



教育部人文社会科学重点研究基地
北京外国语大学中国外语与教育研究中心
大数据视野下的外语与外语学习研究系列丛书

总主编 ◎ 梁茂成

英语中的生命度 等级研究

吉洁 ◎著

collocation
egocentricity Crown/CLOB
prototypicality dependencies
animate corpora lemma keywords text
CLAWS PowerGREP BNC
dynamicity hierarchy prototype
BNC Stanford Parser AnimacyHabiter
WordNet dynamic mechanism prototype features



教育部人文社会科学重点研究基地
北京外国语大学中国外语与教育研究中心
大数据视野下的外语与外语学习研究系列丛书
总主编 ◎ 梁茂成

英语中的生命度 等级研究

A Study of Animacy Hierarchy in English

吉洁 ◎著

外语教学与研究出版社
FOREIGN LANGUAGE TEACHING AND RESEARCH PRESS
北京 BEIJING

图书在版编目 (CIP) 数据

英语中的生命度等级研究 / 吉洁著. -- 北京 : 外语教学与研究出版社, 2017.11
(大数据视野下的外语与外语学习研究系列丛书 / 梁茂成总主编)
ISBN 978-7-5135-9642-8

I . ①英… II . ①吉… III . ①英语－语法－研究 IV . ①H314

中国版本图书馆 CIP 数据核字 (2017) 第 291544 号

出版人 徐建忠
责任编辑 毕争
责任校对 解碧琰 刘伟
封面设计 彩奇风
出版发行 外语教学与研究出版社
社址 北京市西三环北路 19 号 (100089)
网址 <http://www.fltrp.com>
印刷 北京九州迅驰传媒文化有限公司
开本 650×980 1/16
印张 14
版次 2017 年 11 月第 1 版 2017 年 11 月第 1 次印刷
书号 ISBN 978-7-5135-9642-8
定价 52.90 元

购书咨询: (010) 88819926 电子邮箱: club@fltrp.com

外研书店: <https://waiyants.tmall.com>

凡印刷、装订质量问题, 请联系我社印制部

联系电话: (010) 61207896 电子邮箱: zhijian@fltrp.com

凡侵权、盗版书籍线索, 请联系我社法律事务部

举报电话: (010) 88817519 电子邮箱: banquan@fltrp.com

法律顾问: 立方律师事务所 刘旭东律师

中咨律师事务所 殷斌律师

物料号: 296420001

该书得到以下项目及资金的资助，特此鸣谢。

项目名称：近六十年涉京报道中本土化英语的多维变迁：基于语料库的历时研究（16YYC036）

项目类别：北京市社会科学基金项目

资助单位：北京市哲学社会科学规划办公室

项目名称：基于多模态语料库的英语精读辅助教学资源的建设（3162015ZYKD08）

项目类别：中央高校基本科研业务费专项资金

资助单位：外交学院

资金名称：中央高校基本科项业务费专项资金

资助单位：外交学院

总序

一、引言

科学研究方法大致有二：其一，归纳法。归纳法指根据一类事物的部分对象的属性推知该类事物的所有对象皆具有某种属性。比如，早期的人类在多次与狼邂逅的过程中，逐渐意识到这种体型匀称协调、四肢修长、头腭尖形、鼻端突出、耳尖直立、善于快速奔跑的野生动物具有极强的攻击性，不可为伍，需要敬而远之或群起而杀之。显然，人类是在经历了多次这样的邂逅之后才意识到了狼的危险性，每一次邂逅都为人类积累了经验、加深了印象，终于在总结若干次教训之后形成了结论：所有的狼都是危险的。诚然，人类在形成结论之前不可能邂逅了所有的狼，但照样可以得出正确的结论。其二，演绎法。演绎法指从一般性的 (general) 前提出发，通过推导得出具体的 (specific) 结论。比如，在人们把“所有的狼都是危险的”这一命题视作为一般性前提时，每次邂逅一匹狼，必然会立刻意识到眼前这匹狼是危险的。这其中包含了三个论断，即：所有的狼都是危险的；这是一匹狼；这匹狼是危险的。归纳法是由具体到一般的过程，而演绎法是由一般到具体的过程。

