

Computer
Simulation of Visual
Word Reading

词汇阅读的 计算机模拟

杨剑峰 著



陕西师范大学出版社

国家自然科学基金面上项目
视觉词汇加工的动态神经网络及其形成(31171077)

词汇阅读的计算机模拟

杨剑峰 著

陕西师范大学出版总社

图书代号 ZZ16N1355

图书在版编目(CIP)数据

词汇阅读的计算机模拟 / 杨剑峰著. —西安: 陕西
师范大学出版总社有限公司, 2016.9

ISBN 978-7-5613-8688-0

I. ①词… II. ①杨… III. ①计算机模拟—应用
语言学—研究 IV. ①TP302.1

中国版本图书馆 CIP 数据核字(2016)第 247452 号

词汇阅读的计算机模拟
CIHUI YUEDU DE JISUANJI MONI

杨剑峰 著

责任编辑 古洁
责任校对 樊荣
封面设计 鼎新设计
出版发行 陕西师范大学出版总社
(西安市长安南路 199 号 邮编 710062)
网 址 <http://www.snupg.com>
经 销 新华书店
印 刷 陕西金德佳印务有限公司
开 本 787mm×1092mm 1/16
印 张 12.5
字 数 200 千
版 次 2016 年 9 月第 1 版
印 次 2016 年 9 月第 1 次印刷
书 号 ISBN 978-7-5613-8688-0
定 价 30.00 元

读者购书、书店添货或发现印装质量问题, 请与本社高等教育出版中心联系。
电话:(029)85303622(传真) 85307864

前　　言

20世纪90年代，在认知科学领域，联结主义对信息加工论的观点提出了挑战，从而掀起了认知科学的一场革命。信息加工论将人脑等同于电脑，认为认知就是心理符号的运算过程。而联结主义则在真正意义上将认知过程等同于大脑的功能，基于神经元的生理机制来理解人的认知加工过程。计算机模拟是联结主义理论强有力的技术手段，通过建构基于人工神经元网络的计算机模型，尝试模拟正常成人、儿童以及认知障碍等特殊人群的行为表现，借助神经元群或神经网络的功能来解释具体的认知加工过程，以实现对人脑功能的理解。

本书旨在介绍词汇阅读的联结主义取向。

首先，语言认知科学的研究仍然关注着三个基本的理论问题：即知识是如何表征、获得及应用的。心理语言学的研究受乔姆斯基的传统理论影响，认为语言获得是先天的、不可获得的。联结主义的兴起和迅速发展，对语言获得与加工的三个基本理论问题给出了完全不同的答案。本书第一章详细介绍了联结主义对上述三个问题的回答，在这些基本理论问题上的答案也是联结主义理论解释词汇阅读的理论基础。

具体到词汇阅读领域，联结主义提出了阅读的三角模型。本书第二章详细介绍了联结主义理论对词汇阅读的理论解释，基于三角模型的计

算机模型不仅能成功模拟出正常成人阅读的各种行为表现,还能成功模拟获得性阅读障碍以及发展性阅读障碍的不同亚类型,并能给出合理的理论解释。同时,联结主义三角模型最核心的思想是提倡词汇阅读过程中语音和语义加工的交互作用。本书第三章通过多方面的证据来详细阐述这种交互作用的内部机制。

其次,通过几个具体的联结主义计算机模型在汉字阅读中的应用,本书介绍了人工神经元网络在词汇阅读中的具体应用及其方法。第四章详细介绍了不同的字形、语音以及语义表征方案,通过对比已有表征方案的不足,从而确定汉字的字形和语音表征方案。尤其是仍然处于探索阶段的语义表征方案,我们尝试了多种不同的表征方案,在今后的计算机模拟乃至语义信息检索等领域都具有非常重要的借鉴作用。随后在第五章具体介绍了几个汉字阅读的计算机模拟研究,通过对汉字阅读中的基本行为表现的模拟,一方面表明基于拼音文字系统的联结主义模型能成功概括到汉字阅读中;另一方面表明联结主义模型的内部加工机制具有跨语言的普遍性。

最后,认知神经科学已经成为当前心理语言学研究的主流。长期以来,语言认知神经科学致力于揭示语言加工对应的大脑功能脑区,如视觉词形识别区、语音/语义加工脑区等。近年来,研究者开始关注词汇阅读的神经生理模型与认知理论模型的统一。联结主义认知理论立足于神经网络的思想,具有统一认知和神经模型的可能,从而成为当前语言认知神经科学研究的一大重要研究取向。本书在第六章详细介绍了当前词汇阅读的认知神经科学进展,并试图在联结主义取向下理解词汇阅读的大脑神经机制。

