

蒙古文依存句法分析

苏向东 高光来 闫学亮

(内蒙古大学计算机学院 呼和浩特 010021)

摘要 近年来,依存句法分析逐渐成为自然语言处理领域中的研究热点。然而,蒙古文的依存句法分析尚未得到足够的重视。基于最大生成树模型在蒙古文依存关系树库 TMDT 上进行了蒙古文依存句法分析的研究。在简要介绍蒙古文的特点和蒙古文依存关系树库 TMDT 之后,详细讨论了最大生成树模型。为找到该模型在蒙古文依存句法分析中合适的特征,重点通过实验对 8 种特征及其组合在句法分析中的性能进行了比较。结果显示,Basic Unigram Features、Basic Bi-gram Features 以及 C-C sibling Features 这 3 种特征的组合性能最佳。本研究为蒙古文依存句法分析奠定了基础。

关键词 蒙古文,依存句法分析,最大生成树,自然语言处理

中图分类号 TP391 **文献标识码** A **DOI** 10.11896/j.issn.1002-137X.2014.08.021

Dependency Parsing for Traditional Mongolian

SU Xiang-dong GAO Guang-lai YAN Xue-liang

(College of Computer Science, Inner Mongolia University, Huhhot 010021, China)

Abstract Dependency parsing has become increasingly popular in natural language processing in recent years. Nevertheless, dependency parsing focused on traditional Mongolian has not attracted much attention. We investigated it with Maximum Spanning Tree (MST) based model on Traditional Mongolian dependency treebank (TMDT). This paper briefly introduced traditional Mongolian along with TMDT, and discussed the details of MST. Much emphasis was placed on the performance comparisons among eight kinds of features and their combinations in order to find a suitable feature representation. Evaluation result shows that the combination of Basic Unigram Features, Basic Bi-gram Features and C-C Sibling Features obtains the best performance. Our work establishes a baseline for dependency parsing of traditional Mongolian.

Keywords Traditional mongolian, Dependency parsing, Maximum spanning tree, Natural language processing

1 引言

在自然语言处理领域中,句法分析是许多与文档分析相关的的关键步骤,诸如:知识抽取、问答系统、自动文摘、情感挖掘等。作为句法分析的一个分支领域,依存句法分析近年来逐渐成为一个研究热点。然而,蒙古文依存句法分析尚未得到足够的重视。主要原因如下:首先,针对蒙古文自然语言处理的研究起步较晚。其次,作为依存句法分析所必需的标注数据集,蒙古文依存关系树库 TMDT 直到 2011 年才出现。另外,蒙古语的语法和其他语言有较大不同。这些不同体现在语序、格、时态、语态、领属等方面。这使得蒙古文句法分析中依存关系的确定和分析比较复杂。

本文的出发点是在蒙古文依存关系树库 TMDT 上,评测最大生成树模型在蒙古文依存句法分析中的性能,并找到该模型在蒙古文依存句法分析中最佳的特征表示。基于最大生

成树模型的依存句法分析方法是寻找最可能的依存关系问题转化为最大生成树的构建问题。本文首先介绍了蒙古文以及实验的数据源蒙古文依存关系树库 TMDT,接着详细讨论了最大生成树模型并通过实验重点比较了 8 种特征及其组合在该模型中的依存句法分析的性能。这 8 种特征包括:Basic Unigram Features、Basic Bi-gram Features、In Between POS Features、Surrounding Word POS Features、Extended Surrounding Word POS Features、C-P-C Triple Features、C-C Sibling Features 以及 Word Stem Features。实验表明,基于最大生成树模型的蒙古文依存句法分析方法获得了满意的性能;同时,Basic Unigram Features、Basic Bi-gram Features 和 C-C Sibling Features 这 3 种特征的组合在句法分析中的性能最佳。

本文第 2 节回顾了依存句法分析方面的工作;第 3 节描述了蒙古文及其依存关系树库 TMDT;第 4 节讨论了基

到稿日期:2013-05-26 返修日期:2013-07-15 本文受国家自然科学基金资助项目(61263037),内蒙古自然科学基金重大项目(2011ZD11)资助。
苏向东(1984-),男,博士生,助理工程师,主要研究方向为蒙古文信息处理、模式识别与人工智能,E-mail: sxddxs5747@sina.com;高光来(1964-),男,教授,博士生导师,主要研究方向为蒙古文信息处理、模式识别与人工智能、信息检索;闫学亮(1984-),男,硕士,主要研究方向为模式识别与人工智能、信息检索。

