

文章编号: 1003-0077 (2011) 00-0000-00

基于多策略的藏语语义角色标注研究*

龙从军¹, 康才峻¹, 李琳², 江荻¹

(1.中国社科院民族所, 北京市 100081; 2.青海师范大学计算机学院, 青海省西宁市 810004)

摘要: 语义角色标注研究对自然语言处理具有十分重要的意义。英汉语语义角色标注研究已经获得了很多成果。然而藏语语义角色标注研究不管是资源建设, 还是语义角色标注的技术探讨都鲜有报道。藏语具有比较丰富的句法标记, 它们把一个句子天然地分割成功能不同的语义组块, 而这些语义组块与语义角色之间存在一定的对应关系。根据这个特点, 本文提出规则和统计相结合的、基于语义组块的语义角色标注策略。为了实现语义角色标注, 文中首先对藏语语义角色进行分类, 得到语义角色标注的分类体系; 然后讨论标注规则的获得情况, 包括手工编制初始规则集和采用错误驱动学习方法获得扩充规则集; 统计技术上, 选用了条件随机场模型, 并添加了有效的语言特征, 最终语义角色标注的结果准确率、召回率和 F 值分别达到 82.78%、85.71%和 83.91%。

关键词: 藏语; 语义角色标注; TBL; CRFs

中图分类号: TP391

文献标识码: A

Multi-Strategic Research on Semantic Role Labeling of Tibetan

LONG Congjun¹, KANG Caijun¹, Li Lin², JIANG Di¹

(Institute of Ethnology & Anthropology Chinese Academy of Social Sciences, Beijing 100081, China; The Computer College of Qinghai Normal University, Xining, Qinghai 810004, China)

Abstract: To study Semantic role labeling is of great significance for natural language processing. The researches of semantic role labeling about English and Chinese have obtained many achievements. However, the resources construction and technological means of semantic role labeling in Tibetan are still in initial stage. Tibetan has rich syntactic markers which naturally segment a sentence to different semantic chunking, and there are certain relationship between these semantic chunking and semantic roles. According to this characteristic, the authors of this paper propose the semantic role labeling strategy based on semantic chunking by combining two means of rules and statistics. In order to realize the semantic role labeling, the authors design classification system of Tibetan semantic roles and then discuss the acquisition of rules, including a manual initial rule sets and expanded rule sets from Transformation-Based Error-driven Learning (TBL). The authors adopt Conditional Random Fields (CRFs) Model. In processing of experiment, some effective language characteristics are added, and the results of semantic role labeling are 83.91% for precision, 82.78% for call rate and 85.71% for F values.

Key words: Tibetan; Semantic Role Labeling; TBL; CRFs

1 引言

自动语义角色标注 (Semantic role labeling, 缩写为 SRL) 是自然语言处理的重要任务, 对提高语言信息处理系统的性能具有重要的意义。语义角色标注的过程可以表述为: 设立一套标签体系 (角色分类体系), 部分地标注句子的成分结构 (能承载语义角色), 使计算机自动的获得一定的“理解”能力。

最早研究 SRL 的是 Gildea 和 Jurafsky, 他们开发了一套 SRL 系统, 经对不同的两套语料测试, 实验结果准确率分别约为 82%和 65%^[1]。在 CoNLL2004 会议中, 他们提交的论文强调对句法组块进行分类, 在训练语料相同的情况下, 比较了词到词与短语到短语的标注结果,

* 收稿日期:

定稿日期:

基金项目: 国家自然科学基金重点项目(61132009)

标记	语义角色	标记	语义角色
-ས/གིས/ཀྱིས/ཀྱིས/ཡིས	施事(AG)	དང	涉事(CE)
ལ/-ར	领事(GE)	ལ/-ར/སྟ/དུ/ཅ/ཏུ	结果(OE)
-ས/གིས/ཀྱིས/ཀྱིས/ཡིས	使役(CA)	-ས/གིས/ཀྱིས/ཀྱིས/ཡིས/ནས	方式(MR)
ལ/-ར/སྟ/དུ/ཅ/ཏུ	受事(PT)	-ས/གིས/ཀྱིས/ཀྱིས/ཡིས	工具(IT)
ལ/-ར/སྟ/དུ/ཅ/ཏུ	受役(BE)	-ས/གིས/ཀྱིས/ཀྱིས/ཡིས	材料(ML)
ནས/ལས	起源(SE)	ལ/-ར/སྟ/དུ/ཅ/ཏུ	时间(TM)
ལ/-ར/སྟ/དུ/ཅ/ཏུ	对象(TT)	ལ/-ར/སྟ/དུ/ཅ/ཏུ	处所(LC)
ནས/ལས/བས	依据(BS)	ལ/-ར/སྟ/དུ/ཅ/ཏུ	方位(DN)

