

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

فرم گردآوری اطلاعات پایان نامه ها
کتابخانه مرکزی دانشگاه علامه طباطبایی

عنوان: بررسی تاثیر برنامه آموزشی مستقیم (رو در رو) بر سنجش مهارت نگارش در زبان دوم	
نویسنده / محقق: هومن بیژنی	
مترجم:	
استاد راهنما: سرکار خانم دکتر معرفت الدین / جناب آقای دکتر فهیم	استاد مشاور / استاد داور: جناب آقای دکتر تاج
کتابنامه: دارد	واژه نامه: ندارد
نوع پایان نامه: <input type="checkbox"/> بنیادی <input type="checkbox"/> توسعه ای <input checked="" type="checkbox"/> کاربردی	
مقطع تحصیلی: کارشناسی ارشد	سال تحصیلی: ۱۳۸۷-۸۸
محل تحصیل: تهران	نام دانشگاه: علامه طباطبایی دانشکده: ادبیات فارسی و زبانهای خارجی
تعداد صفحات: ۱۸۰	گروه آموزشی: آموزش زبان انگلیسی
کلید واژه ها به زبان فارسی: ارزشیابی، آموزش درجه بندی کننده، پایایی، پایایی ارزیاب ها، اعتبار	
کلید واژه ها به زبان انگلیسی: Assessment, Rater training, Reliability, Inter-rater reliability, Validity	



Allameh Tabataba'i University
Faculty of Persian Literature and Foreign Languages

Evaluating the Effectiveness of a Face-to-Face Training Program on L2 Writing Assessment

A Thesis Submitted in Partial Fulfillment of the Requirements for
the Degree of Masters of Arts in Teaching English as a Foreign
Language (TEFL)

Thesis Advisor:
Dr. Fahimeh Marefat

Thesis Reader:
Dr. Zia Tajeddin

By:
Houman Bizhani

5iv 1388 | July 2009

*We hereby approve that this thesis
By*

Houman Bizhani

Entitled “Evaluating the effectiveness of a face-to-face training program on L2 writing assessment” be accepted in partial fulfillment of the requirements for the degree of Masters of Arts in TEFL.

Committee on the Oral Examination:

Advisor: Dr. Fahimeh Marefat

.....

Reader: Dr. Zia Tajeddin

.....

Examiner: Dr. Mansoor Fahim

.....

Head of the Department of English Language and Literature

Dr. Zia Tajeddin

.....

July 2009

Dedicated to:

*My dear mother
who encouraged and inspired me with
affection*

ACKNOWLEDGMENTS

As I look back on the process of planning, running and writing the research for this thesis, I become aware how fortunate I was in having had the opportunity to work with so many intelligent, generous, and helpful people.

I would like to thank, first of all, Dr. Fahimeh Marefat, my advisor, who had been a wonderful supporter of my research and without her sincere help, approval, and influence, this study could have never been penned.

Moreover I would like to express my deepest gratitude to Dr. Zia Tajeddin, my reader, for his generous contribution, assistance, supports, patience, and constructive comments in this study.

Special thanks go to Dr. Sara Cushing Weigle, UCLA faculty member, for her sincere help and for providing me with the literature in this study.

I wish to acknowledge Dr. Ute Knoch, University of Auckland faculty member, whose great help in providing me with the helpful articles in this research encouraged me a lot.

I am also greatly in debt to Dr. Asadi from Tarbiat Moallem University, whose honest and great help in training raters helped me run this research.

Perhaps my greatest debt is to the 12 raters who gave so generously all their time in rating so many compositions, not once but three times.

Most especially I would like to thank my great mother who has been a constant support and believer in me and has shown incredible patience and good cheer during this study.

Abstract

Rater training is an essential component of any reliable writing assessment in the first or second language studies. However, little is known about the processes by which raters learn to apply the appropriate criteria in making judgments about writing samples. This study examined the process through both qualitative and quantitative approaches. Twelve raters, six inexperienced and six experienced, rated 15 essays on a topic before training, 15 additional essays on the same topic immediately after training, and another 15 essays on the same topic after training using a four-part scoring rubric (The IELTS scale) covering Organization, Structure, Vocabulary, and Punctuation.