因为句子 x 中每个单词 x_j 都是选择对应依存关系分数 $s(i, j)$ 最大的单词 x_i 作为父节点, 所以构建的依存关系树 T_x 的分数 $s(x, T_x)$ 一定是句子 x 所有可能的依存关系树中分数最大的一个。根据生成树的定义可知, 依存关系树也是生成树。因此, 我们可以采用最大生成树模型进行依存句法分析。表 1 列出了本文实验所使用的最大生成树求解算法, 即 Chu-Liu-Edmonds 算法^[5,6]。算法伪代码中 Contract 和 Expand 操作的具体步骤见参考文献[13]。

表 1 Chu-Liu-Edmonds 算法

Chu-liu-Edmonds(G_x, s)
1. 初始化 $G=(V, E), V=V_x, E=\emptyset$
2. 如果 G 是生成树, 返回 G
3. 否则为 V 中每个顶点贪婪查找分数最高的入边, 加入 E 中
4. 如果在 G 中出现圈 C , 对 G 执行 Contract 操作变为 G_c , 令 $y=Chu-liu-Edmonds(G_c, s)$, 利用 C 对 y 执行 Expand 操作, 返回 y'
5. 否则返回第 2 步

现在的问题是如何定义分数 $s(i, j)$ 。因为分数 $s(i, j)$ 决定了单词 x_i 和单词 x_j 是否被认定存在依存关系且在依存关

系中单词 x_i 是单词 x_j 的父节点。句子中依存关系主要与 (1) 句子结构、(2) 单词 x_i 和单词 x_j 本身以及 (3) 单词 x_i 和单词 x_j 的语境信息有关。因此, 我们给出分数 $s(i, j)$ 的定义如下:

$$s(i, j) = w \cdot f(i, j) \quad (2)$$

其中, $f(i, j)$ 是反映上面 3 个因素的单词 x_i 和 x_j 单词各自的特征和组合的特征; w 是反映特征权重的系数向量。本文第 5 节给出了实验中使用的 8 种特征的具体定义。权重系数向量 w 在句法分析器的训练阶段利用边界最大化算法^[7,8] 通过迭代学习获得。

5 特征

本文基于最大生成树模型对蒙古文进行依存句法分析。模型中的特征直接影响句法分析的性能。为了找到合适的特征, 实验中对 8 种特征及其组合的性能进行了比较和分析。图 2 列出了实验中所用的 8 种特征。

Basic Unigram Features	Basic Bi-gram Features	In between POS Features
p-word	p-word, c-word	p-pos, b-pos, c-pos
p-pos	p-pos, c-pos	
p-word, p-pos	p-word, c-pos	Surrounding Word POS Features
c-word	p-pos, c-word	p-pos-1, p-pos, c-pos-1, c-pos
c-pos	p-word, p-pos, c-word	p-pos-1, p-pos, c-pos, c-pos+1
c-word, c-pos	p-word, p-pos, c-pos	p-pos, p-pos+1, c-pos-1, c-pos
	p-word, c-word, c-pos	p-pos, p-pos+1, c-pos, c-pos+1
C-C Sibling Features	p-pos, c-word, c-pos	Extended Surrounding Word POS Features
c1-word, c2-word	p-word, p-pos, c-word, c-pos	p-pos-1, p-pos, p-pos+1, c-pos-1, c-pos, c-pos+1
c1-word, c2-pos		
c1-pos, c2-word	C-P-C Triple Features	Word Stem Features
c1-pos, c2-pos	p-pos, c1-pos, c2-pos	p-stem, p-pos, c-stem, c-pos

p-word: 依存关系树中的父节点对应的单词, c-word: 子节点对应的单词, p-pos: 父节点对应的单词的词性, c-pos: 子节点对应的单词的词性, p-pos+1: 父节点右边单词的词性, p-pos-1: 父节点左边单词的词性, c-pos+1: 子节点右边单词的词性, c-pos-1: 子节点左边单词的词性, b-pos: 父、子节点中间单词的词性

