

目 录

表 目	IX
图 目	XIII
第一章 导论	1
1.1 研究背景	1
1.1.1 《中国英语能力等级量表》	1
1.1.2 《中国英语能力等级量表》之口语量表	2
1.1.3 口语量表的效度研究	3
1.1.4 口语量表的效度取证	5
1.2 研究目的和内容	6
1.3 本书结构	7
第二章 语言力量表及其效度验证	8
2.1 语言力量表	8
2.1.1 主要的语言力量表	9
2.1.2 语言力量表的构成	11
2.1.3 语言力量表的作用	12
2.1.4 语言力量表的问题	15
2.1.5 对口语量表开发的启示	16
2.2 语言力量表的效度验证	17
2.2.1 效度理论和效度验证	17

2.2.2 语言力量表的效度问题	19
2.2.3 语言力量表的效度框架	22
2.2.4 语言力量表的内部效度	24
2.2.5 语言力量表的外部效度	29
2.2.6 语言力量表的后果效度	32
2.2.7 对口语量表效度研究的启示	34
第三章 研究设计	35
3.1 研究框架	35
3.2 研究步骤	38
3.3 研究方法	41
3.3.1 因子分析	42
3.3.2 项目反应理论	42
3.3.3 Rasch 模型	43
3.3.4 IRT 等值	44
3.3.5 标准设定	47
第四章 口语量表的内部效度验证——结构效度、级别效度验证	50
4.1 口语量表的描述参数和能力级别	50
4.2 研究问题和步骤	51
4.3 数据收集	52
4.3.1 调查对象	52
4.3.2 问卷设计	53
4.3.3 分级验证等值设计	55
4.4 结果分析	57
4.4.1 量表结构验证	57
4.4.2 量表级别验证	95

4.5 结语	112
第五章 口语量表的外部效度验证——口语考试对接口语量表	114
5.1 口语量表和大学英语四级口语考试	114
5.2 研究问题	115
5.3 研究方法和步骤	116
5.4 数据收集	120
5.5 结果分析	122
5.5.1 口语量表熟悉阶段结果分析	122
5.5.2 考试内容分析阶段结果分析	125
5.5.3 标准设定结果分析	126
5.5.4 分界分数的确定	128
5.6 结语	131
第六章 口语量表的后果效度验证——口语量表在教学中的应用	133
6.1 量表和教学	133
6.2 研究问题和步骤	136
6.3 数据收集	136
6.3.1 学生口语量表自评阶段	136
6.3.2 设计基于口语量表的课堂活动	138
6.3.3 课堂活动的学生自评、互评和教师评价	140
6.4 结果分析	143
6.4.1 学生口语量表自评阶段	143
6.4.2 设计基于口语量表的课堂活动	148
6.4.3 课堂活动的学生自评、互评和教师评价	153
6.5 结语	160

第七章 结论	162
7.1 研究总结	162
7.2 研究的理论和实践价值	164
7.3 方法论和创新之处	165
7.4 研究局限和未来展望	166
参考文献	168
附录	201
附录一 口语量表对接 CET-SET4 工作手册	201
附录二 口语量表教学应用手册	257
附录三 学生英语口语自评表	273
附录四 学生课堂活动课前指导材料	277
附录五 口头演讲评分表	282
附录六 学生自评互评反馈表	284

表 目

表 2-1	量表效度的分类和研究方法	23
表 3-1	口语量表效度论证的设计思路和效验内涵	37
表 3-2	数据来源和处理	41
表 3-3	对接研究标准设定法	48
表 4-1	分级验证中各级各类口语描述语数量	53
表 4-2	描述语分级验证评分标准	55
表 4-3	全样本因子分析方法：主因子	60
表 4-4	全样本因子载荷得分	61
表 4-5	全样本主因子分析方法：主因子 正交最大方差	62
表 4-6	因子载荷得分：正交最大方差	62
表 4-7	全样本主成分因子分析方法：主成分因子	63
表 4-8	全样本因子载荷得分：主成分因子	63
表 4-9	全样本主成分因子分析方法：主成分因子 正交最大方差	64
表 4-10	因子载荷得分：正交最大方差.....	64
表 4-11	全样本迭代主因子分析方法：迭代主因子法	65
表 4-12	因子载荷得分：迭代主因子.....	65
表 4-13	全样本迭代主因子分析方法：迭代主因子.....	66
表 4-14	因子载荷得分：迭代主因子.....	66
表 4-15	全样本迭代主因子分析方法：迭代主因子 最大正交旋转	67