The major findings, after the ratings were analyzed, were as follows: (1) Inexperienced raters tended to be less consistent than experienced raters before training but not afterwards. (2) Training was effective and moved all the raters towards having more consistency less severity and less bias in rating. (3) Through time, training loses its effectiveness and raters should be retrained in intervals. (4) Despite training, significant differences in severity among raters remained.

The qualitative data analysis along with the help of quantitative analyses showed that inexperienced raters were more willing to accept

authorities' comments in rating and their ratings were more improved than experienced raters because experienced raters are less likely to accept authorities' comments. The results demonstrated that training helps raters understand the crucial criteria to be used in rating essay compositions. On the other hand the results showed that all raters are capable of rating essay compositions regardless of the background or any type of experience in this area; therefore, what makes difference in rating essay compositions is a good training program. Another implication of this study is that decision makers should not charge a huge amount of money in using experienced raters for rating essay compositions, because experience does not increase rating reliability. In fact they should use inexperienced raters but with a good training program in advance.

Dedication	iv
Acknowledgements	v
Abstract	vii
Table of Contents	ix
List of Tables	xiii
List of Figures	xiv
Chapter One: Introduction	1
1.1. Overview of the Study	2
1.2. Statement of the Problem	6
1.3. Significance of the Study	7
1.4. Purpose of the Study	8
1.5. Research Questions	9
1.6. Research Hypotheses	10
1.7. Limitations and Delimitations of the Study	11
1.7.1. Limitations	11
1.7.2. Delimitations	11
1.8. Definition of the Key Terms	12
Chapter Two: Review of Literature	14
2.1. Writing in First and Second Language	15
2.2. Writing as A Social and Cultural Phenomenon	18
2.2.1. Social Aspects of Writing	18
2.2.2. Cultural Aspects of Writing	18
2.3. Writing Assessment	19
2.4. Frameworks for Writing Assessment	22
2.5. Usefulness of Writing Tests	25

2.5.1. Reliability	25
2.5.1.1. Background	26
2.5.1.2. Rating Scale	26
2.5.1.3. Training	26
2.5.2. Construct Validity	27
2.5.3. Authenticity	27
2.5.4. Interactiveness	27
2.5.5. Impact	28
2.5.6. Practicality	28
2.6. Direct vs. Indirect Writing Tests	29
2.6.1. Direct Writing Tests	29
2.6.2. Indirect Writing Tests	29
2.7. Sources of Rater Disagreement	31
2.7.1. Text Variables	32
2.7.2. Rater Variables	33
2.7.3. Rating Scales	35
2.7.4. Context Variables	36
2.7.5. Test-Taker Variables	37
2.8. Writing Instructions	38
2.9. Rating Scales	39
2.9.1. Primary Trait Scoring	40
2.9.2. Holistic Scoring	41
2.9.3. Analytic Scoring	42
2.10. The Scoring Process	44
2.10.1. Procedures for Assuring Reliability	45
2.10.2. Assessing Reliability of Scores	46
2.10.3. Assessing Validity of Scoring Procedure	47
2.11. Rater Training	48
2.12. Rater Variation and Training Effect	52

2.13. Rater Severity and Leniency	53
2.14. Rater’s Verbal Protocols	55
2.15. Computer Mediated and Online Rater Training	56
2.16. Summary	58
Chapter Three: Method	69
3.1. Overview	60
3.2. Design	61
3.3. Participants	61
3.3.1. Students	61
3.3.2. Raters	62
3.3.3. Rater Trainer	64
3.4. Instruments	65
3.4.1. Students’ Compositions	65
3.4.2. Pre-training Questionnaire	66
3.4.3. Post-training Questionnaire	68
3.4.4. Rating Scale	69
3.4.5. Instruction to the Students	70
3.4.6. Instruction to the Raters	70
3.5. Procedures	71
3.5.1. Phase One: Pre-training data collection	71
3.5.2. Phase Two: Data Collection During Training Session	73
3.5.3. Phase Three: Immediate Post Training Data Collection	75
3.5.4. Phase Four: Delayed Post Training Data Collection	77
Chapter Four: Results and Discussion	79
4.1. Results	81
4.1.1. Research Question 1	81
4.1.2. Research Question 2	93