语言研究也不例外，其方法概括起来也不外乎有归纳法和演绎法。演绎法依据可靠的前提进行严密推导，常常可以直击结论。对这种研究方法的运作逻辑我们暂且不做讨论。对于归纳法，其中有若干要素需要考虑。首先，狼有很多特征，哪些特征才具有区别性？哪些属性才是狼的致命属性？比如说，狼嚎是否是我们应该考虑的特征？其次，人类需要与狼邂逅多少次，得出来的结论才是可靠的？返回到语言研究中，前一个问题是我们最为关注的问题。语言分析可以从多种语言特征入手，但哪些语言特征才是最有意义的？我们又该如何选择、提取和分析这些语言特征呢？后一个问题是实证研究中的样本问题，即，我们需要

观察多大的语言样本，才可以得出可靠的结论？

自20世纪后半叶语料库语言学问世以来，研究者越发对自然发生语言数据产生了依赖，因而产生了“经验主义语言学”、“概率语言学”、“数据驱动语言学”等说法，语料库语言学也随之兴起。就其实质而言，语料库语言学采用的是典型的归纳法。语料库是大量自然语言样本的汇集，解决了以上的第二个问题，即实证研究中的样本问题。有了大样本，充分观察成为可能，归纳而得到的结果变得更为可靠甚至可以反复验证。此外，作为方法论的语料库语言学还包含一整套分析方法和分析工具，因而解决了以上第一个问题，即如何提取和分析语言特征的问题。关于选择何种语言特征进行分析，我们将在下面讨论。

总之，有了语料库，我们可能“邂逅”的语言事实更为真实、丰富、全面，这也使得通过归纳法得出的结论更为可靠、经得起验证，不需要像Edward Sapir那样亲力亲为地走入印第安部落之中去采集各式各样的语言数据，也不需要像Charles Fries那样随身携带录音机，甚至不需要像Otto Jespersen那样不失时机地以卡片形式随时记录阅读和日常生活中接触到的各种语言事实。

基于语料库数据进行语言研究，这种方法与演绎法最重要的区别之一在于，研究者在研究中所使用的所有数据均为实际发生过的语言事实，而不是靠想象编造出来的句子：

The rat the cat the dog chased killed ate the malt.

Colorless green ideas sleep furiously.

Sincerity admires John.

Golf admires John.

显然，以依据研究者的直觉编造出来的句子作为研究数据，所得结果需要以语言事实来加以验证。正因为语料库语言学研究中的全部数据皆源于事实，结果也更为可靠，因而受到了越来越多研究者的青睐。在这一理念的主导下，我们近年来进行了若干项研究，目的在于利用语料库和语言大数据，对一些语言理论问题进行深入探讨，并试图解决中国外语教育中的一些现实问题。基于这些研究，我们编辑出版了这一套丛书。

二、语料分析中的语言特征选择

正如狼的所有特征并非同等重要一样，语言特征的选择在语言的量化研究中也至关重要。在前语料库时代，虽有研究者关注语言事实，但

大部分研究者常常根据自己的直觉选择一些特征进行研究。到了语料库时代，特征的选择方法发生了根本性变化。

在语料库时代，人们将语料库中的连续文本制作成词表或多词词表，甚至制作成词类 (POS, part of speech) 列表或词类序列 (POS sequence) 列表，然后对基于不同语料库制作而成的此类列表通过精巧的算法进行频率对比，进而有效地发现语料库中更为有意义的语言特征，特别是词语使用方面的特征。这种方法是语料库语言学研究中常用的主题词分析 (keywords analysis)，研究中几乎总会使用到一个观察语料库和一个参考语料库，并将由这两个语料库析出的词表进行对比，差异较大的词语 (即语言特征) 会自动浮现出来。这种特征选择方法虽有人工参与，但研究者的主观性和偏好得到了有效控制，因而研究结果也更为可靠，研究也可以重复验证。在有些研究中，人们还在两个语料库中查询自己感兴趣的语现象，然后对所得频数进行对比，以发现两语料库间的差异。此外，人们还可以编写复杂的正则表达式，从语料库中提取比词表更复杂的语言特征，如名词短语、介词短语、动宾结构、定中结构、关系从句等，甚至涉及意义单位。