(3) 角色配对原则。句法标记并不能与语义角色构成一一对应关系，仍然存在部分无标记的语义组块。针对这种情况，本文作者在对语义角色分类时充分考虑了无标记与有标记的语义角色块在一个句子中的配对关系，如施事与受事，领事与属事，系事与类事，使役与受役之间存在配对关系；特殊句型与语义角色的关系，如领事、属事与领有句有关，系事与类事与判断句有关，使役与受役与使动句有关等。

综合以上的各种因素，本文作者最终为藏语设计了 22 个语义角色类型，具体如表 2 所示。

表 2 藏语语义角色分类体系

名称	标记	名称	标记	名称	标记	名称	标记
施事	AG	属事	BL	起源	SE	方式	MR
受事	PT	系事	CP	对象	TT	工具	IT
当事	EX	类事	IS	依据	BS	材料	ML
客事	BO	使役	CA	涉事	CE	时间	TM
领事	GE	受役	BE	结果	OE	方位	DN
处所	LC	目的	PU				

3 语义角色标注规则构建

与统计方法相比，规则方法在自然语言处理中并无明显优势，这对于资源丰富、数据获取便利的大语种来说更是如此。但是对于资源少、句法标记较丰富的藏语来说，在现阶段也不失为一种有益的尝试。为此，本文采用手工编制初始规则集和利用基于转换的错误驱动学习算法(Transformation-Based Error-driven Learning, TBL)对规则库进行泛化，从而获得扩充规则集。

3.1 初始规则集

初始规则包括语义块边界规则和语义角色与格标记及助词的对应规则。规则的获得主要由人工总结归纳。语义块边界规则由左边界、右边界、双边界和左右边界例外四个部分组成。其中左右边界例外是一个调节规则，就是对左右及双边界标注结果进行纠错。四个部分规则共有 271 个，右边界特征 114 个，双边界特征 119 个，左边界特征 15 个（包括全部动词和否定副词等特征词），左右边界例外特征 35 个。语义角色与格标记及助词的对应规则 63 个。

3.2 扩充规则集

在初始规则集的基础上，本文作者采用 TBL 算法自动从语料中学习并建立扩充规则集。TBL 算法利用学习器从语料中自动获取转换规则集，因此建立一个高效的学习器是 TBL 算法的关键。学习所需资源主要包括以下三方面：(1) 正确标注语义角色的语料，(2) 经初始标注的语义角色语料，(3) 转换规则模板集合。通过比较资源 (1) 和资源 (2) 之间的标注差异，习得扩充规则集。

4 统计模型及特征选择

4.1 条件随机场模型

条件随机场 (Conditional Random Fields, CRFs) 是一种判别式概率模型, 多用于标注或者分析序列材料。在基于统计的标注方法中, 条件随机场模型具有很好的效果, 其模型思想主要来源于最大熵模型, 但又不存在最大熵模型的数据稀疏问题; 同时也无需对数据进行不必要的独立性假设, 在这个方面也优于隐马尔科夫模型 (HMM)。CRFs 通常采用如图 1 的一阶链式结构。

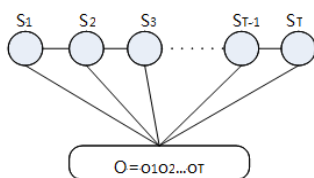


图 1 CRFs 链式图

4.2 标签设计

语义角色标注的标签设计十分关键, 在训练语料不丰富的前提下, 标签的数量会直接影响标注的效果。本文的语义角色标注所使用的语料经人工标注, 标注语料能够提供词性信息、语义组块边界信息和语义角色标记信息。如果把这些信息结合起来考虑, 设计一套联合标签, 在一定程度上, 可能会提高标注效果。边界标记使用了 BIO 标记法, 单词成块或者块外标点为 S, 块开始用 B, 块结尾用 E, 块中间用 I。

表 3 语义角色标注联合标签

词形	ཁོ	མ	ང	འི	གླེང་ལོག་པོ	གསལ་བ	ང	གསལ་བཟུང་	འི་ལོག་པོ	།
词性	rh	wa	rh	wg	ng	a	wd	v	t	xp
边界标记	S	S	B	I	I	E	S	B	E	S
角色标记	AG	M	TT	TT	TT	TT	M	P	P	OT