4.1.3. Research Question 3	102
4.1.4. Research Question 4	109
4.1.5. Research Question 5	130
4.1.6. Research Question 6	137
4.2. Summary of Findings	142

Chapter Five: Conclusion, Implications, and Suggestions for Further Research 145

5.1. Conclusion	146
5.2. Implications of the Study	147
5.3. Suggestions for Further Research	149
References	152
Appendix 1: Pre-training Questionnaire	160
Appendix 2: Post-training Questionnaire	163
Appendix 3: IELTS Rating Scale	166
Appendix 4: Students' Writing Instruction	169
Appendix 5: Rater's Rating Instruction	170
Appendix 6: z-value for Each Discrepancy Score of Each Trait (Pre)	171
Appendix 7: z-value for Each Discrepancy Score of Each Trait (Post)	176

List of Tables

Table 2.1	Types of rating scales used for writing assessment	40
Table 3.1	Raters' characteristics	64
Table 3.2	Table of specification for raters (Pre-training)	67
Table 3.3	Table of specification for raters (Post-training)	69
Table 3.4	Summary of data collection and research procedures	78
Table 4.1	Fixed-choice responses to questions on pre-training	82
Table 4.2	Fixed-choice responses to questions on post-training	88
Table 4.3	Pearson's correlation coefficient among raters (Pre)	95
Table 4.4	Standard deviation among raters before training	97
Table 4.5	Pearson's correlation coefficient (Post Immediate)	99
Table 4.6	Standard deviation among raters after training	101
Table 4.7	Raters' degree of severity or leniency before training	105
Table 4.8	Raters' degree of severity or leniency after training	107
Table 4.9	Raters' biasness in each trait prior to training	112
Table 4.10	Raters' biasness in each trait following training	114
Table 4.11	Raters' dispersion index in rating scale traits (Pre/Post)	117
Table 4.12	Interrater reliability among NEW/OLD raters (Pre/Post)	134
Table 4.13	NEW/OLD raters' dispersion before and after training	135
Table 4.14	Sum of point differences	135
Table 4.15	Pearson's correlation coefficient (Post-Delayed)	139
Table 4.16	Standard deviation among raters 2 months after training	140
Table 4.17	Raters' change of behavior in the 3 phases of the study	141

List of Figures

Figure 2.1	Ruth and Murphy's model of writing assessment (1988)	23
Figure 2.2	Freedman Calfee's model of writing assessment (1983)	25
Figure 2.3	Factors in writing assessment (McNamara, 1996)	31
Figure 4.1	Raters' change of behavior before and after training	109
Figure 4.2	Bias analysis: Rater NEW1	118
Figure 4.3	Bias analysis: Rater NEW2	119
Figure 4.4	Bias analysis: Rater NEW3	120
Figure 4.5	Bias analysis: Rater NEW4	121
Figure 4.6	Bias analysis: Rater NEW5	122
Figure 4.7	Bias analysis: Rater NEW6	123
Figure 4.8	Bias analysis: Rater OLD1	124
Figure 4.9	Bias analysis: Rater OLD2	125
Figure 4.10	Bias analysis: Rater OLD3	126
Figure 4.11	Bias analysis: Rater OLD4	127
Figure 4.12	Bias analysis: Rater OLD5	128
Figure 4.13	Bias analysis: Rater OLD6	129

