

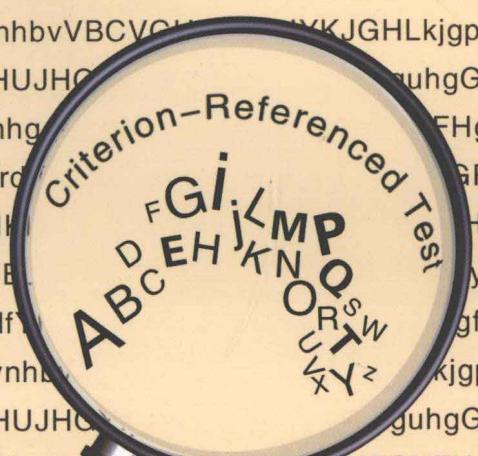
# 标准参照语言测试研究

Criterion-Referenced  
Reference  
Criterion-Referenced Test

黄锐 著



ASCDghgyHFGGghGUYGhggfGGhjgghfGFGHGFgftFGHfvVGHfgghgFG-  
FhgfvGH-  
fghfhgGHFfgfFHGtyGHGtrTEUTIftrtytyrgFTYRYTFtftyrtiyfgFTERTDGty  
ftyfgFTYRTFghfyufjGFTYRTYFghfjghftygFKLGhjguyDTFGjvkgoIGYGH  
GFgfyuglhgyuGYGYKJGHgyughGYJHjlguygJHGyghFGGgFYUFGHFGHD  
FFDajkhflfYFGhfGFGhkFGdfFGDFgdjkhgfhFGHFTgfgfJYFTYFghfjfyhgf  
GHDFGXvhbvVBCVQJLJKNJGHLKjgphjhgGHFTGHFGHFtyrteUIPl  
okjhjkghuiHUJHC  
jKNJKHJghhg  
UYTYFGhfrd  
JGJKHJKHH  
fgVBNMNE  
OFFDajkhflf  
GHDFGXvhb  
kjhjkghuiHUJHC  
KNJKHJghhgaf  
(LjkbhbgV  
BJHGYGY  
yugjhGYJHjlguyg



厦门大学出版社 国家一级出版社  
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位

# 标准参照语言测试研究

黄锐 著



厦门大学出版社 国家一级出版社  
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位

## 图书在版编目(CIP)数据

标准参照语言测试研究/黄锐著. —厦门:厦门大学出版社,2012.10  
ISBN 978-7-5615-4456-3

I. ①标… II. ①黄… III. ①语言-测试-研究 IV. ①H09

中国版本图书馆 CIP 数据核字(2012)第 248919 号

厦门大学出版社出版发行

(地址:厦门市软件园二期望海路 39 号 邮编:361008)

<http://www.xmupress.com>

xmup @ xmupress.com

厦门集大印刷厂印刷

2012 年 10 月第 1 版 2012 年 10 月第 1 次印刷

开本:889×1194 1/32 印张:9.75

字数:259 千字 印数:1~2 000 册

定价:36.00 元

本书如有印装质量问题请直接寄承印厂调换

**本书获“集美大学优秀学术  
著作出版基金”资助**

# 序

黄锐的专著《标准参照语言测试研究》是在她的硕士学位论文《现代教育测量理论在标准参照语言测试中的应用与案例研究》的基础上,经过她多年的努力和孜孜不倦的耕耘、积累、拓展和深入研究,最终撰写而成的。该书现正式出版,作为她的老师,我感到颇为欣慰。

自从上个世纪 60 年代 Glaser 提出“标准参照测量”概念后,在教育测量界有关它的讨论就热门起来,国内对此的研究也在 90 年代迅速开展起来,到本世纪初达到高潮。中国汉语水平考试自 1984 年研制实施以来,语言测试界对标准参照测试的研究也随之发展,但就英语作为外语的语言测试界来说,对标准参照测试的研究却相对落后,国内目前尚无一本有关此内容的专著,研究该方向的文章也不多见。本书作者能在教育测量界和汉语测试界在标准参照测试研究的基础上,结合外语测试的特点撰写此书是一种大胆尝试。

本书从语言测试理论出发,通过与常模参照测试的理论和发展建构过程的比较,探讨了标准参照语言测试的理论和应用。信度和效度研究是语言测试中的难点,而标准参照测试的信度和效度更难把握。作者通过反复研读文献,以及几年的思考和写作,系统地梳理了标准参照语言测试的信度估计和效度研究的理论和方法,为同行们提供了一个比较全面了解标准参照语言测试全貌的专著。此外,作者还结合自身参与普通高等学校招生命题的工作经验,在书中就外语教师如何更好地参与命题、改卷、试题分析等方面的工作,提出了一些建设性意见。