本书努力尝试将联结主义理论在词汇阅读中的应用,以及涉及到的相关信息都涵盖进来。但因水平有限,实感力不从心,只能权当是一次知识的积淀过程,希望能对今后的词汇阅读研究起到一定的借鉴作用。如果此书能为从事词汇阅读研究的学者或研究生提供一个中文参考,也就甚感欣慰了。

杨剑峰

2016年9月于陕西师范大学

目 录

第一章 语言加工及其获得的概率限定理论	(1)
第一节 联结主义语言获得理论	(3)
1.1 基于统计的语言获得研究	(3)
1.1.1 词汇分段的统计学习	(3)
1.1.2 语法的统计学习	(5)
1.1.3 争论	(6)
1.2 语言获得的联结主义计算模型	(7)
1.3 基于语料库的统计研究进展	(10)
第二节 视觉词汇阅读的概率限定取向	(11)
2.1 阅读获得的概率限定——粒度理论	(12)
2.2 阅读获得与阅读加工的关系	(16)
2.3 语音、语义对阅读加工的概率限定	(18)
第三节 小结	(22)
3.1 天性与教养的问题	(22)
3.2 知识的表征问题	(23)

3.3 语言知识的获得问题	(23)
第二章 词汇阅读的理论模型	(25)
第一节 双通道理论	(27)
1.1 传统的双通道理论	(27)
1.1.1 基本原理	(27)
1.1.2 假词阅读的词典效应	(29)
1.1.3 真词阅读的词典效应	(30)
1.1.4 认知神经心理学研究	(31)
1.2 修正的双通道理论	(32)
1.2.1 赛马假说	(32)
1.2.2 多通道理论	(33)
1.2.3 加法假说	(34)
1.2.4 对照模型	(35)
1.2.5 双通道瀑布模型	(36)
第二节 联结主义理论取向	(37)
2.1 词汇阅读的联结主义三角模型	(37)
2.2 联结主义对正常阅读的解释	(38)
2.3 联结主义对获得性阅读障碍的解释	(40)
2.3.1 表层阅读障碍	(40)
2.3.2 语音阅读障碍	(41)
2.3.3 深层阅读障碍	(44)
2.4 联结主义阅读理论的优势	(45)
第三节 阅读理论的分歧与融合	(46)

3.1	两种理论取向的分歧	(47)
3.2	理论整合:两条加工通道的共同作用	(48)
第三章 语音和语义加工在词汇阅读中的作用		(51)
第一节 正常成人阅读		(53)
1.1	正常成人阅读	(53)
1.1.1	行为实验	(53)
1.1.2	计算机模拟	(56)
1.1.3	脑功能成像研究	(56)
1.2	汉字阅读	(57)
第二节 脑损伤患者(获得性阅读障碍)		(60)
2.1	获得性阅读障碍研究	(60)
2.2	获得性阅读障碍的两种理论解释	(62)
2.3	主要系统假说	(64)
2.4	汉语获得性阅读障碍	(69)
第三节 儿童语言发展		(71)
3.1	儿童发展	(71)
3.2	元语言缺陷	(73)
第四节 小结		(74)
第四章 汉字的字形、语音及语义表征		(75)
第一节 汉字字形表征		(76)
1.1	对汉字字形表征的尝试	(76)
1.2	新的汉字表征方案	(78)
第二节 汉字语音表征		(82)

2.1 拼音文字系统的语音表征方案	(82)
2.2 汉语音节表征	(84)
2.3 基于国际音标的汉字语音表征	(85)
第三节 汉字语义表征	(87)
3.1 随机编码	(87)
3.2 基于语义知识库的特征抽取	(88)
3.3 基于语料的共现统计	(90)
3.4 几种语音表征的对比	(93)
第四节 讨论	(93)
第五章 词汇阅读的计算机模拟研究	(95)
第一节 汉字阅读的形—音对应模型	(97)
1.1 模拟 1:汉字阅读的联结主义模型	(97)
1.1.1 计算模型的建构	(97)
1.1.2 模型的训练	(99)
1.1.3 模型的测试	(99)
1.1.4 模拟结果	(101)
1.2 行为实验	(104)
1.2.1 实验方法	(104)
1.2.2 实验结果	(104)
1.3 讨论	(106)
第二节 汉字阅读发展的计算模型	(107)
2.1 模拟 2:儿童阅读发展的模拟研究	(109)
2.1.1 模拟方法	(109)