الف: موضوع و طرح مساله (اهمیت موضوع و هدف): با وجود آنکه اهمیت آموزش در ارزیابی مهارت نگارش مورد بحث قرار گرفته اما تحقیقات ناچیزی در رابطه با درک ارزیاب ها نسبت به امر برنامه های آموزشی انجام شده است. همچنین تحقیقات بسیار کمی در رابطه با تاثیر آموزش ارزیابی مهارت نگارش بر ارزیاب های باتجربه و بی تجربه انجام شده است. از طرف دیگر تحقیقاتی که تاکنون انجام شده اند اکثرا بر مبنای ارزیابی تاثیر آموزش در کوتاه مدت بوده اند و بنابراین تحقیق کمی در رابطه با تاثیر بلند مدت آموزش در ارزیابی مهارت نگارش انجام شده است. به همین ترتیب، تحقیقات و بررسی های بسیار کمی در رابطه با میزان تاثیر آموزش در کاهش درجه سخت گیری و یا آسن گیری ارزیاب ها و همچنین میزان سو گیری آنها انجام شده است. این رساله تمامی مسایل بالا را مورد بررسی قرار داده است.

ب: مبانی نظری شامل مرور مختصری از منابع، چهارچوب نظری و پرسشها و فرضیه ها: لیناکر (۱۹۸۹) بر این باور است که گوناگونی ارزیاب ها امری اجتناب ناپذیر و بخشی از پروسه ی ارزیابی است. او اعتقاد دارد که ارزیاب ها را نمی توان با استفاده از آموزش به میزان یکسانی از درجه سخت گیری رساند. بنابر این عملکرد آموزش نباید به گونه ای باشد که هدف از آن رساندن ارزیاب ها به تطابق با یکدیگر (پایایی ارزیاب ها) باشد بلکه باید به گونه ای باشد که هدف از آن رساندن ارزیاب ها به تطابق با خودشان (پایایی ارزیاب) باشد. بر این اساس پرسشهای زیر قابل طرح اند. ۱- درک ارزیابها از برنامه رو در رو آموزشی در سنجش مهارت نگارش قبل و بعد آموزش چگونه است؟ ۲- آیا پایایی ارزیابها بلافاصله بعد از برنامه ی رو در رو آموزشی افزایش می یابد؟ ۳- آیا بعد از برنامه آموزشی ارزیابها سخت گیر تر یا آسان گیر تر می شوند؟ ۴- آیا کاهشی در میزان سو گیری ارزیابها بعد از برنامه آموزشی وجود دارد؟ ۵- آیا تفاوتی بین ارزیابهای باتجربه و بی تجربه در بخش مهارت نگارش قبل و بعد از برنامه آموزشی وجود دارد؟ ۶- آیا پایایی ارزیابها در طولانی مدت پس از اتمام برنامه ی آموزشی همچنان افزایش داشته است؟

پ: روش تحقیق شامل تعریف مفاهیم، روش تحقیق، جامعه مورد تحقیق، نمونه گیری و روشهای نمونه گیری، ابزار اندازه گیری، نحوه اجرای آن، شیوه گرد آوری و تجزیه و تحلیل داده ها: به منظور پاسخ به سوالات تحقیق از یک طرح تجربی نمای کمیتی استفاده شده است. در این راستا ۶۰ زبان آموز، ۱۲ ارزیاب، ۲ استاد دانشگاه شرکت داشته اند. همچنین ۴۵ نمونه انشاء زبان آموزان، ۲ پرسشنامه، مقیاس درجه بندی آیلنس، دستور نامه زبان آموزان و ارزیابها مورد استفاده قرار گرفت. در مرحله اول این پژوهش از زبان آموزان آزمون نگارش به عمل آمد و پس از تایپ ۱۵ عدد از آنها به ارزیابها برای ارزیابی داده شد. در مرحله دوم، اولین پرسشنامه به ارزیابها داده شد و سپس برنامه آموزشی اجرا شد. در مرحله سوم پرسشنامه دوم به ارزیابها داده شد و سپس ۱۵ برگه نگارش دیگر به ارزیابها داده شد تا ارزیابی کنند. در مرحله چهارم که ۲ ماه بعد از پایان برنامه آموزشی اجرا شد، ۱۵ برگه ی نگارش دیگر به ارزیابها داده شد تا ارزیابی کنند. در این پژوهش از روشهای آماری سنجش میزان همبستگی، سنجش میزان ضریب پایایی و آزمون تی استفاده شد.