上文中描述的基于语料库的语言研究是当今最为常见的语言研究方法之一，其源头至少可以追溯到 20 世纪八九十年代，也有研究者将此种研究范式视为盛行于 20 世纪 50 年代的美国结构主义的延续和发展，甚至也有研究者将语料库之源头追溯到更为久远的时代。笔者认为，基于语料库的研究最早也只能追溯到电子语料库问世之日。正是随着电子语料库的问世，语言研究所需的研究素材在量 (quantity) 和质 (quality) (即语言的真实性) 两方面才有了真正的突破。基于语料库的语言研究是时代发展的必然，也为语言研究带来了新视野和新维度。在研究过程中，文本的质和量是研究的基础，而文本分析技术和对比算法起到了关键的作用，可以帮助我们发现最有意义的语言特征。

到了当今的大数据时代，情况又有了新的变化。计算机技术的发展推进了网络技术和互联网的普及，而网络的普及就意味着越来越多的人会花费更多的时间浏览越来越多的网页、上传越来越多的内容，发帖、回帖、发表评论，等等，这一切几乎无时无刻不在发生。智能手机的出现和普及更加推进了这一进程，登录网络、发表言论不再受时间和空间的限制。而所有这一切活动中最为常见的媒介正是我们研究的对象——语言。如此发展下去，网络上的语言资源会越来越多，沉淀也会越来越深，长尾效应也越来越明显。在这一背景之下，语言学家自然不应该满

足于原来规模的语料库，他们与计算机领域的专家联手，设计出了各种工具（常称为网络爬虫），可以从网络上获取大量的文本，彻底颠覆了传统语料库的概念。如今，语料库规模已经由原来的百万词级增大到动辄几千万词或数亿词级，甚至达到几十亿或百亿词级。如此规模的语料库，其优势自然毋庸置疑，长尾效应更扩展了研究维度，基于这样的语料库所得到的研究结果也更为可靠、更为多样化，对语言变化的预测能力也更强。然而，在这样的语料库中查询语言特征或由如此规模的语料库生成词语、词类、各类序列或结构列表变得不再那么容易，对这些海量语料库通过主题词分析法进行对比则更加困难。在大数据时代，我们所面临的问题已经不再是语言研究素材的不足。恰恰相反，数据量过于庞大为语言特征的提取带来了新的挑战，原来的文本分析技术和对比算法不再适用。研究者不得不另辟蹊径。

三、大数据时代的语言研究

大数据给语料库语言学者带来了新问题和新挑战。数据量（volume）庞大是大数据时代最为显著的特征，但这并不是大数据的唯一特征。数据传输和变化之快，即大数据的速度（velocity）使得研究所依赖的数据几乎没有确定的形态，时刻处于变化之中，体量也不断增大，这也是我们必须面对的另一问题。除此之外，大数据的庞杂性（variety）也是一个棘手的问题。以上三个V被公认是大数据的典型特征。在大数据时代，语料库的创建、语言分析工具的开发、统计分析方法的更新和完善、统计结果的呈现等多个问题都将面临一场革命性的变化。

在语料库创建方面，巨量语料库的提纯是一个至关重要的问题。由于网络文本的多样性，粗暴而盲目地堆砌文本、追求语料库的大容量，会使得语料库变得十分地异质、庞杂，因而是不可取的。为此，人们汲取了网络爬虫技术，并加以改造，推出了Web as Corpus技术并开发了专用软件，依据网络页面中的关键词快速创建各种专题语料库。这种技术必将成为大数据时代语言研究中的重要技术。另外，专题语料库固然重要，但对于语言研究者而言，语体差异性、文本的时代性等问题也是语言研究中必须考虑的因素。与语体差异性、文本时代性等密切相关的問題之一是，我们应该如何通过各种途径有效获取文本的外部属性（即元信息），这也是大数据时代的语言研究中面临的又一重大挑战。只有挖掘网络文本的元信息特征，研究者才可以利用文本的各种社会属性（如语

种、产生年代、作者身份、作者性别、语体特征、领域特征等），使语言研究特别是文本差异（text variation）研究得以深入。

在语言分析工具方面，由于大量文本都存储于网络或云端，加之语料库规模不断扩大，原先广泛使用的WordSmith Tools、AntConc等单机版的文本分析工具逐渐会变得不再适用，基于网络或云端的工具或许将会成为技术开发的重点之一。此外，在语料库加工方面，基于大数据和深度学习（Deep Learning）技术设计的系统（如谷歌公司开发的句法标注工具SyntaxNet）将代表主流的研究方向，标注的准确率也会有明显提高。

