



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

中国学生英语限定性 关系从句中关系代词 取舍的多因素分析

房印杰 ◎著

predictability
logistic regression collostruction
variation usage-based preference
pattern prototypicality
WECCL continuum
annotation multifactorial relativizer omission
schematic construction concrete
competition interlanguage

外语教学与研究出版社

FOREIGN LANGUAGE TEACHING AND RESEARCH PRESS



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

中国学生英语限定性 关系从句中关系代词 取舍的多因素分析

A Multifactorial Analysis of Relativizer
Omission by Chinese EFL Learners

房印杰 ◎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

中国学生英语限定性关系从句中关系代词取舍的多因素分析 / 房印杰著. --
北京 : 外语教学与研究出版社, 2017.9

(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 978-7-5135-9511-7

I. ①中… II. ①房… III. ①英语－从句－研究 IV. ①H314.3

中国版本图书馆 CIP 数据核字 (2017) 第 249356 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 15
版次 2017 年 9 月第 1 版 2017 年 9 月第 1 次印刷
书号 ISBN 978-7-5135-9511-7
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: 295110001

该书得到以下项目的资助，特此鸣谢。

项目名称：基于使用的中国学生英语关系从句认知统计模型研究

项目类别：教育部人文社会科学青年基金项目

项目批准号：17YJC740018

资助单位：教育部

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，前一个问题也是语言学家最为关注的问题。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，因而解决了以上的第一个问题，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音现象，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问题之一是，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：17JJD740003)的研究成果，特此鸣谢。

梁茂成
二〇一七年三月

前言

关系从句一直是语言学研究的一个热点和重点。20世纪60年代以来，针对英语母语者的关系从句研究经历了三个维度的转变。在研究路径上，从基于内省式思辨转向基于语言产出；在研究对象上，从关注不同类型关系从句的认知难度差异转向关系代词与从句整体认知加工负载的相关性；在研究方法上，从单因素分析转向多因素分析。不同于针对英语母语者关系从句的研究，国内针对中国英语学习者关系从句的研究在上述三个维度上相对滞后。回顾国内外的关系从句研究，可以发现该类研究存在如下不足：1) 在研究对象上，对比不同类型关系从句的认知难度一直占据主导地位，对关系代词的选择机制关注不足。虽然针对英语母语者的研究开始关注关系代词与从句整体认知加工负载的相关性，但是该类研究仍处于起步阶段；国内对关系代词的研究则更为鲜见。2) 在研究范畴上，对制约关系从句的潜在因素关注不足，往往聚焦少数几个因素，整体涵盖范畴有限。3) 在统计方法上，现有研究囿于描述统计，对多因素统计方法应用不足。4) 对关系从句的认知加工机制讨论不足，未能充分结合心理语言学、认知语言学、计算语言学、语料库语言学的最新研究成果。

鉴于此，本书以中国英语学习者和英语本族语者产出的非主语型限定性关系从句为研究对象，聚焦两个语言群体在关系代词取舍中的异同，尝试挖掘影响关系代词取舍的认知加工机制，特别是制约中国英语学习者关系代词取舍的认知加工机制。整个研究以基于使用的研究范式为理论基础，通过复杂统计模型开展数据分析，着重回答三个研究问题：1) 中国英语学习者与英语本族语者在关系代词取舍使用上是否存在共同的制约因素？2) 与英语本族语者相比，制约中国英语学习者关系代词取舍的因素有何显著差异？3) 在可比语言特征下，两个群体的关系代词取舍有何差异，尤其是显性/隐性关系代词在两个群体中的组间差异有何表现？

本研究从中国英语学习者语料库（WECCCL 1.0、WECCCL 2.0）和英语本族语者语料库（LOCNESS、NESSIE、ICNALE）中抽取1,726句非主语型限定性关系从句，对其标注多层面语言学特征，然后分别从抽象特征层面和具体词语层面开展四种统计分析：单因素分析、二元逻辑斯蒂回归、双回归分析、共现词分析。