语义角色以角色类型作为标签。联合标签还包括词形和词性, 具体情况如表 3 所示。

4.3 特征选择

条件随机场算法中, 数据的很多重要信息是通过特征函数给出的。CRF 工具以特征模板的方式, 给出这些特征函数的定义, 使选取特征和定义特征函数非常方便。本文所使用的特征包括词形信息、词性信息、语义组块音节长度信息、谓语动词类型信息以及谓语动词与语义组块之间的距离信息。这些信息的具体描述如下:

词形: 指词的形式属性, 如 ཁོ (khong, 他) 与 རྒྱལ་རང་ (khyed rang, 您) 是两个不同的词形。

词性: 指词的词类信息; ཁོ (khong, 他)、རྒྱལ་རང་ (khyed rang, 您) 是人称代词 (rh), 而 གནས་ཚུལ་ (gnas tshul, 情况) 是普通名词 (ng)。

音节长度: 是指一个语义组块的音节数量, 语义组块的平均音节数量可以作为语义组块边界识别的参考。

[གནས་ཚུལ་/ngམི་འདྲ་བ་/iaའི་/wgའོག་/nd] [གོ་བ་/ngམི་འདྲ་བ་/ia][ཡོད་/ve][།/xp] (gnas tshul mi vdra bavi vog go ba mi vdra ba yod, 不同的情况下有不同的理解。)

这个句子中有四个块, 其中[ཡོད་/ve] (yod, 有) 是谓语组块, [།/xp] 是标点符号, 不需要考虑, 其余两个语义组块各自的音节长度分别是 7 个和 4 个。

谓语动词类型: 是指谓语动词的语义类型, 本项研究按照动词必有论元的数量对动词分类, 分为一元动词、二元动词和三元动词。动词语义类型可以影响动词携带语义角色的数量。

谓语动词与语义块的距离: 是指谓语动词与承载语义角色块之间间隔的音节数, 一般来说受事语义角色与谓语动词近, 施事语义角色离动词远, 这些特征可能有效地辅助推断语

语义组块[ཁྱེད་རང་/rhལི་/wgམིང་/ng] (khyed rang gi ming, 你的名字) 和[ང་/rh] (nga, 我) 边界识别错误, 导致{PT}语义角色标注失败。

(2) 边界识别正确, 语义角色标注错误。语义角色标注的错误表现有: 语义角色未标注, 语义角色标注的位置错误和语义角色选择错误。如例句 4、5。

例 4: [གནས་ཚུལ་/ngཚང་མ་/a]{PT}[ཁོང་/rh]{TT}[ལ་/wd][ཤོད་/vo][དང་/y][ཡི་/xp] (gnas tshul tshang ma khong la shod dang. 全部的情况对他说吧。) (参考答案)

[གནས་ཚུལ་/ngཚང་མ་/a][ཁོང་/rh]{TT}[ལ་/wd][ཤོད་/vnདང་/y][ཡི་/xp] (语义角色标注结果)

组块[གནས་ཚུལ་/ngཚང་མ་/a] (gnas tshol tshang ma, 全部情况) 边界识别正确, 但是{PT}语义角色标注失败, 属于语义角色未标注。

例 5: [ཉིན་མོ་/nt][མི་བཞིན་/d][ལྗང་/a]{OE}[དུ་/ub][འགྲོ་/vo][གི་/t][ཡི་/xp] (nyin mo rim bzhin thung du vgro gi. 白天逐渐变短了。) (参考答案)

[ཉིན་མོ་/nt][མི་བཞིན་/d]{EX}[ལྗང་/a]{OE}[དུ་/ub][འགྲོ་/vo][གི་/t][ཡི་/xp] (语义角色标注结果)

语义组块[ཉིན་མོ་/nt]与[མི་བཞིན་/d]边界识别正确, 但是{AG}语义角色标注错误, 其中[མི་བཞིན་/d]不能承载语义角色, 属于非语义角色组块, 本例属于语义角色标注位置错误。

5.3 实验结果改进

利用已经建立的边界规则库和语义角色标注规则来优化统计标注的结果, 在训练语料规模有限的情况下可能会产生一定的效果。因此我们利用手工编制初始规则集和利用 TBL 方法获得的扩充规则集对统计结果进行校正。在进行第二次实验时, 本文选择了统计标注最好结果作为规则校正的对象。图 3 表示 baseline 统计方法标注、加入语义组块音节数统计方法标注以及统计和规则相结合的标注三种实验结果的对比。

从图 3 可以看到, 利用规则方法对统计结果进行校正, 与单纯依靠统计方法相比, 实验结果有大幅度地提升, 准确率、召回率和 F 值分别达到了 82.78%、85.71% 和 83.91%, 可见规则方法的调节效果还是比较明显的。

