

汉语近义词辨析知识库构建研究

李娟

北京大学中国语言文学系, 北京 100871; E-mail: lijuan0_0@pku.edu.cn

摘要 现有近义词辨析词典往往忽略近义词的细微用法, 缺乏对近义词的句法分布、语义特征、组合限制及使用模板的细致描述。针对该问题, 基于真实语料, 以搭配为载体, 以词与词之间的关系为手段, 构建了一个近义词辨析知识库, 为近义词提供用法和语境知识, 包括近义词的搭配词及搭配关系、搭配频率以及近义词在语料库中的句法分布信息、语义特征、语境中的例句等。该知识库用可视化方式呈现近义词的搭配词和近义词的句法、语义和用法知识, 表示更加直观, 对汉语二语学习者更友好。

关键词 近义词辨析; 中文国际教育; 词汇知识库

A study on the Construction of Chinese Near-Synonyms Knowledge Base

LI Juan

Department of Chinese Language and Literature, Peking University, Beijing 100871; E-mail: lijuan0_0@pku.edu.cn

Abstract Near-synonyms discrimination dictionaries often overlook the nuances of near-synonyms, and lack of detailed description of the knowledge of syntactic distribution, semantic characteristics, combination restrictions and usage templates of near-synonyms. The author constructs a knowledge base of near-synonyms based on corpus, which takes collocation as the carrier and the relationship between words as the means to provide rich usage and context knowledge for near-synonyms discrimination. The proposed knowledge base contains usage information: collocation, collocational relationship, collocational frequency, distributional information, semantic features, example sentences, etc. The knowledge of syntax, semantics, and usage of near-synonyms are presented in a visual way, so it is more intuitive and more friendly to Chinese second language learners.

Key words near-synonym discrimination; Chinese international education; vocabulary knowledge base

近义词的使用和辨析是对外汉语教学中的重点和难点, James^[1]将近义词的使用错误分为两类: 意义关系混淆和搭配错误^[2]。张博^[2]指出: 从理论上说, 这两类错误的区别是: 1) “意义关系混淆”是词义不合, 但搭配无误, 如“大约一分钟后, 十岁的小明拿着行李兴高采烈地跳下楼来(应当用‘跑’, ‘跳下楼’是正确的搭配, 但是在此句中不合语义)”; 2) “搭配错误”是搭配不当, 但在意义上没有问题。如“我眼中的太阳光仿佛一万支箭射过来, 使我张不开眼睛来(应当用“睁”, “张眼睛”不能搭配)。但是, 在中文国际教育中, 汉语二语学习者使用近义词时, 意义使用错误和词语的搭配错误经常混合存在, 单纯

区分意义关系或纠正搭配错误并不能够完全解决近义词辨析的问题。

邢红兵^[3]指出, 第二语言词汇习得实际上是一个从意义到用法的实现过程, 其难点在于意义到用法的转变。他提出“搭配知识”的概念, 希望将词语之间的搭配关系视为词汇知识的重要组成部分, 并利用词汇搭配知识体系进行第二语言词汇习得研究。Sinclair^[4]认为在外语词汇学习过程中, 目的语中常用词的主要用法模式以及典型搭配的学习最重要。

对于二语学习者来说, 搭配知识不但有助于提高语言表达的流利程度, 还能帮助学习者克服母语

负迁移,输出更地道的目的语。是否掌握足够的搭配知识,是学习者能否熟练运用二语的关键。对外汉语的教学实践也表明,外国人学习汉语,往往难于掌握汉语词语的搭配规律^[5]。

除了搭配知识,如何获取以及获取哪些近义词的用法知识,也是解决近义词辨析难题的关键。冯志伟^[6]指出,除词典和语法书中的近义词用法知识外,更多的语言学知识隐藏在语料库中,语言学家对局部语言现象归纳出来的语言学知识难免有片面或错误之处,语料库是语言学知识最可靠的来源。

本文基于真实语料(语料来源于CCL语料库和BCC语料库),以近义词的搭配词和搭配关系以及近义词的用法模板为载体,使用依存句法分析工具辅助抽取语料库中的词语搭配,人工筛选高频的典型搭配,归纳近义词的用法信息,希望提供近义词的以下用法信息:1)常用或特定的词汇搭配格式;2)近义词的分布特征,词在组合中充当的句法成分;3)近义词及其搭配词的语义类,在组合关系中受到的句法语义限制;4)近义词使用时的句法模板。