从标注语料库中提取和统计语言特征时，原先广泛使用的统计方法不再适用，主题词分析方法随着语料库规模的增大也必将变得越来越困难，逐渐取而代之的是更为复杂的数据科学（Data Science），聚类、因子分析、复杂回归分析等成为语言分析的常用方法，分析工具也由原来常用的SPSS等工具变成R等更为复杂的系统。R软件的优势不仅在于可以分析大数据，还将编程和统计融合起来，使研究者可以定制各式各样的分析手段。

在统计结果呈现方面，语料库研究常见的图表呈现方式仍然会被广泛使用，但与此同时，随着数据量的增大，数据的可视化将成为呈现研究结果的重要方式，这种呈现方式将更为直观、便于理解。相信在不远的未来，语料库研究的结果将会使越来越多的人受益。

四、结语

随着大数据时代的到来，语料库语言学必将得到更多研究者的重视和青睐，大数据时代的特点将在语言研究中逐渐显现。我们希望通过本系列丛书的出版推进语言研究的不断科学化，推动我国外语与外语教育研究的发展。

本套丛书是教育部人文社会科学重点研究基地北京外国语大学中国外语与教育研究中心“十三五”规划重大项目“大数据视野下的外语与外语学习研究”（编号：17JJD74000）的研究成果，特此鸣谢。

梁茂成
二〇一七年三月

前言

生命度原本作为生物概念，用以区分生命体和非生命体。Jespersen (1922) 最早将生命度与语言研究联系在一起，他认为印欧语系中“性范畴”的起源可能与生命度有关，并进一步探讨了生命度与语序、论元角色和所有格等语言现象之间的关系。之后几十年中，生命度渐渐作为一种语言学范畴得到广泛研究。从语言学层面来看，生命度不再是生物意义上或“有”或“无”的二分，而是具有梯度性的等级。许多语言对生命度等级具有明确的语法规定，语言类型学将这些等级归纳为：“第一、二人称代词>第三人称代词>专有名词>人>动物>无灵名词”。这些等级直接制约着数、一致、格标记、语序、语态、语义角色等多种语言现象，还影响着语言的理解和习得。鉴于生命度的重要作用和广泛影响，语义学、语言类型学、心理语言学、认知语言学和语言习得等多个领域都对生命度进行了大量探讨。

然而，已有研究多是在研究语言变异、语言理解或认知的过程中，间接探讨了生命度。各自带有不同的学科视角，并不针对生命度概念本身，也没有解释生命度等级及其影响的成因。已有生命度等级主要集中于名词和代词中，尤其是有灵词语内部，较少探讨动词等其他词类的生命度，且多来源于类型学的观察举例或认知语言学的内省提炼，缺乏大量语料的验证，等级排序也较为固态，较少考虑语境等变化因素。因此，本研究试图对生命度概念进行深入的理论思考，在其基础上提出生命度等级的分析方法，并使用真实语料对其研究层面进行拓展。

本研究基于英国国家语料库 (British National Corpus) 与 Crown CLOB 语料库，从横向、纵向两个维度拓展了生命度的研究范围。在横向维度中，将英语中的生命度等级从原有的名词和代词扩展到动词和形容词，并从词语层面扩展到了搭配和语篇层面。在纵向维度，将

生命度从静态等级扩展到动态等级。研究结果表明，词语层面存在静态生命度等级，且各词类的等级分布不同：名词个数呈“沙漏型”分布，频数呈“蝴蝶型”分布；动词个数呈“橄榄型”分布，频数呈“陀螺型”分布；形容词的个数与频数均呈“金字塔型”分布。并且不同等级词语所带有的原型特征也不同：词语等级越高，所带特征就越多，特征的权重也越大。而在搭配和语篇层面，既存在静态等级，又存在动态等级。生命度等级之所以会发生动态变化，主要是由于转喻、隐喻、限定及多等级性这四种作用机制。有九种搭配形式经常发生动态变化，它们比其他搭配形式具有更强的动态性。不同语体语篇的生命度等级和动态性也不相同：其中小说语体的生命度最高，动态性最弱；而学术语体的生命度最低，动态性最强。