ت: یافته های تحقیق: بعد از آموزش تمامی ارزیاب ها تا حد بسیار بالایی به مطابقت با یکدیگر رسیدند. آموزش همچنین درجه سخت گیری و آسان گیری ارزیاب ها و همچنین میزان سو گیری آنها را تا حد بسیار زیادی کاهش داد. همچنین آموزش بر ارزیاب های بی تجربه بسیار موثر تر از ارزیاب های با تجربه بوده و نهایتا اینکه گذر زمان باعث کمرنگ شدن اثر آموزش بر ارزیاب ها می شود.

ث: نتیجه گیری و پیشنهادات: این پژوهش نهایتا به افزایش میزان پایایی بین ارزیاب ها و کاهش میزان سخت گیری و یا آسان گیری آنها انجامید. این پژوهش نشان می دهد که استفاده صرف از مقیاس باعث افزایش میزان ضریب پایایی نمی شود بلکه آنچه که میزان پایایی را افزایش می دهد آموزش مناسب است. این پژوهش نشان داد که تجربه و زمینه باعث افزایش میزان ضریب پایایی نمی شود بنابراین سیاستگذاران آموزشی بجای صرف هزینه برای استفاده از ارزیاب های باتجربه باید از ارزیاب های بی تجربه استفاده کنند و در عین حال از برنامه های آموزشی بهره گیرند. همچنین اینکه آموزش باید به طور تناوبی برای ارزیاب ها تجدید شود. بر این اساس پیش نهاداتی برای تحقیق های بیشتر مطرح می شود که عبارتند از: استفاده از پیش نویس های کلامی بدست آمده از ارزیاب ها برای بررسی ذهنیت آنها در ارزیابی. بررسی تاثیر ویژگی های شخصیتی ارزیاب ها بر ارزیابی ایشان. استفاده از ارزیاب های بومی و مقایسه آنها با ارزیاب های غیر بومی.