1 相关研究

《现代汉语语法信息词典》^[7]采用关系数据库文件格式描述词语的语法功能分类和词语的语法属性,分类体系可操作性强,对语法属性的描述非常深入和丰富。词典设置的许多语法特征十分珍贵,是汉语自动句法语义分析的基础,是自然语言处理研究中非常重要的词汇知识库^[8]。但是《现代汉语语法信息词典》是一部机读词典,只描述了词语的组合规则和语义限制,缺少丰富的词语组合示例和词语使用时的用法信息(如缺少词语在使用时的句子模板和丰富例句),不能作为词汇学习工具书直接用于汉语的二语学习。

《现代汉语搭配词典》^[9]和《现代汉语实词搭配词典》^[10]中的词语搭配准确,同时含有拼音、词性、词义、短语结构和搭配类型等信息。但是,这种专家词典缺少例句,只能应用于汉语本体研究或作为机读词典使用,难以应用到汉语的二语学习中。

外向型学习词典是专门为汉语二语学习者编写的词典,收录常见的难以辨析的词语,可从词义解释、用法分析或例句的角度为辨析近义词的差异提供参考。《1700对近义词语用法对比》^[11]是一部提供近义词用法比较信息的词典,共收录1713组近义词(其中80多个是3个词为一组),是目前近义词

辨析词典中收词数量最多的,在词典中提供了近义词的词语搭配实例,用来对比两个近义词能否与同一个搭配词搭配。表1展示《1700对近义词语用法对比》词典中“矮”和“低”的搭配差异。限于词典篇幅,词典提供的搭配数量有限,很难对词的搭配做出全面的明确的解释。

借助语料库来学习和研究搭配可以为学习者提供充分的搭配用法知识和恰当的语境。“中文搭配助手”^[12]是一个在线的词语搭配查询网站,它基于汉语教材语料库抽取教材中的典型搭配,可以自动获取搭配,能够提供词语频率和互信息等计量信息,提供例句帮助理解词语。但是,机器自动抽取的搭配数据中存在一定量的错误搭配。如“去+地”这样错误的搭配,是从“留学生要去绍兴、杭州等地进行语言实践活动。”中抽取的;“地+有”这样错误的搭配,是从“防治‘网瘾’三地各有招数。”中抽取的。在查询量名搭配中,出现“圈+地”的错误搭配。

综上所述,现有的资源直接用于近义词辨析时存在以下不足之处:1)已有的词汇知识库只展示词语和搭配词,不能对比近义词的异同,规模太小,信息有限;2)已有的资源只展示词语和搭配词之间的线性搭配,缺乏可视化词汇搭配网络的展示;3)缺乏系统性的描述框架,没有整体的语言知识描写框架,对词语的特征描写不够系统,只展示词语的搭配词和搭配类型,不标注其他语法信息,没有语义信息和用法信息;4)在线词语搭配网站存在一定的错误搭配。

2 近义词集及语料

2.1 近义词集

《对外汉语常用词语对比例释》^[13]以汉语二语学习中最常用的、意义或用法易混淆的词语为对比对象和选词范围,涉及254组近义词,630多个词。

《1700对近义词语用法对比》选择汉语学习者容易出错的1700对近义词或近义词(包括关系密切

表1 《1700对近义词语用法对比》中“矮”和“低”的词语搭配

Table 1 Collocation between “dwarf” and “low” in 1700 Groups of Frequently Used Chinese Synonyms