本书从选题、架构、设计到撰写，始终倾注着我的导师梁茂成教授的心血。还要感谢计算语言学家冯志伟教授、浙江大学的刘海涛教授、中国人民解放军外国语学院的易绵竹教授、北京外国语大学的李文中教授和许家金教授，感谢他们从繁忙的工作中抽出时间审读了本书的初稿，提出了许多宝贵的意见和建议。同时感谢外交学院中央高校基本科研业务费专项资金和北京市社会科学基金项目对本研究的资助。

由于笔者水平有限，书中难免有纰漏之处，恳请各位读者不吝赐教。

目 录

绪论	1
0.1 研究背景	1
0.2 研究意义	3
0.2.1 理论意义	3
0.2.2 方法论意义	4
0.2.3 实践意义	5
0.3 研究概述	6
0.4 本书结构	6

第一章 语言学中的生命度研究综述	8
1.1 什么是生命度	8
1.1.1 生命度是一种语义属性	9
1.1.2 生命度是一种语用知识	9
1.1.3 生命度是一种语法规定	10
1.1.4 生命度是一种认知参数	13
1.1.4.1 共情等级	13
1.1.4.2 认知凸显等级	13
1.1.4.3 施事性等级	14
1.1.4.4 定指性等级	15
1.2 生命度的重要性	16
1.2.1 影响语言多个方面	16
1.2.1.1 横向延拓	17
1.2.1.2 纵向延拓	18

1.2.2 影响多种语言	19
1.2.2.1 生命度硬限制	19
1.2.2.2 生命度软限制	20
1.3 已有研究的总结与问题	22
1.4 本研究的不同之处	25
1.5 小结	27
第二章 对生命度概念的理论思考	28
2.1 对生命度概念的再认识	28
2.1.1 生命度的体验性	28
2.1.2 生命度的原型范畴性	30
2.1.3 生命度的共性与动态性	31
2.2 生命度的原型特征分析框架	33
2.2.1 哲学中关于“人的本质”的阐释	33
2.2.1.1 理性	33
2.2.1.2 感性	34
2.2.1.3 能动性	35
2.2.1.4 社会性	36
2.2.2 心理学中关于“生命识别线索”的研究	37
2.2.2.1 形态与动作线索	37
2.2.2.2 思维线索	39
2.2.3 生物学中对于“生命”的定义	41
2.2.4 生命度原型特征分析框架的构建	42
2.2.4.1 生命度原型特征的确定	42
2.2.4.2 特征之线索有效性的分析	44
2.3 小结	48
第三章 研究方法	49
3.1 研究设计	49
3.1.1 研究问题	49
3.1.2 研究流程与步骤	50
3.1.2.1 词语层面	50
3.1.2.2 搭配和语篇层面	52

3.1.2.3 本章所涉及的研究步骤	53
3.1.3 研究工具	53
3.1.3.1 WordNet	53
3.1.3.2 PowerGREP	54
3.1.3.3 Survey Monkey	55
3.1.3.4 SPSS	55
3.1.3.5 Stanford Parser	56
3.1.3.6 AnimacyHabiter	57
3.1.3.7 AnimacyViewer	57
3.1.3.8 Microsoft Excel	57
3.1.4 语料	57
3.1.4.1 BNC	57
3.1.4.2 Crown/CLOB	58
3.2 名词生命度的计算	59
3.2.1 名词词表的建立	59
3.2.2 名词生命度的问卷评分	62
3.3 动词和形容词生命度的计算	63
3.3.1 常用动词和形容词词表的建立	63
3.3.2 AnimacyHabiter的设计	65
3.3.2.1 总体思路	65
3.3.2.2 前期准备	67
3.3.2.3 操作流程	67
3.4 搭配和语篇生命度的计算	69
3.4.1 生命度分级词表和原型特征词表的建立	69
3.4.2 AnimacyViewer的设计	71
3.4.2.1 静态生命度的计算	71
3.4.2.2 动态生命度的计算	72
3.5 小结	73
第四章 词语的生命度等级	74
4.1 名词的生命度等级	74
4.1.1 基于原型特征分析的名词生命度等级	74
4.1.2 名词生命度等级的验证及确定	78
4.2 动词的生命度等级	80