在抽象特征层面，二元逻辑斯蒂回归和双回归分析显示：1) 关系代词取舍构式受到一系列因素的制约，并非随机分布，中国英语学习者和英语本族语者在关系代词取舍中存在九个共同的制约因素：先行词语义丰富度、先行词的生命性、先行词确定性、先行词修饰语数量、先行词在主句中充当的句法成分、关系从句主语词性、关系从句主语确定性、从句肯定/否定、从句长度；2) 二者存在四个差异性制约因素：关系从句长度、先行词具体性、从句主语具体性、主句谓语动词类型；3) 隐性/显性关系代词构式在两个群体中的组间差异并不相同：显性关系代词的组间差异显著大于隐性关系代词的组间差异，显性关系代词的组间差异与显著制约变量的认知加工难度呈现负相关性。

在具体词语层面，共现词分析显示：1) 中国英语学习者和英语本族语者对隐性关系代词的使用较为相似，对显性关系代词的使用差异较大；2) 在两个语言群体共有的共现词中，中国英语学习者对高搭配强度的共现词依赖程度更加显著。

综合上述研究结果，可以发现：关系代词取舍构式中存在一个由具体词语构式和抽象图式构式组成的连续统；从具体词语构式到抽象图式构式存在竞争—浮现的认知加工机制；与英语本族语者相比，中国英语学习者的关系代词取舍构式在构式连续统中更为靠近具体词语层面。

通过挖掘关系代词取舍构式的认知加工机制，本研究实证检验了中介语的两个理论特征，尝试确认了关系代词取舍构式的默认原型，并论证了基于使用的研究视角在二语习得研究中的适用性。

本书的撰写得到了诸多前辈、学友的指导和帮助。首先，我要感谢我的导师梁茂成教授，如果没有他引领我进入统计编程的世界，我便难以完成本研究；其次，我要感谢李文中教授、许家金教授、熊文新教授、刘振前教授、苗兴伟教授、王立非教授、徐昉教授，他们在百忙中审读了本书的初稿，并提出了诸多宝贵的意见和建议。同时，我要感谢Stefan Gries博士、Sandra Deshors博士、Roger Levy博士为我解答诸多技术层面的困惑。

由于本人水平有限，本书难免有疏漏不妥之处，恳请广大读者批评指正！

目 录

结论	1
0.1 研究背景	1
0.2 选题缘由	3
0.2.1 理论层面	3
0.2.2 方法论层面	4
0.2.3 实践层面	4
0.3 研究概述	5
0.3.1 语料抽取	5
0.3.2 语料标注	6
0.3.3 统计分析	6
0.3.4 理论阐释	7
0.4 本书结构	7
 第一章 限定性关系从句研究回顾	 9
1.1 关系从句的界定	9
1.2 限定性关系从句研究现状	10
1.3 关系从句研究中的理论假设	12
1.3.1 基于内省或实验法的假设	12
1.3.2 基于真实语料的假设	17
1.3.3 简评	20
1.4 关系代词相关研究回顾	21
1.4.1 以内省/实验法为主导的研究	21
1.4.2 基于语料的多因素分析的萌芽	22
1.5 前人研究的贡献与不足	26

1.5.1 前人研究的贡献	26
1.5.2 前人研究的不足	26
1.6 小结	28

第二章 中介语与基于使用研究视角下的关系从句研究 29

2.1 中介语研究	29
2.1.1 中介语的理论特征	29
2.1.2 中介语研究与统计的结合	32
2.1.3 对本研究的启示	34
2.2 中介语中的限定性关系从句研究	34
2.2.1 Gass (1979) 的中介语限定性关系从句研究	34
2.2.2 中国学生英语关系从句相关研究	35
2.2.3 前人研究的贡献与不足	37
2.3 基于使用的限定性关系从句研究视角	38
2.3.1 基于使用的研究范式	38
2.3.2 作为构式的限定性关系从句	41
2.3.3 限定性关系从句的认知加工理论	42
2.4 小结	45

第三章 语言学研究中的多因素分析 47

3.1 多因素分析的界定	48
3.2 语言学中多因素分析的发展阶段	48
3.2.1 基于人工的特征分析	49
3.2.2 基于统计模型的多因素分析	51
3.3 现有多因素分析研究分类	52
3.3.1 按照研究对象分类	52
3.3.2 按照本族语/中介语分类	53
3.4 多因素分析中的统计方法	54
3.4.1 线性分类器	55
3.4.2 逻辑斯蒂回归	55
3.5 小结	55