表 4-16	全样本迭代主因子分析方法：因子载荷得分 最大正交旋转	67
表 4-17	全样本最大似然因子分析方法：最大似然	68
表 4-18	因子载荷得分：最大似然	69
表 4-19	全样本最大似然因子分析方法：最大正交旋转	69
表 4-20	因子载荷得分：最大正交旋转	70
表 4-21	不同因子分析方法估计结论及比较	70
表 4-22	五级 10 套问卷因子分析	75
表 4-23	六级 14 套问卷因子分析	75
表 4-24	五级 10 套问卷因子载荷得分和高载荷分类	77
表 4-25	六级 14 套问卷因子载荷得分和高载荷分类	78
表 4-26	B2-1 因子结构模型估计——最简化条件法对比最小平均 偏相关法.....	84
表 4-27	B2-1 因子分析结构方程模型估计	85
表 4-28	B3-14 因子结构模型估计——最简化条件法对比最小平均 偏相关法.....	91
表 4-29	因子分析结构方程模型估计 B3-14	92
表 4-30	口语量表的描述等级	95
表 4-31	教师评学生和学生自评数据模型拟合结果	96
表 4-32	剔除的教师评学生不拟合模型的描述语	97
表 4-33	描述语在 Rasch 等值后的难度值描述统计	98
表 4-34	验证前后级别相差 3 级或 4 级的描述语示例	99
表 4-35	五、六级问卷教师样本数量的省市自治区分布	102
表 4-36	五、六级问卷调查的教师数和所评学生数	102
表 4-37	五级和六级 24 套平行问卷的 cronbach 系数	103
表 4-38	B2-1 IRT-GRM 模型项目参数估计	105
表 4-39	B2-2 IRT-GRM 模型项目参数估计	106
表 4-40	五~六级问卷的等值转换常数	108
表 4-41	问卷参数估计以及转换后的参数估计	110

表 4-42	五~六级口语描述语难度值相关度	112
表 5-1	四级口语话题	119
表 5-2	抽取的 CET-SET4 考试话题	121
表 5-3	各话题 CET-SET 4 考生分数统计	121
表 5-4	CET-SET 4 各级别对应的分数段	122
表 5-5	抽取 30 名考生的分数分布信息	122
表 5-6	三级至八级口语典型特征	122
表 5-7	专家准确性判断统计	124
表 5-8	一致性检验——Cohen's Kappa 系数	125
表 5-9	对照口语量表的 CET-SET4 考试任务分析	126
表 5-10	中点分析法	130
表 5-11	对接结果	131
表 6-1	试测后剔除的描述语	137
表 6-2	学生自评问卷 36 条描述语分布表	138
表 6-3	基于口语量表的教学设计(示例)	139
表 6-4	有疑问的描述语列表	144
表 6-5	与课堂演示活动相关的口语能力描述语	148
表 6-6	教师的教学设计	149
表 6-7	班级学生 CET-SET4 考试成绩分布	153
表 6-8	班级学生 CET4 和学期末成绩分布	153
表 6-9	八组口头演讲评分均值	154
表 6-10	学生自评、互评 Cronbach Alpha 信度	154
表 6-11	学生自评、互评组间 t 检验	155
表 6-12	学生互评的三个评分维度和各外部考试的相关分析	156
表 6-13	学生自评的三个评分维度和各外部考试的相关分析	157
表 6-14	教师评分和学生互评 pearson 相关分析	158
表 6-15	教师评分和学生自评 pearson 相关分析	158
表 6-16	学生自评互评反馈	159



图 目

图 2-1	《欧框》的语言能力描述体系	11
图 2-2	《欧框》交际策略分类参数	11
图 3-1	研究操作设计	39
图 4-1	《量表》分级验证问卷回收统计	52
图 4-2	分级验证等值设计	56
图 4-3	口语量表框架	58
图 4-4	全样本因子分析	72
图 4-5	B2-1 分技能秩相关分析	82
图 4-6	B2-1 最大正交旋转及加权最小二乘法最大正交旋转	82
图 4-7	B2-1 分层因子结构模型和双层次因子结构模型	83
图 4-8	B2-1 验证性因子分析隐变量结构关系（标准化后）	86
图 4-9	B2-1 验证性因子分析协方差结构关系（标准化后）	86
图 4-10	验证性因子分析协方差结构关系（标准化后）	87
图 4-11	B3-14 分技能秩相关分析	88
图 4-12	B3-14 最大正交旋转和加权最小二乘法旋转得分	89
图 4-13	B3-14 分层次因子结构模型和双层次因子结构模型	90
图 4-14	验证性因子分析隐变量结构关系（标准化后）	94
图 4-15	验证性因子分析协方差结构关系（标准化后）	94
图 4-16	验证性因子分析协方差结构关系（标准化后）	95
图 4-17	描述语 logit 值和 CTT 难度对比的分布情况	99
图 4-18	五级和六级 24 套问卷的 IRT 等值设计	101