4.2.1 基于语料概率计算的动词生命度等级	80
4.2.2 高与较高生命度动词的原型特征	83
4.2.2.1 理性	83
4.2.2.2 感性	85
4.2.2.3 能动性	88
4.2.2.4 社会性	90
4.2.2.5 身体动作	93
4.2.2.6 生理过程	96
4.2.2.7 动词的原型特征与其主宾构式	97
4.2.3 中与较低生命度动词的原型特征	99
4.2.3.1 理性	100
4.2.3.2 社会性	102
4.2.3.3 能动性	103
4.2.3.4 身体动作	104
4.2.3.5 中低动词带有高级特征的原因	106
4.3 形容词的生命度等级	107
4.3.1 基于语料概率计算的形容词生命度等级	107
4.3.2 高与较高生命度形容词的原型特征	109
4.3.2.1 理性	110
4.3.2.2 感性	111
4.3.2.3 能动性	111
4.3.2.4 社会性	113
4.3.2.5 身体形态	114
4.3.2.6 生理过程	114
4.3.2.7 形容词的原型特征与其定表用法	115
4.3.3 中与较低生命度形容词的原型特征	115
4.3.3.1 理性	115
4.3.3.2 感性	117
4.3.3.3 能动性	119
4.3.3.4 社会性	120
4.3.3.5 身体形态	122
4.3.3.6 中低形容词带有高级特征的原因	122
4.4 名词、动词、形容词的生命度等级之总结与对比	123
4.5 小结	126

第五章 搭配与语篇的生命度等级	127
5.1 搭配的生命度等级	127
5.1.1 什么是搭配的生命度等级	127
5.1.1.1 静态生命度等级	127
5.1.1.2 动态生命度等级	129
5.1.2 生命度动态等级的作用机制	130
5.1.2.1 转喻	130
5.1.2.2 隐喻	131
5.1.2.3 限定	132
5.1.2.4 多等级性	133
5.1.3 生命度动态等级的常见搭配形式	135
5.2 语篇的生命度等级	138
5.2.1 什么是语篇的生命度等级	139
5.2.1.1 静态生命度等级	139
5.2.1.2 动态生命度等级	139
5.2.2 不同语体语篇的生命度对比	140
5.2.2.1 静态生命度对比	140
5.2.2.2 动态生命度对比	142
5.3 词语、搭配、语篇的生命度等级之联系	144
5.4 小结	145
<hr/>	
第六章 结论	146
6.1 研究发现	146
6.1.1 词语层面的生命度等级	146
6.1.2 搭配和语篇层面的生命度等级	147
6.2 应用展望	148
<hr/>	
参考文献	150
<hr/>	
附录	162

表 目

表 1-1 生命度等级的跨语言语法规规定	12
表 1-2 生命度相关研究总结	23
表 1-3 本研究与已有研究的不同	26
表 2-1 生命度的共性与动态性	32
表 2-2 生物学家所提出的生命特征之前五个	42
表 2-3 哲学、心理学与生物学中关于人、动物和生命的属性汇总	44
表 3-1 WordNet中25类名词的序号、名称及译名	54
表 3-2 Stanford Parser中四种依存关系的输出格式及示例	56
表 3-3 Crown/CLOB语料库的库容信息（单位：词次）	59
表 3-4 AnimacyHabiter所输出的动词生命度结果	66
表 3-5 AnimacyHabiter所输出的形容词生命度结果	67
表 3-6 AnimacyHabiter的检索、判定与统计阶段之间的对应关系	68
表 3-7 AnimacyViewer所输出的语篇静态生命度结果	72
表 3-8 AnimacyViewer所输出的语篇动态生命度结果	73
表 4-1 25类名词的生命度原型特征分析	75
表 4-2 25类名词的问卷平均得分及其内部一致性信度	78
表 4-3 各生命度等级的名词类别及其平均得分	79
表 4-4 各生命度等级的名词个数、频数及用法	80
表 4-5 名词与动词生命度等级的差异	81
表 4-6 由人发出该动作（主语为人）的概率等级	82
表 4-7 使人发出该动作（宾语为人）的概率等级	82
表 4-8 各生命度等级的动词个数及示例	83