第四章 研究方法与数据	57
4.1 研究问题	57
4.2 研究工具	57
4.2.1 Stanford Parser	58
4.2.2 PowerGREP	58
4.2.3 Tregex	58
4.2.4 统计编程语言R	59
4.3 语料	60
4.3.1 中国英语学习者语料库	60
4.3.2 英语本族语者语料库	61
4.3.3 非主语型限定性关系从句的抽取	61
4.3.4 非主语型限定性关系从句的筛选	62
4.4 非主语型限定性关系从句的标注	66
4.4.1 先行词相关因素标注	67
4.4.2 从句主语相关因素标注	75
4.4.3 主句相关因素标注	80
4.4.4 从句相关因素标注	84
4.5 统计方法	86
4.5.1 单因素分析	87
4.5.2 二元逻辑斯蒂回归分析	87
4.5.3 双回归分析	87
4.5.4 共现词分析	89
4.6 小结	89
第五章 多因素统计模型分析	90
5.1 数据概述	91
5.2 统计分析流程	92
5.3 单个语言学因素与关系代词取舍的相关性	93
5.3.1 共有显性相关自变量	93
5.3.2 两个语言群体各自独有显著变量	109
5.4 语言学因素的交互效应	114
5.4.1 模型检验	115
5.4.2 主效应分析	116

5.4.3 交互效应分析	120
5.4.4 因素交互效应小结	124
5.5 可比语言特征下关系代词取舍的组间差异	125
5.5.1 第一轮回归	125
5.5.2 第二轮回归	126
5.6 共现词分析	132
5.6.1 先行词分析	133
5.6.2 从句主语分析	135
5.7 小结	138
第六章 分析与讨论	140
6.1 中国英语学习者和英语本族语者在关系代词取舍中的制约因素	140
6.1.1 共有制约因素	141
6.1.2 差异性制约因素	143
6.1.3 可比语言特征下关系代词的组间差异	145
6.2 对关系代词取舍的认知加工解读	146
6.2.1 关系代词取舍中的多因素竞争	147
6.2.2 关系代词取舍构式的浮现	148
6.2.3 关系代词取舍构式的能产性	151
6.3 中介语理论特征的实证检验	152
6.3.1 中介语的体系性	153
6.3.2 中介语的可渗透性	154
6.4 关系代词取舍的原型性构式	155
6.4.1 隐性/显性关系代词的频次	156
6.4.2 隐性/显性关系代词构式的类型数	156
6.4.3 隐性/显性关系代词构式的认知加工差异	156
6.5 小结	157
第七章 结论	159
7.1 研究发现	159
7.1.1 中国英语学习者和英语本族语者对关系代词取舍的相似性	159

7.1.2 中国英语学习者和英语本族语者在关系代词取舍中的差异性	160
7.1.3 中国英语学习者和英语本族语者在关系代词取舍中的差异程度	160
7.2 研究价值与创新	161
7.2.1 理论层面	161
7.2.2 方法论层面	163
7.2.3 实践层面	163
7.3 研究局限与今后拓展	165
7.3.1 研究局限	165
7.3.2 今后拓展方向	165
参考文献	167
附录	190

表 目

表 1-1 内省/实验法理论假设在本研究中的应用	12
表 1-2 Wiechmann (2007) 关系从句部分标注因素	25
表 1-3 隐性关系代词的显著相关因素 (Wiechmann 2007)	25
表 2-1 Jarvis (2000) 干扰变量	33
表 3-1 词汇层面的已有多因素分析研究	52
表 3-2 句法层面的已有多因素研究	53
表 4-1 非主语型限定性关系从句在各个语料库中的分布	62
表 4-2 非主语型限定性关系从句标注体系	66
表 4-3 限定词与关系代词that的取舍关系 (Wasow et al. 2011: 179)	69
表 4-4 先行词句法成分分类	74
表 5-1 非主语型限定性关系从句数据分布	91
表 5-2 NS与NNS共有显著因素	94
表 5-3 英语本族语者独有显著自变量	109
表 5-4 中国英语学习者独有显著自变量	109
表 5-5 二元逻辑斯蒂回归最终模型	115
表 5-6 线性回归模型显著变量	128
表 5-7 隐性/显性关系代词结构中先行词共现词 (中国 英语学习者)	134