图 4-19	B2-1 锚题的 ICC 曲线和 CCC 曲线·····	106
图 4-20	B3-1 锚题的 ICC 曲线和 CCC 曲线·····	107
图 4-21	连接后各套问卷项目的两两测试特征曲线（以五级为例） ·····	110
图 5-1	专家评分总层面图 ·····	127
图 5-2	两种计算方法计算的分界分数 ·····	129
图 6-1	成功实施口语量表需要兼顾的几个方面 ·····	138
图 6-2	基于口语量表的课堂活动设计实施方案 ·····	141
图 6-3	学生自评描述语总层面图 ·····	144
图 6-4	描述语调查难度分布图 ·····	145
图 6-5	学生基于 CET-SET4 和自评洛基量尺的量表级别分布 ·····	147

第一章 导论

1.1 研究背景

1.1.1 《中国英语能力等级量表》

语言力量表相关的研究已经有 60 多年的历史，主要集中在北美、欧洲、澳大利亚等国家和地区。其中，《欧洲语言共同参考框架：学习、教学、评估》(*Common European Framework of Reference for Languages: Learning, Teaching, Assessment*) (简称《欧框》) 是很有代表性和影响力的语言力量表，它为语言能力的描述提供一个等级标准框架，同时使用者可以根据自身需要进一步细分等级，具有一定的灵活性 (邹申等, 2015)。《欧框》的研制吸收了当代语言学、语言教学和语言习得等领域的研究成果。自《欧框》问世以来，我国学者便一直关注和研究《欧框》，相关研究包括对《欧框》的解读和介绍 (韩宝成, 2006; 张建琴、邹为诚, 2010; 付桂芳、Broeder, 2011) 以及分析、对比和应用研究 (方绪军, 2007; 白乐桑、张丽, 2008; 岑海兵、邹为诚, 2011; 黄婷、贾国栋, 2012; 刘静观, 2012)。同时，我国学者开始呼吁建立中国的语言能力标准，并提出诸多的建设性建议 (杨惠中、桂诗春, 2007; 杨惠中等, 2012; Jin & Wu, 2014)。

2014 年 9 月，国务院印发了《关于深化考试招生制度改革的实施意见》(下称《意见》)，明确提出要加强我国外语能力测评体系建设，

标志着新一轮考试招生制度改革全面启动。这是我国教育领域一项具有基础性、全局性意义的重大制度改革。为了贯彻实施《意见》，教育部于2014年10月启动“国家外语能力测评体系建设”项目。作为该项目的一项基础性工作，《中国英语能力等级量表》（下称《量表》）的研发也同时启动。目前，《量表》已由教育部和国家语言文字工作委员会发布，并于2018年6月1日正式实施（教育部考试中心，2018）。截至目前，已有不少研究从不同的方面介绍了《量表》建设的过程和成果（刘建达，2015，2019a；刘建达、吴莎，2019；金艳、揭薇，2017；朱正才，2015，2016；邹申等，2015）。

《量表》由六部分组成，分别描述语言理解能力、语言表达能力、语用能力、中介能力、语言知识和语言策略。语言理解能力包括理解口头语言信息的能力（听力能力）和理解书面语言信息的能力（阅读能力）；语言表达能力包括用口头语言表达信息的能力（口语能力）和用书面语言表达信息的能力（写作能力）；语用能力分为语用理解能力和语用表达能力；此外，《量表》还描述了中介能力，即口译和笔译能力。各类语言能力都是围绕描述、叙述、说明、论述、指示、互动六项功能所需的语言能力进行描述。《量表》的语言知识框架包括语言结构知识（即语法知识、篇章知识）和语用知识。《量表》的语言策略指的是语言使用者完成一项语言交际活动所采取的有组织、有计划、目标明确的行动步骤，具体包括规划、执行、评估、补救四个步骤（刘建达、韩宝成，2018；刘建达、吴莎，2019）。