词语	~个子	~水平	~空	~年级	~墙	~树	声音~	~头
矮	√	×	×	×	√	√	×	×
低	×	√	√	√	×	×	√	√

的词或词语结构),从语义和语用等方面进行对比分析,每个词条下面分为词义说明、词语搭配和用法对比三部分。

这两本词典选词丰富,涵盖大部分易用错近义词,我们从中筛选 1558 组常用近义词作为搭配库的词语。

2.2 语料来源

本文所需语料应满足:无语病和大规模两个条件,因此选择 CCL 语料库和 BCC 语料库中的文学作品语料作为语料来源。

3 近义词辨析知识库

3.1 句法结构知识

词语搭配和搭配之间的语法关系是近义词辨析时需要的一项非常重要的语言信息,本文通过近义词的句法分布特征、近义词与搭配词的句法关系展示近义词的句法结构知识。关于词语搭配,学界没有统一的定义,Firth^[14]从共现(co-occurrence)的角度探讨搭配问题,认为要通过一个词的共现词来理解这个词的意义。Sinclair^[4]将搭配定义为两个或者两个以上的词在文本中短距离内的共现。Chomsky^[15]认为词语搭配要满足两个选择限制(selection restriction)的条件,即语法规则和词语的语义特征。本研究中的词语搭配指符合一定的语法、语义组合规则且共现频率较高的词语组合。

3.1.1 通过依存句法分析结果获取近义词的搭配词

依存语法(dependency parsing)通过分析词与词之间的依存关系揭示句子的句法结构,认为句子中存在一个核心词,该核心词不受其他词的支配,且能够以某种依存关系支配其他词,处于支配地位的词称为支配者(head),处于被支配地位的成分称为从属者(dependency)。

本文使用 Stanza^[16]获取词与词的语法关系,融合共现频次和互信息等统计特征,自动过滤低频、互信息值低的搭配对,基于语言学知识,经过人工筛选后,可以较快捷地构建词语搭配库。搭配抽取

流程如下:1)用 Stanza 对语料库进行句法分析;2)根据搭配类型抽取共现词对、计算共现频次;3)获取共现词对互信息值;4)根据搭配类型、频次和互信息值筛选共现词对,构建高频搭配词矩阵。

图 1 和表 2 展示来源于 CCL 语料库的例句“现在,保护野生动物已受到世界各国的重视。”的分析结果。

3.1.2 获取共现频率

邢红兵等^[17]认为词语的搭配频率是词汇知识系统的重要特征,是词汇知识存储和加工过程中的重要信息。Wray^[18]认为二语学习者受社会文化以及认知等因素影响,对搭配频率缺乏敏感性,无法像母语者一样熟练地处理高频搭配。但是,目前越来越多的实证研究表明,二语搭配加工存在频率效应^[19],高水平英语学习者运用搭配的速度和准确性与搭配频率成正比,低水平英语学习者主要依赖单个词的频率^[20],受搭配频率的影响不大。

因此,频率信息不仅有助于自动过滤搭配强度低的共现对,对于二语学习教学参考也非常重要。我们将依存句法分析结果中含有近义词的共现词对抽取出来,获取它们的共现频率。例如,从句子“现在,保护野生动物已受到世界各国的重视。”中抽取“受到+保护”、“保护+动物”共现频率各记为 1 次,与其他语料中抽取的相同词语组合累加计算共现频率。

本研究从含有“保护”“爱护”“爱惜”“珍惜”4 个近义词的 2600 句语料中抽取到 3120 条共现对。

3.1.3 获取互信息值

互信息(mutual information, MI)是信息论中的一个信息度量,可以看做一个随机变量中包含的另一个随机变量的信息,反映两个随机变量之间的关系及其强弱。互信息通过计算间隔的词语的关联程度,度量在一个词出现的情况下另一个词出现的概率,从而在一定程度上反映搭配的程度,互信息值越高说明搭配强度越大。两个随机变量的互信息可以定义为

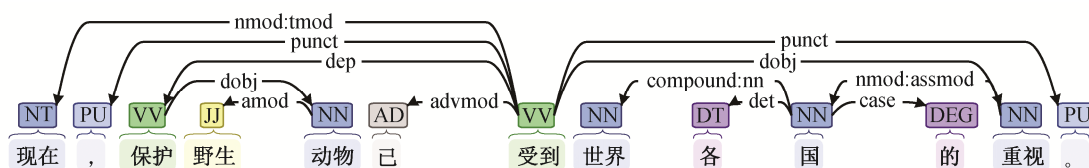


图 1 依存句法分析结果

Fig. 1 Results of dependency parsing

表 2 依存句法分析结果
Table 2 Results of dependency parsing

id	word	head id	head	deprel
1	现在	7	受到	nmod:tmod
2	,	7	受到	punct
3	保护	7	受到	advcl
4	野生	5	动物	amod
5	动物	3	保护	dobj
6	已	7	受到	advmod
7	受到	0	root	root
8	世界	9	各国	nmod
9	各国	11	重视	nmod
10	的	9	各国	case:dec
11	重视	7	受到	obj
12	。	7	受到	punct

