IN THE NAME OF GOD


THE COMPASSIONATE


AND


THE MERCIFUL

بسمه تعالی

**تاییدیه اعضای هیات داوران حاضر در جلسه دفاع از رساله دکتری**

خانم / آقای سارا جلالی            رساله            واحدی خود را با عنوان:  مقایسه نظری و عملی دیدگاه‌های روان‌سنجی کلاسیک و سوال-پاسخ و طراحی مدل آزمون انطباقی رایانه‌ای براساس نظریه سـوال-پاسخ

در تاریخ            ارائه کردند.

اعضای هیات داوران نسخه نهایی این رساله را از نظر فرم و محتوا تایید کرده اسـت و پـذیرش آن‌را بـرای تکمیـل درجه دکتری آموزش زبان انگلیسی            پیشنهاد می‌کنند.

| امضاء | رتبه علمی | نام و نام خانوادگی | اعضای هیات داوران |
|---|---|---|---|
|  | دانشیار | غلامرضا کیانی | 1- استاد راهنمای اصلی |
|  | – | – | 2- استاد راهنمای دوم |
|  | استادیار | ولی‌الله فرزاد | 3- استاد مشاور اول |
|  | – | – | 4- استاد مشاور دوم |
|  | استادیار | سیده سوسن مرندی | 5- استاد ناظر |
|  | دانشیار | سید محمد علوی | 6- استاد ناظر |
|  | استاد | پرویز بیرجندی | 7- استاد ناظر |
|  | استادیار | محمدرضا عنانی سراب | 8- استاد ناظر |
|  | استادیار | رضا غفارثمر | 9- نماینده شورای تحصیلات تکمیلی |

2

## آیین‌نامه حق مالکیت مادي و معنوي در مورد نتايج پژوهشهاي علمي دانشگاه تربيت مدرس

مقدمه: با عنايت به سياستهاي پژوهشي و فناوري دانشگاه در راستاي تحقق عدالت و كرامت انسانها كه لازمه شكوفايي علمي و فني است و رعايت حقوق مادي و معنوي دانشگاه و پژوهشگران، لازم است اعضاي هيأت علمي، دانشجويان، دانش‌آموختگان و ديگر همكاران طرح، در مورد نتايج پژوهشهاي علمي كه تحت عناوين پايان‌نامه، رساله و طرحهاي تحقيقاتي با هماهنگي دانشگاه انجام شده است، موارد زير را رعايت نمايند:

ماده 1- حق نشر و تكثير پايان نامه/ رساله و درآمدهاي حاصل از آنها متعلق به دانشگاه مي باشد ولي حقوق معنوي پديد آورندگان محفوظ خواهد بود.

ماده 2- انتشار مقاله يا مقالات مستخرج از پايان‌نامه/ رساله به صورت چاپ در نشريات علمي و يا ارائه در مجامع علمي بايد به نام دانشگاه بوده و با تاييد استاد راهنماي اصلي، يكي از اساتيد راهنما، مشاور و يا دانشجوي مسئول مكاتبات مقاله باشد. ولي مسئوليت علمي مقاله مستخرج از پايان نامه و رساله به عهده اساتيد راهنما و دانشجو مي باشد.

تبصره: در مقالاتي كه پس از دانش‌آموختگي بصورت تركيبي از اطلاعات جديد و نتايج حاصل از پايان‌نامه/ رساله نيز منتشر مي‌شود نيز بايد نام دانشگاه درج شود.

ماده 3- انتشار كتاب و يا نرم افزار و يا آثار ويژه حاصل از نتايج پايان‌نامه/ رساله و تمامي طرحهاي تحقيقاتي كليه واحدهاي دانشگاه اعم از دانشكده ها، مراكز تحقيقاتي، پژوهشكده ها، پارك علم و فناوري و ديگر واحدها بايد با مجوز كتبي صادره از معاونت پژوهشي دانشگاه و براساس آئين نامه هاي مصوب انجام شود.

ماده 4- ثبت اختراع و تدوين دانش فني و يا ارائه يافته ها در جشنواره‌هاي ملي، منطقه‌اي و بين‌المللي كه حاصل نتايج مستخرج از پايان‌نامه/ رساله و تمامي طرح‌هاي تحقيقاتي دانشگاه بايد با هماهنگي استاد راهنما يا مجري طرح از طريق معاونت پژوهشي دانشگاه انجام گيرد.