1.1.2 《中国英语能力等级量表》之口语量表

《量表》之口语力量表（下称口语量表）从口语能力、口头交际策略和口语文本特征等多维度描述各等级英语口语能力的典型特征。其研发的总体路径是基于中国英语学习者口语交际能力理论模型，在交互理论的指导下界定口语能力构念，采用文献法、采样法以及专家论证等方法撰写和改编描述语，建立口语能力描述语库并划分等级，采集全国大样本教师和学生评价数据为等级划分获取依据，提出具有区别性特征的口语交际能力等级描述。

已发布的口语量表由两部分组成，分别描述口头表达能力和口头表达策略，包括口头表达能力总表、口头表达策略、口头表达能力自我评价量表等 12 张量表，分为基础、提高、熟练三个阶段，设立九个能力等级，总共 379 条描述语。

口语量表将为我国英语口语的教学、测评和学习提供一个统一的参考标准。口语量表对教学的指导作用体现在进一步明确英语口语教学的目标，有利于教学的组织和衔接，有助于教师把教学划分为较小的、具有明确目标的阶段，为各阶段选择合适的教学方法，更有效地组织课堂教学。同时，量表也将有利于教材的编写和开发，提高教育系统的透明度。

口语量表对测评的指导作用体现在为英语口语考试的研发提供科学、系统、清晰的评价指标体系，有助于教师设计检测学习成果的评估办法和工具，提高口语能力测试和评价的质量，更好地掌握和了解学生的实际口语水平。口语量表的实施将使英语口语测试的分数解释更加清晰，分数使用更加公平。考试机构可以通过考试与量表的衔接实现不同考试之间的可比性，增强考试的透明度，以便于考试用户更好地理解和使用考试结果，有利于教育主管部门对不同地区、不同学校的口语教学进行更为科学的评估，还将有助于我国的英语口语测评与国际标准接轨。

口语量表对学习的指导作用体现在帮助学习者了解自己的知识掌握程度和技能达到的水平，跟踪自己的学习过程，开展自我评估，从而进一步明确学习目标，思考和调整自己的学习策略。此外，口语量表有助于学习者根据学习目标，挑选合适的学习材料，确定适合自己的学习内容、任务和测评，提高学习效率。

1.1.3 口语量表的效度研究

量表的用途广泛，因而其效度尤为重要。语言力量表只有经过效度论证，才能证明其有效性和可操作性（杨惠中等，2012）。效度是评价量表质量的重要标志。因此，口语量表的效度研究贯穿研发的始终，且在口语量表发布后，依然需要不断搜集证据，证明口语量表在

实际应用中的效度。

纵观教育测量和语言测试领域，效度理论的发展脉络大致经历了分类论、整体论和论证论三个阶段（Chapelle, 1999；韩宝成、罗凯洲, 2013）。效度理论的发展推动和影响了考试和量表的开发和效度验证模式，特别是效度论证论，研究者将自己的研究结果解释为在论证的基础上支持特定假设，有助于研究者明晰研究方法，为不同的研究主题提供共同的以论证为基础的效度研究元素。效度论证论除了广泛地用于考试的效度验证研究中，还被用于评分量表的开发和验证研究，如 Knoch（2007）、Knoch & Chapelle（2017）提出的基于论证的表现性测试评分量表的效度验证框架。

力量表的效度有别于评分量表以及考试的效度，三者本身是不同的研究对象。因此，口语量表的效度验证有别于口语考试和评分量表的效度验证。为保证口语量表的稳定性、有效性和生命力，其效度验证需要建立一个科学的、完整的、有针对性的验证框架，不能完全沿用考试的效度验证框架。

目前，国内外关于力量表效度验证的研究较少，且缺乏系统连贯的理论框架指导，即整合各个方面的经验证据或理论依据的力量表效度论证框架。存在这样的现状有诸多原因：一方面，力量表的效度研究不如测试的效度研究那样悠久，世界各地对于力量表的关注度也是近年才逐渐升温；另一方面，力量表开发者比较注重量表的研制和开发效度，对力量表效度的检验往往更关注内容效度和结构效度，很少有研究将量表的效度延伸或拓展到量表的外部效度和后果效度。量表的外部效度是指量表的外延性，量表的后果效度检验指通过运用量表进行效度检验，即量表开发基本完成以后进行的后验。此外，量表的开发者和量表的使用者往往从不同的立场和角度看待量表，从量表的开发到使用，效度验证存在断层，量表的开发者不一定是量表的使用者。如 ACTFL 和《欧框》之类影响深远的语言力量表，在其开发伊始和投入使用后的很长一段时间，学界对于其质量和应用研究的关注度一直很高，推崇者有之，批判者有之。即便如此，