该书对语言测试研究者、外语教师、高校师范生和语言测试的硕士生都有较高的参考和应用价值，是一本很值得阅读的专著。  
是为序。

史秋衡

2012年9月21日

于厦门大学颂恩楼

# 前 言

本书以语言测试的基本理论为基础,通过比较常模参照测试的理论和发展建构过程,探讨了标准参照语言测试的理论和应用,重点说明了标准参照语言测试的编制、项目分析以及标准参照语言测试的信度、效度和分数的解释。对标准参照语言测试的难点,即信度和效度进行了比较全面的探讨和研究,并通过理论联系实际的案例,采用了NEATs和TEM等国内大型的标准参照语言测试的统计数据,介绍了经典测量理论、概化理论和项目反应理论在标准参照语言测试的信度估计和效度验证中的应用,为同行学者提供了较多可参考和借鉴的依据。本书的主要读者对象为语言测试研究人员和外国语言学及应用语言学专业的硕士生、英语教师、英语专业的师范生以及对语言测试感兴趣的人士。

本书得以完稿和出版,首先要感谢我的授业恩师,教育部评估专家、国家“985工程”创新基地副主任、厦门大学高等教育质量与评估研究所所长、教育研究院副院长、博士生导师史秋衡教授,尽管毕业多年,但他始终鼓励我在自己的研究领域里进行不断的探索和研究,并为本书作了序言;其次要感谢集美大学资助本书的出版并为我提供到广东外语外贸大学访学一年的机会,让我不但能全身心地完成书稿,还能直接地向语言测试界的专家们学习;感谢我的访学指导老师,广东外语外贸大学副校长、博士生导师刘建达教授和广东外语艺



术学院院长、博士生导师曾用强教授,是他们的教导和鼓励,让我能顺利地完成书稿;刘建达教授在百忙中通读了本书,并对书中的不足提出了宝贵的修改意见;广东外语外贸大学的博士生导师亓鲁霞教授对本书的目录和书名也提出过修改意见,在此一并表示感谢。

由于本人才疏学浅,书中难免有错误和不足,望读者批评指正。

黄 锐

2012年7月于广东外语外贸大学



## 缩略语词汇

AERA	American Educational Research Association
APA	American Psychological Association
CI	Confidence Interval
CET 4 & 6	College English Test, Band 4&6
CRT	Criterion-referenced Test
CRLT	Criterion-referenced Language Test
CTT	Classical Testing Theory
DI	Difference Index
ETS	Educational Testing Service
ICC	Item Characteristic Curve
ID	Item Discrimination
IF	Item Facility
IELTS	The International English Language Testing System
IRT	Item Response Theory
IIFs	Item Information Functions
LAD	Language-acquisition Device
LID	Local Item Dependence
MLAT	Modern Language Aptitude Test
NCME	National Council on Measurement in Education
NEAT	National English Achievement Tests
NRT	Norm-referenced test
PETS	Public English Test System
PLAB	Pimsleur Language Aptitude Battery
SEM	Standard Error of the Measurement
TCC	Test Characteristic Curve
TEM 4 & 8	Test English for Major, Grade 4&8
TOEFL	The Test of English as a Foreign Language

# 目 录

## 第一章 绪论 /1

- 1.1 研究的目的和意义 /1
- 1.2 基本概念界定 /4
  - 1.2.1 测试、考试、测量、评估和评述 /4
  - 1.2.2 语言测试 /8
  - 1.2.3 标准参照测试与标准参照语言测试 /9
  - 1.2.4 项目反应理论 /10

## 第二章 语言测试的发展和理论架构 /12

- 2.1 语言教学发展的四大阶段 /12
  - 2.1.1 科学前阶段 /12
  - 2.1.2 结构主义语言学 /13
  - 2.1.3 认知法与转换生成语法语言学 /14
  - 2.1.4 交际能力语言教学法 /15
- 2.2 语言测试的发展史 /16
  - 2.2.1 短文写作—翻译测试法 /16
  - 2.2.2 结构主义—心理测量法 /16
  - 2.2.3 综合测试法 /17
  - 2.2.4 交际测试法 /17
- 2.3 语言测试的实质 /18
- 2.4 语言测试的功能及其种类 /19
  - 2.4.1 测试、考试、测量、评估和评述的关系 /19
  - 2.4.2 测试的功能 /21
  - 2.4.3 测试的种类 /24



