

面向信息处理的蒙古语代词语义的分类初探

白玲迪

(内蒙古呼和浩特市 010021)

摘要:蒙古文信息处理的很多工作迫切需要蒙古语词语语义属性的描述体系^②。蒙古语代词语义的形式化描述是蒙古文信息处理语义处理中的一项基础工程^③。本文论述了面向信息处理的代词语义分类的迫切性即代词语义分类的目的、原则、目标及具体语义分类体系与标记集。

关键字:信息处理;蒙古文;代词;语义分类

1. 引言

蒙古语信息处理研究工作从上世纪八十年代着手建立语料库开始,基本完成了了字、词处理阶段的工作,目前已全面开展句子、篇章处理阶段并取得了一定成效。然而,语义分析是句子处理阶段的重要任务之一。蒙古语信息处理的实际需求,词性语义分类目前已完成动词、名词、形容词并都付诸使用。在标注蒙古语语料库语义信息时,代词相关信息极少,遇到很多麻烦,因此完善语料库语义标注中代词语义研究不可或缺。代词是指代事物、现象,具有抽象概括性。它指代的事物、现象体现在实际上下文指定的语言单位。代词的指代范围宽,只有从语境中才能找到它的实际意义。蒙古语代词具有语法形式格和数的形态变化,且使用率高。从动词谓语句简单句的动词(含一价、二价、三价动词)与名词语义搭配研究来看,对句子进行语义加工时,对动词标注语义分类和价数,对名词标注语义分类和介量量(语义角色),而对代词只标了语义角色标注,所以不确定代词在本句中具体代替的内容。为此,着眼蒙古文信息处理的实际需求,迫切需要代词语义进行分类。例如:

① 1e3 TERE/Rj{UJ} SERI/Ve2{VSIHOVY,1}+L_E/Fs31 ./Wp1
含代词句句!

② 1e3 NARA/Ne2{Nbbu11,UJ}
MANDV/Ve2{VHOSIMA,1}+L_A/Fs31 ./Wp1 含名词句句!

2 代词语义分类的目的与原则

2.1 目的

(1)从语言学方面看,语义分类是有利于描述代词的共性和特性。词汇的搭配是语言词汇结合而成的现象,与语义有紧密联系。

(2)蒙古语代词与汉语代词相比语法变化较多。比如汉语代词中“他、她、它”等三种指定代词,但蒙古语中只有一词“@k-r”来表示上面汉语的三种形式,没有阴阳性、人和物的区别,所指的语言单位只能在它的语境中查明。语

① 白玲迪,蒙古族,毕业于内蒙古大学,硕士研究生。

② 海银花《面向信息处理的蒙古语名词语义研究》[D],内蒙古大学,博士论文,2010年。

③ 额尔敦朝鲁《面向信息处理的蒙古语动词语义研究》[D],内蒙古大学,博士论文,2005年。

1

料库出现代词,必须分析其上下文,并标注语义分类,最后归纳代词语义搭配的特征。

(3)代词在词汇体系词数虽少,但使用频率高且复杂。它包含实体(有无生命体)和抽象物(时间、地点、状态、数、范围、性质)等,使机器识别代词时必须要有指定的符号表示语义分类。

(4)通过语义分类能明确知道代词所指代的事物并理解句子语义关系。代词本身没有实际意义,必须从它所在的语境中了解代词的实际意义。

(5)代词语义分类标记有利于完善动词与名词语义搭配研究。

① 1e5 BI/Rb {UD} CU/Sq YEHE/Ac

BAYARLA/Ve2{VSEVRBA,1}+L_A/Fs31 ./Wp1



② 1e5 BI/Rb{Rbe11}{UD} CU/Sq YEHE/Ac

BAYARLA/Ve2{VSEVRBA,1}+L_A/Fs31 ./Wp1

2.2 原则

本文基于语料库代词指代的语境意义进行语义分类。以下为代词语义分类的原则:

(1)具有层次性。代词语义类分为大语义类,子语义类等清晰的层次性。原因在于“大语义类”是针对刻画意义相同或相似的一些代词的共性,而“子语义类(包括中类,小类)”是为了描述属于同一“大语义类”代词之间的差异。