ماده 5- اين آيين‌نامه در 5 ماده و يك تبصره در تاريخ87/4/1 در شوراي پژوهشي و در تاريخ 87/4/23 در هيأت رئيسه دانشگاه به تاييد رسيد و در جلسه مورخ 87/7/15 شوراي دانشگاه به تصويب رسيده و از تاريخ تصويب در شوراي دانشگاه لازم‌الاجرا است.

**English Department**
**Faculty of Humanities**
**Tarbiat Modares University**


Thesis
Submitted in Partial Fulfillment of the
Requirements for the Degree of Doctor of Philosophy in
Teaching English as a Foreign Language


**Theoretical and Practical Comparisons of Two Psychometric
Frameworks Classical Test Theory and Item Response Theory,
and Designing the Computer Adaptive Test on the Basis of Item
Response Theory**


**By**
Sara Jalali

**Supervisor**
Dr. Gholam Reza Kiany

**Advisor**
Dr. Valiollah Farzad

**2010**

# DEDICATION

*To the loving memory of my father who was not able to see me complete my doctoral studies.*

*To my mother who supported me through this endeavor.*

*&*

*To my dearest sister, Dr. Pooneh Jalali, for her overwhelming support and encouragement.*

# Acknowledgements

I would like to express my profound gratitude to my supervisor, Dr. Kiany for his interest in the project and for his valuable comments and guidance from the selection of the topic up to the completion of the research. He has been there to answer questions and offer guidance. He has always asked the hard questions, and allowed me to grow as a researcher. I will be eternally grateful to him.

I am also indebted to my advisor, Dr. Farzad who kindly read and commented on the first draft of the dissertation and provided me with helpful guidelines.

I owe to my internal and external readers Dr. Ghafar Samar, Dr. Marandi, Dr. Alavi, Dr. Birjandi and Dr. Anani for their valuable advice. Their comments and discussions with me have made me think and reexamine my work. I should thank them for listening to me and helping guide my research ideas and questions.

I also owe special thanks to the staff of National Organization of Educational Testing (NOET) who helped me a lot in data gathering and handling the software.

# Abstract

There are two major theories of measurement in psychometrics: Classical Test Theory (CTT) and Item-Response Theory (IRT). Despite its widespread and long use, CTT has a number of shortcomings, which make it problematic to be used for practical and theoretical purposes. IRT tries to solve these shortcomings, and provide better and more dependable answers. One of the applications of IRT is the assessment of Differential Item Functioning (DIF). DIF tells the test developer whether the test item functions differently for different groups. Another important use of IRT is in the area of Computer Adaptive Tests (CAT). CAT is based on IRT, and the stepping-stone in preparing a CAT is the preparation of an item bank. Item banking is based on IRT. When IRT is ignored, item banking will not be applicable and consequently there will be no CAT.

The present study first provided a thorough comparison of CTT and IRT from both theoretical and practical perspectives. For this part of the study, the scores of 3000 testees were used. After that, IRT was utilized to estimate DIF between two gender groups and three fields of study i.e. mathematics, science and humanities in the specific English language part of the foreign language university entrance exam questions of the year 2006. For this part, the data of 15486 participants were used for finding gender DIF and the data of 3924 participants for field DIF. Then, IRT was used to prepare an item bank of the specific English language part of the mock foreign language entrance exam questions for the years 2006 and 2007. This mock exam is administered by an institute related to National Organization of Educational Testing (NOET). For preparing the item bank, specific new software i.e. FastTEST, was utilized. Finally, this item bank was utilized for preparing the CAT version of the English exam, which was the final goal of the dissertation.

The findings of this study showed that CTT- and IRT-based person statistics correlated highly across the three IRT models. Also, it was found that item difficulty and item discrimination indexes from CTT correlated highly with those from all IRT models. The DIF analysis showed that there were a number of DIF items in the exam and these items were analyzed in order to find the source of DIF. Finally, a suitable item bank along with the CAT version of the English exam was prepared. The findings of the present study can be of great importance for the educational system. The researcher proposed some suggestions as to the use of IRT and English CAT in Iran.