图 2 依存句法分析中全部特征

Basic Unigram Features 和 Basic Bi-gram Features 代表依存关系中父子节点对应的单词本身和单词的词性; In between POS Features、Surrounding Word Pos Features 和 Extended Surrounding Word Pos Features 代表依存关系中父子节点对应的单词的语境信息; C-C Sibling Features 和 C-P-C Sibling Features 代表依存关系中父节点以及同一父节点的相邻子女对应的单词的信息; Word Stem Features 代表的是依存关系中父子节点对应单词的词干信息。

6 实验

本文的出发点是测试最大生成树模型在蒙古文依存句法分析中的性能, 并找出适合蒙古文依存句法分析的最优特征。本节给出了实验结果并进行了分析。

6.1 实验结果

在蒙古文依存句法分析中, 本文采取了 5 组交叉验证的测试方法。评价句法分析性能的 4 个指标分别为: 根精度 (RA)、带标记依存关系精度 (LAS)、完全匹配精度 (CM) 和无标记依存关系精度 (DA)。这 4 个指标的定义如下:

$$RA = \frac{\text{根被正确识别的句子数目}}{\text{全部句子的数目}} \quad (3)$$

$$LAS = \frac{\text{正确标注了父节点和依存类别的单词的数目}}{\text{单词的数目}} \quad (4)$$

$$CM = \frac{\text{依存分析完全正确的句子的数目}}{\text{全部句子的数目}} \quad (5)$$

$$DA = \frac{\text{正确标注了父节点的单词的数目}}{\text{所有单词的数目}} \quad (6)$$

为了表示方便, 我们使用单个字母来代表每个特征。其中:

a: Basic Unigram Features

b: Basic Bi-gram Features

c: In between POS Features

d: Surrounding Word POS Features

e: Extended Surrounding Word POS Features

f: C-C Sibling Features

g: C-P-C Triple Features

h: Word Stem Features

句法分析实验分为带类别标注和无类别标注两个子类。带类别标注的子类的评测指标为 RA、LAS 和 CM; 无类别标注的子类的评测指标为 RA、DA 和 CM。根据特征所表示的信息的不同, 进一步将带类别标注和无类别标注句法分析实验各分为 5 组, 如表 2 所列: 1-3, 4-10, 11-13, 14 和 15-17。

表2 8种特征及其组合在蒙古文依存句法分析中的性能

Seq.	Features	Unlabeled			Labeled		
		RA	DA	CM	RA	LAS	CM
1	a	73.0%	83.1%	19.0%	72.0%	80.4%	13.0%
2	b	72.0%	81.6%	24.0%	71.0%	79.5%	17.0%
3	a+b	74.0%	83.4%	28.0%	73.0%	81.0%	20.0%
4	a+b+c	67.0%	81.7%	21.0%	66.0%	79.2%	13.0%
5	a+b+d	68.0%	84.2%	27.0%	67.0%	82.0%	19.0%
6	a+b+e	67.0%	82.2%	28.0%	66.0%	80.1%	20.0%
7	a+b+c+d	70.0%	82.5%	21.0%	69.0%	80.0%	13.0%
8	a+b+c+e	70.0%	82.2%	21.0%	69.0%	79.8%	14.0%
9	a+b+d+e	70.0%	83.1%	26.0%	69.0%	80.9%	19.0%
10	a+b+c+d+e	73.0%	82.4%	20.0%	73.0%	79.8%	14.0%
11	a+b+f	74.0%	85.0%	28.0%	73.0%	82.6%	19.0%
12	a+b+g	73.0%	84.4%	27.0%	73.0%	82.1%	19.0%
13	a+b+f+g	70.0%	84.2%	29.0%	69.0%	81.8%	20.0%
14	a+b+h	66.0%	82.4%	28.0%	65.0%	80.3%	22.0%
15	c+d+e+f+g+h	72.0%	82.5%	27.0%	71.0%	79.8%	16.0%
16	a+b+d+f+h	67.0%	83.8%	28.0%	66.0%	81.5%	20.0%
17	a+b+c+d+e+f+g+h	77.0%	83.2%	23.0%	76.0%	81.0%	18.0%

6.2 分析

表2列出了蒙古文依存句法分析的实验结果。我们把第3组实验的结果作为基准,把性能相同或优于它的结果用粗体字来显示。在第11组的实验中,DA和LAS达到了最大值,分别为85.0%和82.6%,其他4项指标接近于最优性能。因此,Basic Unigram Features、Basic Bi-gram Features和C-C sibling Features这3种特征的组合在蒙古文依存句法分析中整体性能最佳。