说明: id 代表词语在句中的位置, word 是句子中的词, head id 表示与词语依存的词的位置, head 为与该词语具有依存关系的词语, deprel 为两个词语之间的依存关系。例如,“动物”的 head 是“保护”,“保护”指向“动物”,“动物”是“保护”的宾语(dobj),这两个词是动宾关系。

$$I(w_1, w_2) = \log_2 \frac{P(w_1, w_2)}{P(w_1) \times P(w_2)}$$

其中, w_1 和 w_2 表示语料库中的词语, $P(w_1, w_2)$ 表示 w_1 和 w_2 在语料库中同时出现的概率, $P(w_1)$ 和 $P(w_2)$ 分别表示 w_1 和 w_2 单独在语料库中出现的概率。如果 w_1 和 w_2 没有同时在语料库中出现过, 则 $I=0$ 。 I 与 w_1 和 w_2 的相关性成正比, I 越大, w_1 和 w_2 越有可能是合理的搭配。

3.1.4 优化筛选和抽取准确率计算

由于互信息抽取过程中完全不考虑语法语义组合的限制, 因此抽取的词语搭配准确率不高。在抽取共现词对的过程中, 不但要考虑词对出现的概率, 还需要考虑词对是否符合相应的语言学规则。Hunston^[21]提出“MI 值大于等于 3 的搭配可视为显著搭配”, 因此我们只从共现频率和 $MI \geq 3$ 的数据中

抽取符合语言学规则的搭配对进入知识库。

搭配是一种心理组块, 每个人的心理词典不一样。为了处理时更方便, 除依据个人语感外, 本文在人工筛选时保持从宽的筛选标注, 只筛选搭配合适且频率高的共现条目作为搭配条目。

表 3 列出从含有“保护”“爱护”“爱惜”“珍惜”的 2600 条语料中抽取的共现条目及筛选后共现条目的相关数据。

3.1.5 抽取结果分析

如表 4 所示, 本文将抽取的结果按依存关系进行分类, 可以发现: 1) “保护”“爱护”“爱惜”“珍惜”的搭配词主要分布于宾语、主语和状语位置, 但是在补语位置, 这些词的用法差异较大, “爱惜”可以作状态补语, “爱护”没有此用法; 2) 主语和宾语位置的搭配词的语义类和语义特征不同。表 5 记录主语和宾语位置搭配词的实例。

不符合作者语感的共现对如表 6 所示, 本文认为不合理的共现对有以下几类: 1) 搭配限制, <爱护, 个>、<珍惜、自>、<珍惜、交往>、<爱护、对于>、<珍惜、晚上>不能搭配; 2) 搭配接受程度不同, <爱护, 眼珠>在语料中出现 4 次, 但本文认为可能是语料中出现错别字或个别作者的个别用法, 不应该把这个搭配展示给汉语二语学习者; 3) 分词错误, <珍惜、加>、<珍惜、懂>、<珍惜、值>可能是分词错误导致的不合理搭配。

3.2 语义知识

词语在组合时必须要有语义限制, 比如可以说“保护现场”, 但不能说“爱护现场”。不同的近义词在词义组合时有不同的限制, 本文从格语法的角度描写每个论旨角色的语义类, 从语义分类的角度描述名词的上下位语义关系和语义特征^[22]。

3.3 用法模板

本文使用用法模板呈现近义词的用法知识, 比

表 3 近义词的共现数据及抽取准确率
Table 3 Co-occurrence and extraction accuracy of near-synonyms

近义词	句子数/句	筛选前的共现条目/条	筛选后的共现条目/条	人工判定为搭配共现条目/条	筛选后的准确率/%
保护	650	795	64	63	98.44
爱护	650	736	91	83	91.21
爱惜	650	753	71	69	97.18
珍惜	650	836	122	109	89.34
合计	2600	3120	348	324	93.10 (平均值)