**Keywords:** Classical Test Theory (CTT); Item-Response Theory (IRT); Computer Adaptive Test (CAT); item bank; Differential Item Functioning (DIF); specific English part, foreign language university entrance exam

**Table of Contents**

| Title | Page |
|-------|------|

# List of abbreviations

| | |
|---|---|
| **ASVAB** | Armed Services Vocational Aptitude Battery |
| **a2PL** | Discrimination Parameter for Two-Parameter Logistic Model |
| **a3PL** | Discrimination Parameter for Three-Parameter Logistic Model |
| **b1PL** | Difficulty Parameter for One-Parameter Logistic Model |
| **b2PL** | Difficulty Parameter for Two-Parameter Logistic Model |
| **b3PL** | Difficulty Parameter for Three-Parameter Logistic Model |
| **BISER** | Biserial Correlation Coefficient |
| **CAST** | Computerized Adaptive Screening Test |
| **CAT** | Computer Adaptive Test |
| **CTT** | Classical Test Theory |
| **df** | degrees of freedom |
| **DIF** | Differential Item Functioning |
| **EFL** | English as a Foreign Language |
| **FCE** | First Certificate in English |
| **ICC** | Item Characteristic Curve |
| **IIF** | Item Information Function |
| **IRF** | Item Response Function |
| **IRT** | Item Response Theory |
| **LR** | Logistic Regression |
| **MH** | Mantel-Haenszel |
| **NOET** | National Organization of Educational Testing |
| **P & P** | Paper-and-Pencil |
| **PBISER** | Coint-Biserial Correlation Coefficient |
| **PCTT** | Proportion Correct based on Classical Test Theory |
| **1PLM/1PL** | One-Parameter Logistic Model |
| **2PLM/2PL** | Two-Parameter Logistic Model |
| **3PLM/3PL** | Three-Parameter Logistic Model |
| **SEE** | Standard Error of Estimation |

**SEM**    Standard Error of Measurement

**SD**     Standard Deviation

**TEFL**   Teaching English as a Foreign Language

**TIF**    Test Information Function

**TOEFL**  Test of English as a Foreign Language

# List of symbols

| Symbol | Explanation |
|--------|-------------|
| **a** | discrimination parameter |
| **b** | difficulty parameter |
| **$b_F$** | item difficulty parameter estimate for females |
| **$b_M$** | item difficulty parameter estimate for males |
| **c** | guessing parameter |
| **D** | discrimination index |
| **e** | constant 2.718 |
| **P** | proportion correct |
| **p-value** | probability of success/significance |
| **$p\,(\theta)$** | probability of a testee's success |
| **q** | proportion of test takers who get the item incorrect |
| **r** | Correlation |
| **$r_{bis}$** | biserial correlation coefficient |
| **$r_{pbi}$** | point-biserial correlation coefficient |
| **$s_x$** | standard deviation of test scores |
| **$s^2_x$** | observed score variance |
| **$s^2_t$** | true score variance |
| **$s^2_e$** | error score variance |
| **$SE_F$** | standard error of estimation for females |
| **$SE_M$** | standard error of estimation for males |
| **SQRT** | square root |
| **t** | t-value |
| **$\theta$** | theta/latent trait |
| **$\theta_T$** | ability estimation based on the whole set of items |
| **$\theta_N$** | ability estimation based on a sub-set consisting of the items with no DIF |
| **x** | observed score |
| **$x_t$** | true score |

$\mathbf{x_e}$    error score

$\overline{x}$    mean score

$\overline{x}_{\mathbf{p}}$    mean score on the test for those who get the item correct

$\overline{x}_{\mathbf{q}}$    mean score on the test for those who get the item incorrect

$\mathbf{z}$    ordinate of normal curve corresponding to p

$\sqrt{}$    square root

$\chi^2$    Chi-square

$\phi$    phi

## List of tables

# List of figures

# CHAPTER I

# INTRODUCTION

## 1-1- Introduction

In this chapter, first, a background to the study is presented. Then, an overview of five important concepts is provided i.e. classical test theory, item response theory, differential item functioning, item banks and computer adaptive tests. In the second section, the problems are stated. The third section describes the significance of this study. The remaining sections of the chapter state the research questions and hypotheses (sections four and five), elaborate on the definitions of pertinent terms (section six), and describe limitations of the study (section seven).

## 1-2- Background

In the theory of measurement, there are two major measurement frameworks: classical test theory (CTT) and item response theory (IRT). Differences are most evident in the statistical analysis underlying each theory.

### 1-2-1- Classical test theory (CTT)

Classical test theory (CTT) is best suited for traditional testing situations, either in group or in individual settings, in which all the members of a target population, e.g. persons seeking college admission, are administered the same or parallel sets of test items. CTT has four underlying assumptions (Bachman, 1990):