4个度量指标的一致性也可以从表2观察到。DA与LAS比较一致,而DA与RA、CM这两者不完全一致。也就是说,DA和LAS的上升和下降是同步的,而DA与RA、CM的升降是不同步的。比较第3组和第13组实验可以看到,DA增大了而RA减小了。在第17组(所有特征的组合)带类别标注和无类别标注的实验中RA均取得最大值,而CM并未取得最大值。部分原因是蒙古文语序为SOV结构,句中非动词类单词可以作谓语,这很容易与其他一些依存关系混淆。

相比其他组实验,第14组带类别标注实验中CM达到最大值。但这并不能证明词干特点有助于提高蒙古文依存句法分析的性能。因为相比于第3组实验,它的其他指标RA、DA和LAS均有所下降。在第13组无类别标记实验中CM达到最大值。

Basic Unigram Features和Basic Bi-gram Features是所有特征中的最重要的两种特征。首先,它们的组合(第3组实验)性能优于所有其他特征的组合(第15组实验)。其次,它们组合的性能接近于最优性能。

从第4组实验到第10组实验我们观察到,增加上下文特征并未能明显提升依存句法分析的性能,反而在一些情况下降低了性能。这说明蒙古文依存句法分析中父子节点对应的单词的语境信息对依存关系的判定作用并不明显。这也证实了对基于最大生成树模型的句法分析器而言,特征的选择非常重要。

结束语 本文基于最大生成树模型在蒙古文依存树库TMDT上进行了蒙古文依存句法分析的研究。文中详细讨论了基于最大生成树模型的依存句法分析方法,重点通过实验比较了8种特征及其组合在蒙古文依存句法分析中的性能,并进行了相应的分析。实验结果显示,Basic Unigram Features、Basic Bi-gram Features和C-C sibling Features这3种特征的组合在蒙古文依存句法分析中的性能最佳。无标记依存关系精度DA和带标记依存关系精度LAS分别达到了

85.0%和82.6%。考虑到蒙古语语法与其他语言的显著差异以及训练数据的有限性,这个结果是令人满意的。

下一步,我们将尝试在句法分析器中加入蒙古语特性知识,用以提高依存句法分析的精度。

参考文献

- [1] Hideki H. Semantic Dependency Analysis Method for Japanese Based on Optimum Tree Search Algorithm[J]. Transactions of Information Processing Society of Japan, 2002, 43(3): 696-707
- [2] McDonald R, Crammer K, Pereira F. Online Large-Margin Training of Dependency Parsers[C]// Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics. 2005: 91-99
- [3] Eisner J M. Three New Probabilistic Models for Dependency Parsing: an Exploration[C]// Proceedings of the 16th Conference on Computational Linguistics. 1996: 340-345
- [4] McDonald R, Pereira F, Ribarov K, et al. Non-projective Dependency Parsing Using Spanning Tree Algorithms[C]// Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing. 2005: 523-530
- [5] Chu Y-J, Liu T-H. On the Shortest Arborescence of a Directed Graph[J]. Science Sinica, 1965, 14: 1396-1400
- [6] Jack E. Optimum Branchings[J]. Journal of Research of the National Bureau of Standards, 1967, 71B: 233-240
- [7] Crammer K, Singer Y. Ultraconservative Online Algorithms for Multiclass Problems[J]. Journal of Machine Learning Research, 2003, 3: 951-991
- [8] Crammer K, Dekel O, Shalev-Shwartz S, et al. Online Passive-Aggressive Algorithms[C]// Proceedings of the Sixteenth Annual Conference on Neural Information Processing Systems (NIPS). 2003
- [9] Meľuk I A. Levels of Dependency in Linguistic Description: Concepts and Problems[J]. Dependency and Valency, 2003(1): 188-230
- [10] Hudson R. An Introduction to Word Grammar[M]. Cambridge: Cambridge University Press, 2010
- [11] Nivre J. Dependency Grammar and Dependency Parsing [R]. School of Mathematics and Systems Engineering, Växjö University, 2005
- [12] 清格尔泰. 蒙古语语法[M]. 呼和浩特: 内蒙古人民出版社, 1992
- [13] Georgiadis L. Arborescence Optimization Problems Solvable by Edmonds' Algorithm[J]. Theor. Comput. Sci., 2003, 301(1-3): 427-437