2.1.2 模拟结果	(110)
2.2 模拟 3: 规则性效应的本质	(112)
2.2.1 模拟方法	(113)
2.2.2 模拟结果	(113)
2.3 讨论	(114)
第三节 语义信息在出声阅读中的作用	(117)
3.1 语义在阅读中的作用	(118)
3.1.1 脑损伤阅读障碍	(118)
3.1.2 反应时数据	(119)
3.2 模拟 4: 语音、语义信息的交互作用机制	(120)
3.2.1 模拟方法	(121)
3.2.2 模拟结果	(123)
3.3 讨论	(126)
第六章 联结主义视角下的阅读神经网络	(129)
第一节 跨语言特异的阅读机制	(130)
1.1 语言特异的阅读认知机制	(130)
1.2 语言特异的阅读神经机制	(131)
1.2.1 视觉词形加工区	(131)
1.2.2 形 - 音转换加工脑区	(134)
1.2.3 汉字阅读的特异性脑区	(135)
第二节 基于神经网络的阅读脑机制研究	(137)
2.1 跨语言统一的阅读理论	(137)
2.1.1 联结主义阅读理论	(137)

2.1.2 汉字阅读的联结主义模型	(140)
2.2 阅读神经网络研究的新进展	(142)
2.2.1 阅读脑区的动态激活	(142)
2.2.2 阅读神经网络	(143)
2.3 中文阅读神经网络	(146)
2.3.1 刺激和任务的交互作用	(146)
2.3.2 汉字特异的阅读脑区可能是实验任务夸大的结果	(148)
第三节 小结	(149)
3.1 文化经验对阅读神经网络的影响	(150)
3.2 阅读脑区的动态模式	(151)
参考文献	(153)
后记	(187)

第一章 语言加工及其获得 的概率限定理论

■ 第一节 联结主义语言获得理论

- ◆ 1.1 基于统计的语言获得研究
 - 1.1.1 词汇分段的统计学习
 - 1.1.2 语法的统计学习
 - 1.1.3 争论
- ◆ 1.2 语言获得的联结主义计算模型
- ◆ 1.3 基于语料库的统计研究进展

■ 第二节 视觉词汇阅读的概率限定取向

- ◆ 2.1 阅读获得的概率限定——粒度理论
- ◆ 2.2 阅读获得与阅读加工的关系
- ◆ 2.3 语音、语义对阅读加工的概率限定

■ 第三节 小结

- ◆ 3.1 天性与教养的问题
- ◆ 3.2 知识的表征问题
- ◆ 3.3 语言知识的获得问题

第一章 语言加工及其获得 的概率限定理论

近年来的心理语言学研究取得了大量的成果,研究探讨的科学问题仍然集中在三个经典问题(Chomsky, 1986):首先是语言知识的本质,也就是在大脑中存储的语言知识的实质是什么?其次是这种语言知识是如何获得的?最后,这种知识又是怎么被使用的,也就是如何使用这些语言知识来理解和产生语言的?

自20世纪50年代以来,Chomsky的生成理论一直主导着心理语言学的研究。这一理论认为语言是由语法规则组成的复杂规则系统,语言知识就是一系列的语法规则,语言加工及其获得就是对语法规则的使用及其获得的过程。这一理论对于语言获得的一个重要主张是“刺激贫乏理论”(Poor Stimulus Argument),认为儿童输入的语言材料是相当零散、不完整的,还充满了矛盾,既有符合语法也有不符合语法的句子;语言环境没有提供可靠的肯定或否定证据说明哪种句子是符合语法规则的;就算是同一种语法结构下,往往有多种表达形式。而语法规则是完整的、复杂的,因而,该理论认为婴儿的语言经验不可能是语言知识的来源,唯一能解释儿童能迅速地获得语言的可能原因就是存在先天的语法知识,这种知识不通过学习(Learning)获得,而是人类的本能行为,并假定存在一种先天的语言获得装置(Language Acquisition Device, LAD),认为婴儿的语言经验主要是获得词典系统,并给LAD设定一系列的语言特异性的

参数。

近 10 年来,随着语言学之外的研究领域的重要发现的增加,逐渐形成了语言获得及其加工的“概率限定(Probabilistic Constraints)”取向,这种新的理论取向代表着语言研究历史上的重大转变(Seidenberg, 1997)。它的核心观点认为语言加工及其获得是对多种同时并存的语言或非语言线索(或信息)的概率限定的使用。

概率限定理论是在联结主义取向下对语言加工及其获得的全新解释,也被称为联结主义语言获得理论。本章第一节介绍联结主义的语言获得理论的基本观点,以及来自三方面的研究证据。第二节介绍联结主义取向下的词汇阅读研究。

第一节 联结主义语言获得理论

联结主义语言获得理论对传统的 Chomsky 的语言获得理论提出了新的挑战,引起了关于语言获得研究的新的热潮。联结主义语言获得理论主要得到了来自三方面的证据支持(Bates & Elman, 1996),即基于统计的语言学习、语料库的统计分析以及联结主义的计算机模拟研究,下面将从这三个方面对这一语言获得理论加以阐述。