说明: 筛选后的共现条目指频率和 MI 值均大于等于 3 的条目, 筛选后的准确率等于人工判定为搭配共现条目除以筛选后的共现条目。

表 4 依存关系类别及搭配实例
Table 4 Dependency relationship and collocation of near-synonyms

deprel	依存关系	共现词数量	搭配
obj	宾语	78	爱护身体、爱惜名誉、珍惜幸福
nusbj	名词主语	43	法律保护、警察保护、大家爱护
advmod	副词修饰	43	很爱护、不爱惜、非常珍惜
aux	助动词	29	能爱护、应当爱惜、要保护
advcl	状语修饰	19	一样爱惜、不知爱惜
det	限定词	11	该珍惜
cop	系动词	7	是珍惜
case:aspect	时态	7	珍惜着、保护了、保护过
conj	并列	4	关怀爱护
cc	中心词+连词	4	爱护和
parataxis	并列关系	3	爱惜爱惜
mark:advb	副词	3	爱惜地
aux:pass	被	2	被珍惜、被保护
mark:comp	补语	1	得珍惜

表 5 近义词在主语和宾语位置的搭配实例
Table 5 Collocation of near-synonyms in the position of subject and object

近义词	搭配词
保护	法律、视力、眼睛、财产、现场、自己
爱护	公物、环境、身体、孩子、荣誉、生命
爱惜	身体、时间、用得、衣裳、羽毛、名声
珍惜	时间、光阴、机会、名誉、爱情、感情

表 6 不符合作者语感的共现对
Table 6 Co-occurrence that are considered ungrammatical

近义词(head)	共现词(word)	依存关系(deprel)	共现频率
爱护	个	nsubj	4
爱护	爱护	xcomp	4
爱护	对于	det	4
爱护	眼珠	obj	4
爱惜	爱惜	advcl	18
爱惜	爱惜	parataxis	4
珍惜	交往	obj	3
珍惜	自	advcl	3
珍惜	加	advcl	3
珍惜	个	obj	3
珍惜	此	obl	3
珍惜	晚上	nmod:tmod	3
珍惜	懂	cop	7
珍惜	值	cop	5

如容易误用的近义词“等待”和“等候”都可以用在“主体_{施事}+V+客体_{目的}”格式中,但它们的用法模板有差异。

“等待”的用法模板:主体_{施事}+等待+客体_{目的}, {主体:[语义类:人|动物|人工物], 客体:[语义类:信息|信息承载物|交通工具]}。

“等候”的用法模板:主体_{施事}+等候+客体_{目的}, {主体:[语义类:人|动物|人工物], 客体:[语义类:人|交通工具]}。

“等待”和“等候”的客体的语义类不同,“等待”的客体可以是信息承载物(如报告、书信和文件),“等候”的客体不能是信息承载物,如“通过学习,一些地区疫情报告人员变被动等待报告卡为主动上门收卡,使病例报告数大为增加。”中的“等待”不能替换为“等候”。

3.4 语法语义知识可视化

3.4.1 语法知识可视化

数据可视化可以为使用者提供交互式的学习体验,基于上面提取的搭配数据,本文使用 R 语言 networkD3 包实现近义词的语法知识搭配网络图。在搭配网络图中,词是节点,语法关系是边。搭配词之间是三元组的关系, <词 1, 语法关系, 词 2>。语法关系指两个搭配词之间的主谓、动宾、定中、状中、述补、联合关系。搭配网络图是 html 格式,可以拖拽,单击某个词语时,可以只显示该词及其相邻词语。

图 2 是“保护”“爱护”“爱惜”“珍惜”这 4 个近义词与其搭配词组成的搭配网络图,可以看出这些词语同时能与哪些词搭配以及分别能与哪些词搭配。比如,“爱护”“爱惜”“珍惜”都能跟“很”搭配,但是“保护”不能与“很”搭配。

