作为测试集, 其他作为训练集, 共进行五次实验。在训练集 1 中, 又将其划分成五份, 使用其中四份作为训练集, 一份作为验证集, 进行参数选择, 最终将公式(4)中的全局信息权重 λ 调整为 2。

本文选取五组特征集组合来进行, 基准系统使用表 3 中的基础组织结构特征, 第 2 组在基准系统的基础上附加二叉化之前的结构树特征 *originalStructure*, 第 3、4、5 组分别在基准系统的基础上附加词向量特征 *Vec_{w2v}*、*Vec_{global}* 和 *Vec_{w2v+global}*, 第 6 组同时使用基础组织结构特征, 加上全局信息补正的词向量特征和二叉化之前的结构树特征。

本文使用的测评指标为正确率(Accuracy)、准确率(Precision)、召回率(Recall)和 F1 值(F1-Score), 其中每个类别的测评指标按照标准的正确率, 准确率, 召回率和 F1 值的计算公式计算, 整体性能的测评指标分别由式(4)、(5)、(6)、(7)计算所得。

¹ <http://scikit-learn.org>

$$Accuracy = \frac{TP}{N} \quad (4)$$

$$Precision = \frac{\sum_{c \in C} Precision(c) \times support(c)}{N} \quad (5)$$

$$Recall = \frac{\sum_{c \in C} Recall(c) \times support(c)}{N} \quad (6)$$

$$F1-Score = \frac{\sum_{c \in C} F1-Score(c) \times support(c)}{N} \quad (7)$$

其中 TP 表示五次实验中分类正确的样本总数， N 表示样本集所有样本的总数， $Precision(c)$ 、 $Recall(c)$ 、 $F1-Score(c)$ 分别表示类型 c 的准确率、召回率和F1值， $support(c)$ 表示样本集中属于类型 c 的样本数量。

5.2 实验结果

从表3中可以看到，使用了特征集中所有特征的第六组实验相比基准系统在正确率，准确率，召回率，F1值上分别有了4.08%，6.27%，4.08%和4.17%的提升，在六组实验中四项指标均达到了最优。

表3 实验结果

编号	特征集	正确率	准确率	召回率	F1值
1	基准系统	67.44%	63.23%	67.44%	64.05%
2	基准系统+ <i>originalStructure</i>	68.58%	63.99%	68.58%	65.09%
3	基准系统+ <i>Vec_{w2v}</i>	69.53%	66.67%	69.53%	65.63%
4	基准系统+ <i>Vec_{global}</i>	69.71%	67.53%	69.71%	66.00%
5	基准系统+ <i>Vec_{w2v+global}</i>	70.02%	68.03%	70.02%	66.71%
6	基准系统+ <i>originalStructure +Vec_{w2v+GloVe}</i>	71.52%	69.50%	71.52%	68.22%

从表3中1, 2两组实验结果的对比和1, 3两组实验结果的对比可见本文提出的二叉化之前树的结构特征以及宏观篇章的词向量表示对于宏观篇章的关系识别均有积极作用。而从3、4、5三组实验的对比则能看出全局信息和局部词向量信息间的相互补充。

本文还对特征集2和5的两个模型的预测结果进行了配对样本t检验，结果显示两个模型具有显著差异($p < 0.01$)，说明二叉化前树的结构特征和词向量表示的宏观语义分别从两个层面对基准系统做出了优化。为探究本文提出的两个特征分别对基准系统在哪几个方面进行了优化，本文取出五折交叉验证实验中的一组，在这组样本上，六个特征集在三个类别上的具体表现如表4所示。

对比表4中的实验数据，可以得出以下结论：

- (1) 对比1, 2两组实验结果可知，二叉化前树的结构主要起到了提高解说类召回率和并列类准确率的作用，即减少了解说类中误分到并列类中的样本。这是因为并列类中包含许多多元关系，而解说类中以二元关系为主。
- (2) 对比1, 3两组实验结果可知，词向量表示的宏观语义提高了因果类的召回率和解说类的准确率，使得很多因果类关系从解说类中区分出来，这是因为因果类中的关系具有更强的语义上的连贯性，本文提出的宏观语义表示方法表达了篇章的语义信息，对于识别因果类关系有帮助。
- (3) 对比3, 4, 5三组实验可知，结合了全局信息和局部信息的模型比起单独使用其中一个取得了更好的效果。结合的模型一方面保持了全局信息对因果类关系的识别率，

另一方面，进一步提高了在并列类识别的正确率，说明局部语义和全局信息间有着互相补充、互相约束的关系。

表 4 五组特征在测试集 2 上的实验结果

类别	特征集编号	正确率	准确率	召回率	F1 值
解说类	1	67.67%	60.62%	67.67%	63.95%
	2	75.00%	61.92%	75.00%	67.84%
	3	68.97%	62.99%	68.97%	65.84%
	4	68.10%	63.45%	68.10%	65.70%
	5	68.10%	63.97%	68.10%	65.97%
	6	68.97%	65.04%	68.97%	66.95%
并列类	1	81.73%	73.13%	81.73%	77.19%
	2	81.11%	77.06%	81.11%	79.03%
	3	82.97%	73.22%	82.97%	77.79%
	4	83.59%	73.77%	83.59%	78.37%
	5	84.21%	74.73%	84.21%	79.18%
	6	85.14%	75.76%	85.14%	80.17%
因果类	1	5.06%	28.57%	5.06%	8.60%
	2	5.06%	30.77%	5.06%	8.70%
	3	8.86%	50.00%	8.86%	15.05%
	4	11.39%	47.37%	11.39%	18.37%
	5	11.39%	39.13%	11.39%	17.65%
	6	13.92%	44.00%	13.92%	21.15%