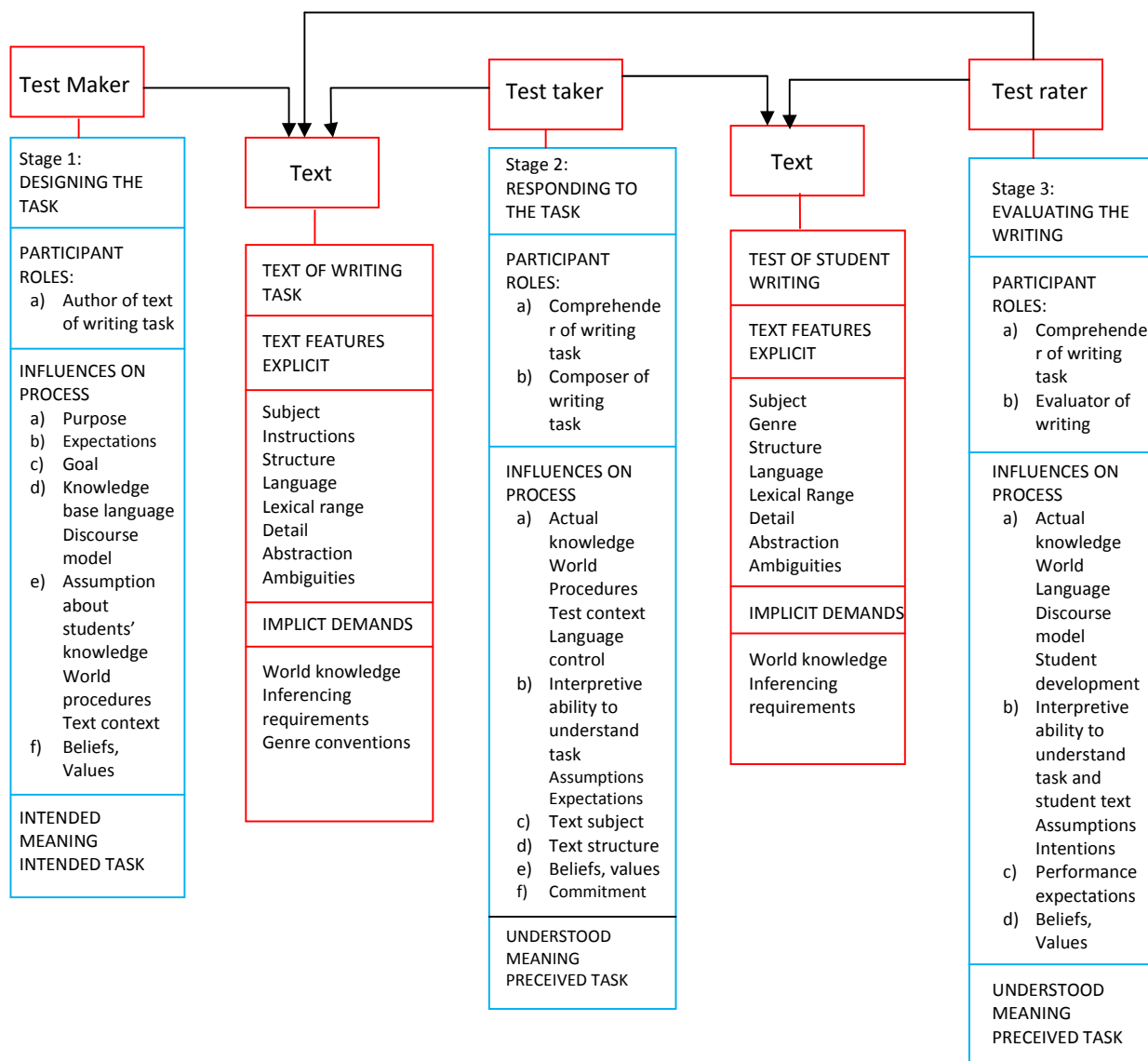
صحت اطلاعات مندرج در این فرم بر اساس محتوای پایان نامه و ضوابط مندرج در فرم را گواهی می نمایم.

استاد راهنما:

سمت علمی:

نام دانشکده:

رییس کتابخانه:



Chapter One

Introduction

1. Introduction

1.1. Overview of the study

The ability to write effectively is becoming increasingly important in our global community, and instruction in writing is thus assuming an increasing role in both second and foreign language education. As a result, the ability to write in a second language is widely recognized as an important skill for educational, business, and personal reasons (Weigle, 2002). With the significance attached to writing in SL and FL contexts, testing writing is considered an important issue. As the role of writing in second language learning increases, there is an even greater demand for valid and reliable ways to test this ability.

In recent years, essay examinations have become a standard practice in assessing the writing skills of both first and second languages. Because such tests require subjective evaluations of writing quality, a great deal of research emphasis has been placed on achieving an acceptable level of interrater reliability in order to show that essays can be scored as fairly and constantly as possible. However, this emphasis on reliability has been at the expense of decreasing test validity (Charney, 1984); that is, the procedure for achieving higher reliability may not lead to valid judgments of writing quality. One issue which is at the heart of both reliability and validity in essay scoring is that of rater training.

In the past 30 years, holistic writing assessment has become the norm in evaluating the writing skills in both first and second language. In holistic assessment, examinees are asked to write compositions on one or more topics, and then they are scored by raters. Because these tests are scored subjectively, it is essential that raters be carefully trained to conform to some standards (Weigle, 1994b). These standards are given to raters through scoring rubrics that describe the typical characteristics of writing samples at different levels. It is evident that without training, ratings tend to be highly unreliable. A large body of literature, beginning with the work of Diederich, French, and Carlton (1998) and continuing through the present day with the work of Elder,