2.5 语言测试的理论构架/31

2.5.1 信度/31

2.5.2 效度/36

2.5.3 真实性/44

2.5.4 交互性/48

2.5.5 影响/56

2.5.6 可操作性/57

**第三章 标准参照语言测试的理论探讨/59**

3.1 常模参照测试概述/59

3.1.1 常模参照测试的含义/59

3.1.2 常模参照测试编制原则/60

3.2 标准参照测试的兴起与内涵/61

3.2.1 标准参照测试理论的兴起/62

3.2.2 标准参照测试的概念及内涵/64

3.2.3 标准参照与常模参照的异同/67

3.2.4 标准参照语言测试的含义和作用/73

3.3 标准参照语言测试的编制及使用/75

3.3.1 编制的基本原则/75

3.3.2 项目分析参数/77

3.3.3 及格的标准水平/79

3.3.4 信度估计/90

3.3.5 效度验证/91

**第四章 测试成绩分析与标准参照/93**

4.1 测试成绩/94

4.1.1 原始分数/94

4.1.2 转换分数/94

4.1.3 Z 分数和 T 分数/95

4.2 分数的频数分布/96

4.2.1 分数的整理/96

4.2.2 分数的频数分布/98

4.3 分数频数的图形显示/99

4.4 分数的集中量/102

  4.4.1 平均数/102

  4.4.2 众数/102

  4.4.3 中位数/103

  4.4.4 集中量的比较/103

4.5 离散量/104

  4.5.1 全距/105

  4.5.2 四分位区间距/105

  4.5.3 平均差/106

  4.5.4 方差与标准差/106

4.6 分数的分布/108

4.7 偏态值和峰值/109

## 第五章 标准参照语言测试应用理论的分析/111

5.1 经典测试理论的建立背景/111

5.2 经典测试理论在标准参照语言测试中的应用/113

  5.2.1 项目分析基本原理/114

  5.2.2 传统项目分析对常模参照测试的影响/115

5.3 标准参照语言测试的项目分析/123

  5.3.1 概述/123

  5.3.2 测试项目的差异性指标/126

5.4 小结/132

## 第六章 标准参照语言测试与项目反应理论/134

6.1 经典测试理论的局限性/134

6.2 现代测试理论的基本原理/136

  6.2.1 项目反应理论的基本概念/136

  6.2.2 项目反应理论的基本假设/140

6.3 基本的试题反应模型介绍/141

  6.3.1 单参数模型/142

  6.3.2 双参数模型/144



6.3.3 三参数模型/145

6.3.4 其他常用的模式/146

6.4 三种不同模型的优点和不足/147

6.5 IRT 对 CRLT 问题的实际应用/148

6.6 小结/154

## 第七章 项目反应理论在标准参照语言测试中的应用/156

7.1 能力与试题参数的估计/157

7.1.1 能力参数的估计/158

7.1.2 其他估计方法与计算机程序/160

7.2 信息函数/161

7.3 多面 Rasch 项目反应理论/167

7.4 IRT 在 CRLT 中应用的案例研究/168

7.4.1 数据背景说明/169

7.4.2 整套试题测试成绩分布情况/170

7.4.3 结论与启示/177

## 第八章 标准参照语言测试的信度、可靠性和单维性/179

8.1 标准参照语言测试的可靠性/180

8.1.1 基本概念:一致性、信度和可靠性/180

8.1.2 阈限损失一致性方法/181

8.1.3 概化理论和领域分数可靠性/186

8.2 概化理论在标准参照语言测试中的应用/195

8.2.1 概化理论的单面交叉设计模型/195

8.2.2 概化理论在标准参照语言测试中的具体应用/198

8.3 信度和可靠性系数的关系比较/199

8.4 项目反应理论一致性问题:单维性、局部独立性及  
模型拟合/201

8.4.1 单维性/204

8.4.2 局部独立性/207

8.4.3 模型拟合/208

8.5 IRT 模型与局部独立性假设问题的认识/213

- 8.5.1 局部独立性假设/213
- 8.5.2 局部依赖问题/216
- 8.5.3 题组中局部试题依赖问题的解决/220

## 第九章 标准参照语言测试的效度研究/223

- 9.1 概述/224
- 9.2 内容效度/227
  - 9.2.1 内容效度的理论论证法/227
  - 9.2.2 内容效度的专家判断法/231
- 9.3 构念效度/234
  - 9.3.1 干预组群构念效度研究/234
  - 9.3.2 差异群体构念效度研究/237
  - 9.3.3 分层结构的构念效度研究/240
- 9.4 内容效度和构念效度的关系/244
  - 9.4.1 关系/244
  - 9.4.2 联合内容和构念效度/245
- 9.5 效度扩展观/247
  - 9.5.1 Messick 的效度观/247
  - 9.5.2 Cronbach 的效度观/252
  - 9.5.3 对标准参照语言测试作决策/254
- 9.6 小结/264