目前还没有研究系统地开展这些量表的效度检验。

能力量表的效度验证框架需要综合各项指标,以整合取自各个方面的经验证据或理论依据(朱正才,2016)。因此,口语量表的效度验证是复杂和多维的,需要在效度验证框架的指导下,长期地、动态地开展,不断收集和积累各个方面的效度证据,通过具体的主张和理据把效度研究框架中各个重要属性有机地联系起来,并明确各个主张和理据的证据来源,从而有助于口语量表核心成分的解释,并探讨口语量表使用的后效。

1.1.4 口语量表的效度取证

《教育与心理测量标准》(*Standards for Educational and Psychological Testing*)(简称《标准》)指出,合理的考试效度论证是要将各种证据整合成连贯的体系,论证的可靠性取决于证据和理论在多大程度上支持考试成绩的预期释义及用途(AERA, APA & NCME, 2014: 21),可以说效度研究就是收集尽可能多的证据去支持考试(谢小庆,2013)。虽然《标准》源自美国,但其影响力已经远超北美各国(Zumbo, 2014)。在开展语言测试的效度验证时,研究者需要依据效度论证框架采集证据和报告研究结果,如果这一技术问题没有解决,那么效度验证最后可能就会流于核对清单式(checklist)验证(Fulcher, 2015: 116-120),而不是真正的效度论证(Chapelle *et al.*, 2010)。

口语量表的效度需要来自理论和实践各个角度的证据支撑:第一、来自口语量表内部的证据,如口语量表的构念和级别效度;第二、口语量表的外部效度证据,如量表与考试之间关系的证据;第三、口语量表的后果效度证据,口语量表的使用和影响,如在考试或教学中的应用。口语量表开发过程中,内部效度证据是非常重要的。North & Schneider(1998)、De Jong(1990)认为,基于专家的经验 and 直觉或者借鉴现有的量表而开发的语言能力量表在实践中会产生各种问题,特别是在高风险的环境中使用这样的能力量表。Fulcher(1996)回顾了ACTFL量表效度的研究,指出如果量表在开发伊始缺乏扎实的实证基础,就很难在其投入使用后对其进行效度研究,研究结果可能

与量表的假设及构念完全背离。量表初步构建后，研究者需要采集基于量表内部结构的证据，可以使用统计方法如因子分析或者结构方程建模分析量表的结构，检验量表等级描述语的可量测性（scalability）（Henning，1992）。量表发布后，研究者要进行科学的研究设计和数据分析，收集量表的外部效度和后果效度证据，为使用和改进量表提供依据。

总之，口语量表的效度取证是效度验证的关键，必须有效地收集和正确地使用口语量表效度证据。同时，国外有关能力量表的相关研究也可以为口语量表效度研究提供参考。

1.2 研究目的和内容

本研究是口语量表研制的一个组成部分，目的是构建口语量表的效度验证框架，并应用该框架综合分析口语量表的效度。本研究的内容包括以下三个方面：

第一，理论构建。本研究旨在从效度理论的最新视角构建量表的效度验证框架并将该框架应用于口语量表的效度研究，探索口语量表效度验证的内涵，发展和丰富语言能力量表的效度理论研究。

第二，实证研究。本研究将根据所构建的效度验证框架，收集相应证据，分析各方面的证据，从而将多个效度研究视角有机地整合到口语量表的效度论证中。此外，本研究尝试使用不同的研究方法，使不同效度证据的相互关系通过实证检验得到解释。

第三，应用探索。本研究将探索如何在实际应用中进一步论证量表的效度，从而将口语量表的效度研究范畴拓展到后期的实际应用阶段。

在研究的总体设计上，本研究采用了定量、定性和经验方法，既基于较为成熟的语言测试效度理论并吸收语言测试效度验证的长处，同时结合了本研究的需求和现状，利用现代测量学技术对口语量表进行系统、深入的效度检验。