(2)有逻辑性。在同一层次基于统一标准分类。

(3)语义分类既不能太过复杂也不能太过简单。如果过于复杂,语义项多导致研究复杂。太过简单,语义项又不能明确描述代词之间细微的差异。

在上述原则基础上进行代词语义分类并系统地划分代词语义,此次分类与以往分类相比优点在于,首先,分类对象的规模较大从而其结果产生的分类体系具有多层次、多类型、多关系的特征,将会提供更详细、更深入的语义信息。其次,由于分类的深度和广度取决于语法分析的需要。应用语义知识解决那些只靠语法知识难以解决的问题,因此为代词提供语法和语义相结合的全面的语言知识^④。

3 代词语义分类目标

(1)将代词语义分类作为代词属性描述基础,同时语义分类的最终目标是完善蒙古语语义知识库。代词语义是整个蒙古语语言知识库的组成部分,为机器翻译、文本校对、文音转换、文字识别、内容提取、编码转换等应用开发提

④ 海银花那顺乌日图《面向“蒙古语语义信息词典”的名词语义分类体系》[J], <http://www.cnki.net>.

2

供语法和语义相结合的、全体的语言知识^⑤。

(2)代词语义分类不仅揭示每一个代词的意义,而更为重要的是需要对各类代词进行语义搭配研究,对语义搭配信息逐一进行描述的系统的语义分类和语义标记集。

(3)代词语义分类及确定标记集为研究代词的价量、价质的确定,价量、价质之间对应关系的判明等代词配价提供有效知识并在代词语义分类形式化基础上为代词与其它词性词的语义搭配生成由此导出形式规则有着重要意义。

4 代词语义分类

代词语义分类是代词基本词汇意义为依据兼语法信息及结合语料库代词使用情况为基础来划分的。本次语义分类的重点是在于,因为蒙古语代词有格和数的形态变化,所以要结合这些代词特征来决定语义分类的细略略度。由于代词应用广泛复杂,同一个代词可能有不同的归类,而产生不同的语义分类体系。

划分的代词语义类为人称代词、反身代词、指示代词、疑问代词、不定代词、范围代词、代动词等七类。以下是七类

分别划分的子类。子类含中类和小类等。代词语义总体分类为四层,第一层为7个大语义类、第二层为31个中类,第三、四层为30个小类。例如:

人称代词:基于数的变化分为第一,第二,第三等三类,再分中类为单数,复数等子类。反身代词:基于数的变化分为单数、复数等小类。

面向信息处理的代词语义标记是相当于代词语义形式化,用相应的符号来表示语义,语义通过符号标记才能使机器识别。设定代词语义标记时我们参考了国家标准《信息处理技术——面向信息处理的蒙古语词性分类与标记集》与动词、名词,形容词的语义标记体系。代词符号标记第一层采用了本层指代名的首音节,第二层分别在上层上同样加本层指代名的首音节,第三、第四层分别在上层上加阿拉伯数字来表示。所有符号标记前加在《信息处理技术——面向信息处理的蒙古语词性分类与标记集》里里设定的《R》字母。最后标记长度为第一层3位数,第二层是5位数,第三层是6位数,第四层是7位数。

代词语义分类体系及标记集:

5. 结语

面向信息处理的代词语义分类是研究蒙古语信息处理语义发展的基本需求。

我们对于蒙古语代词语义进行具体分类时,基于前人研究成果,运用计算语义学的理论与方法为依据初步完成了此项工作。由于时间缘故和研究尚未完全成熟,还有待完善。下一步在代词与其他词性的搭配规律、代词配价研究中将语义分类体系不断充实、调整完善。

参考文献:

- [1]《现代蒙古语》,内蒙古人民出版社,2005年年。
- [2]秦小锋《语义形式化研究的利与弊初探》[J],语言文字探索,2009年年,第5期。
- [3]林林杏光《词汇语义和计算语言学》[M],语文出版社,1999年年。
- [4]《现代蒙古语100万词级语料库》,内蒙古大学蒙古学学院。