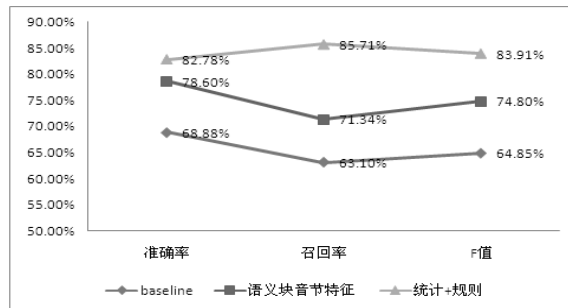


图 3 三种方法标注结果对比图

6 结束

语义角色标注研究属于浅层句法分析的主要内容之一, 在当前完全句法分析存在诸多困难的情况下, 开展浅层句法分析可以有效地提高机器分析与理解自然语言的能力。语义角色标注研究的成果在机器翻译、自动问答、信息抽取等诸多领域都可以得到广泛使用。本文探讨了藏语语义角色标注研究, 通过利用统计和规则相融合的策略, 提升了语义角色标注的效果, 实验结果准确率达到 82.78%。但是本项研究中, 对于嵌套语义组块和长距离语义组块的标注效果并不理想, 这类错误拉低了标注的准确率。在后续研究中, 除了扩充大规模的训练语料和精细化规则集之外, 还需要对嵌套语义组块和长距离语义组块进行专门的纠错处理。

参考文献

[1] Daniel Gildea, Daniel Jurafsky: Automatic Labeling of Semantic Roles, Computational Linguistics, Volume 28, Number 3, p1-45.

- [2] Kadri Hacioglu, Sameer Pradhan, Wayne Ward, James H. Martin and Daniel Jurafsky, Semantic Role Labeling by Tagging Syntactic Chunks. In Proceedings of ConNLL-2004.
- [3] <http://www.cs.jhu.edu/EMNLP-CoNLL-2007/>.
- [4] <http://www.clips.ua.ac.be/conll2008/>.
- [5] Ting Liu, Wanxiang Che, Sheng Li, Yuxuan Hu and Huaijun Liu, Semantic Role Labeling System using Maximum Entropy Classifier, In Proceedings of ConNLL-2005.
- [6] Yu Jiande, Fan Xiaozhong et. Semantic role labeling based on conditional random fields, Journal of Southeast University. Vol 23, No. 3, p 361- 364.
- [7] 王步康、王红玲等：基于依存句法分析的中文语义角色标注[J]，中文信息学报，2010年1月。
- [8] 刘挺、车万翔等：基于最大熵分类器的语义角色标注[J]，软件学报，2007年3月，Vol18, No.3.
- [9] 丁伟伟、常宝宝：基于语义组块分析的汉语语义角色标注，中文信息学报，2009年9月。
- [10] 江荻：现代藏语的句法组块与形式标记[A]. 语言计算与基于内容的文本处理, 孙茂松, 陈群秀主编. 北京: 清华大学出版社. 2003: p160-166.
- [11] 江荻：面向机器处理的现代藏语语法规则和词类、组块标注集[A]，江荻、孔江平主编，中国民族语言工程研究新进展，北京：社会科学文献出版社，2005，p13-106.
- [12] 李琳、龙从军、江荻. 藏语句法功能组块的边界识别[J]. 中文信息学报, 2013年第6期.
- [13] 龙从军、江荻. 现代藏语带助动词谓语句法的识别方法[A]. 第2届青年计算语言学会议论文[C]. 2004.
- [14] 袁毓林：语义角色的精细等级及其在信息处理中的应用[J]，中文信息学报，2007年7月，P10-20.
- [15] 周强、詹卫东、任海波：构建大规模的汉语语块库[A]，清华大学出版社：自然语言理解与机器翻译，2001，pp102-107.
- [16] 周强、孙茂松：汉语句子的组块分析体系[J]，计算机学报，1999, 22(11), pp1158-1165.
- [17] 杨敏、常宝宝：基于北京大学中文网库的语义角色分类[J]，中文信息学报 2011年3月，vol.25, No.2.
- [18] <http://www.keenage.com/>.
- [19] 鲁川：汉语语法的意合网络，商务印书馆，2001年，P111.
- [20] 林杏光：词汇语义和计算语言学，语文出版社，1999年5月，P184.
- [21] <http://CRFspp.googlecode.com/svn/trunk/doc/index.html>.

作者简介：

- 龙从军（1978——），男，博士，主要研究领域：藏语语法、藏语信息处理。Email: longcj@cass.org.cn.
- 康才峻（1980——）男，博士，主要研究领域：藏语信息处理。
- 李琳（1980——），女，博士，主要研究领域：藏语信息处理。
- 江荻（1954——）男，研究员，主要研究领域：藏语语法、藏语信息处理。