## 第十章 标准参照语言测试的分数报告、反馈和管理/265

- 10.1 概述/265
- 10.2 标准参照语言测试的开发/265
  - 10.2.1 团队开发/266
  - 10.2.2 加强教师间的协作关系/267
  - 10.2.3 分配充足的资源/271
  - 10.2.4 均衡标准参照形式/271
- 10.3 提供标准参照反馈/273
  - 10.3.1 谁应获得反馈/273
  - 10.3.2 双向交流式反馈/276

10.4 报告标准参照结果 / 278

10.4.1 考试不想及格的学生 / 278

10.4.2 对只参与了前测的学生的处理方法 / 279

10.4.3 提高分的解释 / 279

10.4.4 CRLT 成绩报告的困难 / 280

参考文献 / 282



# 第一章 絮 论

## 1.1 研究的目的和意义

中国教育心理界的权威人士张厚粲教授(1992)在她的《考试改革与标准参照测试》的开篇这样说道：“任何教育和心理测试的所得分数都必须有所参照才有实际意义，即测试分数只有与测试以外的某种指标进行比较，才能对应试者的分数给予适当的解释。”可见，分数的解释必须有个特定参照，才能全面地进行解释。在教育测量文献中，自 Glaser 在 20 世纪 60 年代提出“标准参照测量(Criterion-referenced measurement)”至今大概半个世纪的时间，标准参照测试对语言测试界来说仍然是新领域，因此，许多教师和考试管理者对其理论了解并不多，也不够深入。在国内的 CNKI 网站上，以“标准参照”和“语言测试”为关键词的文章，不足 5 篇，检索“标准参照”主题词，则多为教育或心理测量及医学类文献。故而，有必要从语言测试角度来探讨标准参照测试。尽管有学者认为要“告别标准参照测验和常模参照测验二元划分”(罗莲 2007)，因为从同一测试得到的分数可作出常模参照和标准参照两种解释，两者是从分数解释的意义上划分的，并非两种不同的考试，但是，经过多年的实践和探讨，测量学界的研究者们基本上持两种不同的观点，绝大多数人认为，虽然标准参照和常模参照测试有一定区别，但两者是从不同角度对分数进行解释的。另外一个观点则认为，需要明确区分常模参照和标准参照测试，将标准参照测试的内容范围进行精确的定义后，其产生的分数可以进行常模参照的解释，但反过来是不行的(张凯 2002)。不管持哪种观点，我们知道，标准参照测试是在常模参照测试(Norm-



referenced Tests, NRT)的理论基础上发展起来的,可以认为是两种不同性质的考试,但又有许多相似之处。在应用标准参照测试的理论处理大规模考试时,又必须引入一些常模参照测试的办法,同时,在常模参照测试中,也借用标准参照测试中的一些思路来深入细致地分析题项和数据。两者互相借鉴和补充,从而丰富了测试理论。

上世纪 80 年代初,在我国影响和规模最大的普通高等学校招生全国统一考试中,引进了国外现代教育测量的理论和技术,对传统考试的办法进行了改革,形成了考试科学化、现代化的热潮。高考作为我国高等学校选拔新生的考试,最早它取用的是一种常模参照测试(张厚粲 1992),即学生从试卷上得到的成绩要和所有参加考试的考生成绩相比较,从而反映他在考生总体(常模总体)中的相对位置;随着教育改革的不断深化,考试改革已不仅局限在考试形式的改革,而是在考试内容、考试制度方面改革的深入发展,人们进一步要求完善考试的功能。随着高中新课改的不断深化,任何一名普通高中学生都需面对高中毕业会考,而这种旨在客观地评价中学生的学业情况,促进学生德、智、体、美、劳的全面发展,并用已定的标准为对照检验学生是否达到培养目标的考试,不再只是“常模参照测试”,而是标准参照测试(Criterion-referenced Tests, CRT)。如果参加会考的师生们还用常模参照测试中的理论和技术选定题项,他们就会发现很好的试题会被弃置不用。因此,靠单一的常模参照测试理论来选定题项、解释分数已远远不能满足当今日益完善的考试制度。

从国外教育测量史来看,20 世纪初,在实验和教育统计的基础上形成了教育测量的理论,在 30 年代,教育测量的理论应用于选拔和人员安置的考试之中,到了四五十年代,这种着重于研究人员之间差异的常模参照性的测试运动达到了高潮。在这之后,随着教育、教学和教育测量研究的不断深入,人们已不满足于只用来反映考生做得好不好好的测试结果,而是希望通过测试进一步了解考生在什么方面好,什么方面差以及优劣的程度,这样美国心理学家 R. Glaser 和 D. J. Klaus 首先于 1962 年提出了“标准参照测验”这一概念,其目的