3.4.2 语义知识可视化

语义关系指两个词语之间的语义限制,表现为搭配词语限于某种特定的语义类,如在汉语普通话中,“吃”搭配的宾语语义类为食物、“喝”搭配的宾语语义类为可食用的液体。对学习来说,既要知道近义词可以与哪个词语搭配,又要明白为什么应该这样搭配。但是,近义词的语用条件很难描述清楚。获取近义词的搭配信息后,就可以使用关系数据库 Neo4j,通过词与词之间的关系,把词语搭配和语义和用法联结起来,展示词语的用法信息。

图 3 展示“爱惜”和“爱护”与其搭配词的语义类

别,可以看出它们的搭配词的语义区别。同样是动宾结构,“爱护”和“爱惜”的对象不同,“爱护”搭配的对象可以是人(“爱护孩子”),也可以是具体物(“爱护森林”);“爱惜”搭配的对象一般不能是人(一般不说“爱惜孩子”),但可以是具体事物(“爱惜衣裳”),也可以是抽象事物(“爱惜时间”)。

4 结语

本文针对汉语近义词辨析困难的问题,基于语料库,使用人机互助的方式构建汉语近义词辨析知识库。该知识库考虑了词汇搭配层面的知识、语义

层面的词义组合规律、用法层面的句法模板以及词汇组合频率等因素,有助于汉语二语学习者掌握近义词的具体用法。该知识库具有如下特点:1)从语料库中抽取归纳词语的高频典型搭配、分布特征和词语的使用模式;2)描述了近义词与搭配词的语法选择限制关系、近义词所属的语义类以及近义词与近义词之间的语义层级关系;3)提供了词语的句法语义可视化结果;4)知识库描述了近义词实际使用时的句法模板。

本文构建的汉语近义词辨析知识库可以提供近义词词表、词语搭配的结果、频率、近义词的某些

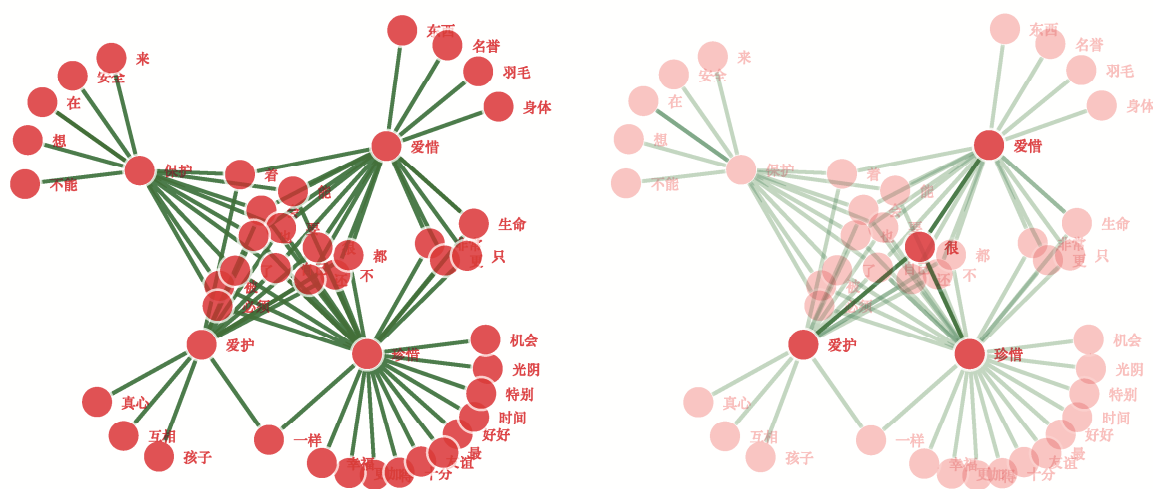


图2 “爱惜”“珍惜”“保护”“爱护”与其搭配词的网络图

Fig. 2 Network of “cherish/treasure”, “cherish”, “protect”, “cherish/treasure/take good care of” and their collocative words