1.1 基于统计的语言获得研究

对婴儿的语言获得研究发现,出生不久的婴儿就已经具有了很强的统计学习能力。

1.1.1 词汇分段(word parsing)的统计学习

对于所有的语言学习者来说,一个必须面对的问题就是把连续的语言流切分成单独的词汇。这是一个相当困难的任务,因为一些离散的听

觉事件(如停顿等标志)使得连续语流中的词边界线索不一致。尽管已有研究揭示 8 个月大的婴儿已经能从连续发音中区分出词汇边界,并能在词汇单独呈现时再认出来(Jusczyk & Aslin, 1995),但是婴儿是利用什么信息或线索来发现这些词汇边界的呢?考虑到不同语言的发音结构差异,这个问题就变得更加复杂。与词汇边界相关的发音线索可能与婴儿的母语相关,这时婴儿可能利用这些线索来发现它们。但是,当语言材料中没有与词汇边界相关的线索起作用时,婴儿又如何发现这些词汇边界呢?

Saffran 等人在 *Science* 上发表的一项研究(Saffran, Aslin, & Newport, 1996)中,研究者对 24 个生活在英语环境中的 8 个月大的婴儿进行测试。他们给这些婴儿呈现以随机顺序反复播放的 4 个 3 音节构成的无意义词汇(如,bidaku/padoti/golabu/bidaku…),语音流是由一个单调的电子装置发出的女声,速度为每分钟 270 个音节,共呈现 180 个单词。发音装置不会提供任何有关字词汇边界的声音信息,没有停顿、重音或其它指示词边界的韵律线索。而唯一可能的词汇边界线索就是音节之间的转换概率,如,在词汇内部的两个音节间的置换概率固定为 1(如 bida),要高于词汇间的两个音节概率(如 kupa 的 ku 与 pa 之间的置换概率为 0.33)。在对语料的熟悉阶段,仅仅给婴儿播放了 2 分钟的语音。测试时,实验 1 使用在熟悉阶段婴儿听过的音节,但组合顺序不同的“非词”。实验 2 使用人工的“部分词”(padoku)作为测试材料,由一个词(bidaku)的后一个音节(ku)和另一个词(padoti)的前两个音节(pado)组成。学习者必须能提取出同时呈现的声音对的相应的频率信息,相应的低转换概率暗示着词边界,才能作出对词与“非词”、词与“部分词”的区分。

实验结果表明,婴儿在词与“非词”的刺激材料上表现出了显著的差别,在“非词”上注意听的时间更长。说明 8 个月的婴儿已经能够识别出新异的和熟悉的三音节组合。即使是在需要更为困难的统计计算的实验

►►► 第一章 语言加工及其获得的概率限定理论

2 中,婴儿也表现出了对“部分词”的听觉偏好,表明婴儿的确能够从 2 分钟的语料输入中抽取出音节之间的这种统计属性。

同样的实验范式在成人的实验(Saffran, Newport, & Aslin, 1996)以及其它的语音实验(Aslin, Saffran, & Newport, 1998; Newport & Aslin, 2000; Saffran, Johnson, Aslin, & Newport, 1999)中也得到了类似的基于统计学习的结果。这些研究重新挑起了人们对于儿童语言获得的思考。

当然,这种基于语料统计学习的思想受到了多方面的挑战。基于统计的学习理论首先面临的问题就是这种统计学习机制究竟是认知某一领域(听觉)特殊的学习机制,还是一种更为一般的学习机制(Pesetsky et al., 1997)。在关于儿童视觉统计学习的研究中(Kirkham, Slemmer, & Johnson, 2002),研究者使用习惯化的实验范式,对 2、5、8 个月大的婴儿呈现离散的视觉刺激,这些视觉刺激的呈现顺序具有统计的规律。婴儿在观看了熟悉的视觉图形序列之后,紧接着观看新异的图形序列,三个年龄段的婴儿都表现出了对新异图形序列的转头偏向。表明婴儿能从连续呈现的视觉图形中抽取出刺激间的过渡概率,而且这种视觉统计学习在 2 个月的婴儿身上就已经表现出来了。这些在多种通道发现的统计学习机制似乎表明人类具有更为一般的统计学习能力。

甚至,在对动物的语言学习实验中,使用与听觉统计学习相同的刺激材料(Jusczyk & Aslin, 1995),卷毛猴也能在基于统计概率的基础上分辨出词与“非词”、词与“部分词”的区别,表明动物也具有这种简单的基于统计的学习机制。

1.1.2 语法的统计学习

人类语言的大量规则并不像词汇切分这么简单,可以从简单的转换概率来完成任务。语言包含有大量的语法规则,人类能依赖这种统计概率来抽取复杂的、繁多的语法规则吗?对更复杂的语言问题,如语法,是