同时，表 4 还反映出在不同的类别上，本文提出的模型表现的差异也较大。即使在最佳的第 5 组中，因果类的表现仍是比较差的。究其原因，这一方面是因为样本集不平衡，因果类的样本数相较其他两类少很多；另一方面，从上述分析也可看出，因果类关系对于语义信息极为敏感，仅其中占比最多的背景关系而言，与解说类关系的区别仅在于是对事物本身的属性进行解说还是对事物相关的环境要素进行解说，是很难通过仅由词汇集成的语义来完全区分的。

6 总结与展望

本文提出了一种基于词向量的宏观篇章语义表示方法和一组适用于宏观篇章关系类型识别的结构特征，并在 MCDTB 语料库上进行了一系列实验。实验结果证明，在宏观篇章关系识别的任务上，本文提出的二叉化前树的结构特征提高了解说类和并列类关系的区分度，而基于词向量的宏观篇章表示方法提高了解说类和因果类关系的区分度，在两个不同的维度上为系统做出了贡献。在将来的工作中，一方面我们将进一步探究宏观篇章的语义表示，寻找类似于微观篇章分析时地句法信息等更高层面的宏观语义特征，另一方面将寻找方法解决样本集不平衡带来的问题，并在此基础上开展更细粒度的小类识别工作。

参考文献

- [1] Liakata M, Dobnik S, Saha S, et al. A discourse-driven content model for summarising scientific articles evaluated in a complex question answering task[C]//Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing. 2013: 747-757.
- [2] Zhou L, Li B, Gao W, et al. Unsupervised discovery of discourse relations for eliminating intra-sentence polarity ambiguities[C]//Proceedings of the Conference on Empirical Methods in Natural Language Processing. Association

- for Computational Linguistics, 2011: 162-171.
- [3] Huttunen S, Vihavainen A, Yangarber R. Relevance prediction in information extraction using discourse and lexical features[J]. 2011.
- [4] Xue N, Xia F, Chiou F D, et al. The Penn Chinese TreeBank: Phrase structure annotation of a large corpus[J]. Natural language engineering, 2005, 11(2): 207-238.
- [5] Carlson L, Okurovski M E, Marcu D. RST discourse treebank[M]. Linguistic Data Consortium, University of Pennsylvania, 2002.
- [6] Mann W C, Thompson S A. Relational propositions in discourse[J]. Discourse processes, 1986, 9(1): 57-90.
- [7] Mann W C, Thompson S A. Rhetorical structure theory: A theory of text organization[J]. Text-Interdisciplinary Journal for the Study of Discourse, 1987, 8(3):243-281.
- [8] Hernault H, Prendinger H, Ishizuka M. HILDA: A discourse parser using support vector machine classification[J]. Dialogue & Discourse, 2010, 1(3).
- [9] Joty S, Carenini G, Ng R T. A novel discriminative framework for sentence-level discourse analysis[C]//Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. Association for Computational Linguistics, 2012: 904-915.
- [10] Joty S, Carenini G, Ng R, et al. Combining intra-and multi-sentential rhetorical parsing for document-level discourse analysis[C]//Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2013, 1: 486-496.
- [11] Feng V W, Hirst G. Text-level discourse parsing with rich linguistic features[C]//Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1. Association for Computational Linguistics, 2012: 60-68.
- [12] Feng V W, Hirst G. A linear-time bottom-up discourse parser with constraints and post-editing[C]//Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2014, 1: 511-521.
- [13] Wang Y, Li S, Wang H. A two-stage parsing method for text-level discourse analysis[C]//Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). 2017, 2: 184-188.
- [14] PDTB Research Group. The penn discourse treebank 2.0 annotation manual[J]. December, 2007, 17: 26-37.
- [15] Li Y, Kong F, Zhou G. Building Chinese discourse corpus with connective-driven dependency tree structure[C]//Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP). 2014: 2105-2114.
- [16] Lin Z, Kan M Y, Ng H T. Recognizing implicit discourse relations in the Penn Discourse Treebank[C]//Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1-Volume 1. Association for Computational Linguistics, 2009: 343-351.
- [17] Park J, Cardie C. Improving implicit discourse relation recognition through feature set optimization[C]//Proceedings of the 13th Annual Meeting of the Special Interest Group on Discourse and Dialogue. Association for Computational Linguistics, 2012: 108-112.
- [18] Qin L, Zhang Z, Zhao H, et al. Adversarial connective-exploiting networks for implicit discourse relation classification[J]. arXiv preprint arXiv:1704.00217, 2017.
- [19] 李艳翠. 汉语篇章结构表示体系及资源构建研究[D]. 苏州: 苏州大学博士学位论文, 2015.
- [20] Kong F, Wang H, Zhou G. A CDT-Styled End-to-End Chinese Discourse Parser[M]//Natural Language Understanding and Intelligent Applications. Springer, Cham, 2016: 387-398.
- [21] 蒋峰, 褚晓敏, 徐昇, 等. 基于主题相似度的宏观篇章主次关系识别方法[J]. 中文信息学报, 2018, 32(1): 43-50.

- [22] Chu, X., Jiang, F., Xu, S., Zhu, Q.: Building a macro chinese discourse treebank. In: In the 11th edition of the Language Resources and Evaluation Conference (LREC). The European Language Resource Association (2018)
- [23] Mikolov T, Chen K, Corrado G, et al. Efficient estimation of word representations in vector space[J]. arXiv preprint arXiv:1301.3781, 2013.
- [24] Pennington J, Socher R, Manning C. Glove: Global vectors for word representation[C]//Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014: 1532-1543.

作者联系方式:

周懿, 江苏省苏州市干将东路 333 号苏州大学, 215006, 18506123318, yzhou0928@stu.suda.edu.cn