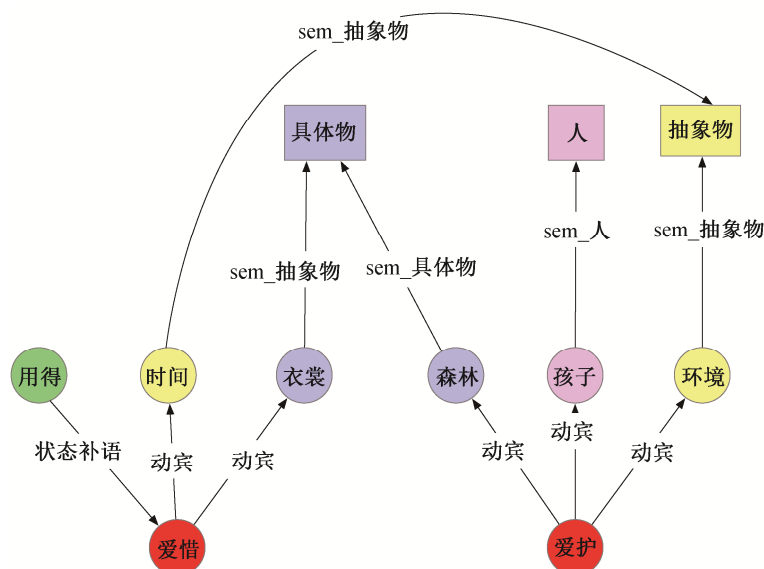


图3 “爱惜”和“爱护”与其搭配词的语义类别

Fig. 3 Semantic categories of “cherish/treasure” and “cherish/treasure/take good care of” and their collocative words

用法差异和语境信息等(目前仅完成部分词语的语义信息和句法模版的标注),但这些功能方面的优势主要体现在实词(如名词、动词和形容词)上,对其他词类(如副词和介词)的描述存在缺陷。

参考文献

- [1] James C. Errors in language learning and use: exploring errors analydid. London: Addison Wesley Longman Limited, 1998
- [2] 张博. 不同母语背景的汉语学习者词语混淆分布特征及其成因研究. 北京: 北京大学出版社, 2016
- [3] 邢红兵. 词语搭配知识与二语词汇习得研究. 语言文字应用, 2013(4): 117-126
- [4] Sinclair J. Corpus, concordance, collocation. Shanghai: Shanghai Foreigner language Education Press, 1991
- [5] 林杏光. 论词语搭配及其研究. 语言教学与研究, 1994(4): 18-25
- [6] 冯志伟. 从语料库中挖掘知识和抽取信息. 外语与外语教学, 2010(4): 1-7
- [7] 俞士汶, 朱学锋, 王惠, 等. 现代汉语语法信息词典详解. 北京: 清华大学出版社, 1998
- [8] 冯志伟, 曹右琦. 评《现代汉语语法信息词典详解》. 中文信息学报, 1999(1): 66
- [9] 梅家驹. 现代汉语搭配词典. 上海: 汉语大词典出版社, 1999
- [10] 张寿康, 林杏光. 现代汉语实词搭配词典. 北京: 商务印书馆, 2002
- [11] 杨寄洲, 贾永芬. 1700对近义词语用法对比. 北京: 北京语言大学出版社, 2005
- [12] 胡韧奋, 肖航. 面向二语教学的汉语搭配知识库构建及其应用研究. 语言文字应用, 2019(1): 135-144
- [13] 卢福波. 对外汉语常用词语对比例释. 北京: 北京语言文化大学出版社, 2000
- [14] Firth J R. Modes of meaning. London: Oxford University Press, 1957
- [15] Chomsky N. Syntactic structures. Hague: Mouton, 1957
- [16] Qi Peng, Zhang Yuhao, Zhang Yuhui, et al. Stanza: a python natural language processing toolkit for many human languages // Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, 2020: 101-108
- [17] 邢红兵, 辛鑫. 第二语言词汇习得的中介语对比分析方法. 华文教学与研究, 2013(2): 64-67
- [18] Wray A. Formulaic language and the lexicon. Cambridge: Cambridge University Press, 2000
- [19] Kim S H, Ji H K. Frequency effects in L2 multiword unit processing: evidence from self-paced reading. TESOL Quarterly, 2012, 46(4): 831-841
- [20] Wolter B, Yamashita J. Word frequency, collocational frequency, L1 congruency, and proficiency in L2 collocational processing. Studies in Second Language Acquisition, 2017, 40(2): 1-22
- [21] Hunston S. Corpora in applied linguistics. Cambridge: Cambridge University Press, 2002
- [22] 梅家驹. 同义词词林. 上海: 上海辞书出版社, 1983