

Detecting Gender DIF with an English Proficiency Test in EFL Context

Seyed Mohammad Reza Amirian^۱, Seyed Mohammad Alavi^۲, Angel M. Fidalgo^۳

Abstract

The aim of the present study is twofold. First, the paper investigated whether University of Tehran English Proficiency Test (UTEPT) manifested substantial gender Differential Item Functioning (DIF). Second, the flagged DIF items were subjected to a content analysis to determine underlying sources of DIF. Mantel-Haenszel (MH) and Logistic Regression (LR) as two popular methods of DIF detection were employed to analyze the data obtained from ۱۵۵۰ test takers in ۲۰۱۰. The findings indicated that even though ۲۸% of items were initially detected by MH and LR as displaying gender DIF, the effect size of DIF was mostly negligible. Moreover, the content analysis phase of the study showed that sometimes it is difficult to hypothesize the linguistic element causing DIF in items. However, humanities-oriented subjects were rated as favoring females and science-oriented subjects were rated as favoring males. Finally, a correlation index of .۹۰ manifested that MH and LR produce highly consistent DIF results. These findings are discussed and implications for test developers and DIF researchers are provided.

Key terms: Fairness, DIF, Uniform DIF, Non-Uniform DIF, MH, LR, UTEPT

۱. Introduction

The present Differential Item Functioning (DIF) study is an attempt to examine University of Tehran English Proficiency Test (UTEPT) for gender DIF in order to ensure test fairness and detect the potential sources of bias. Test fairness is an issue of utmost importance in language testing which is closely related to test validity and test validation (Kunnan, ۲۰۱۰; Xi, ۲۰۱۰). DIF is said to be present when examinees from different groups show differing probabilities of success on an item after they are matched on the underlying ability that the item is intended to measure (Zumbo, ۱۹۹۹).

It is to be emphasized that finding DIF in an item does not necessarily imply that the item is biased, that is, unfair to one of the groups (Angoff, ۱۹۹۳). DIF is a necessary but not sufficient condition for bias. An item may show DIF but not be biased if the difference is because of actual differences in the groups' ability needed to answer the item, for example, if one group is high proficiency and the other low proficiency: The low proficiency group would necessarily score much lower. Only where the difference is caused by construct-irrelevant factors can DIF be viewed as bias. In other words, an item that shows DIF needs to be investigated further to uncover the reasons for its differential functioning.

^۱Hakim Sabzevari University, Sabzevar, Iran. Email: Sm.Amirian@hsu.ac.ir

^۲University of Tehran, Tehran, Iran. Email: Mohammed.Alavi@gmail.com

^۳Universidad de Oviedo, Plaza de Feijoo, s/n, Oviedo, ۳۳۰۰۳, Spain. Email: Fidalgo@umiovi.es

As one of the earliest and most popular DIF detection techniques, Mantel-Haenszel (MH) has been widely used in DIF studies (Dorans & Kulick, ۲۰۰۶; Hambleton & Rogers, ۱۹۸۹; Roever, ۲۰۰۷). Interest in using the MH procedure was shown after publication of the paper by Holland and Thayer (۱۹۸۸). The primary reasons for the popularity of MH are attributed to its computational simplicity, ease of implementation, associated test of significance and capacity for detecting DIF with small sample sizes (Fidalgo, Alavi, & Amirian, ۲۰۱۴; Fidalgo, Ferreres, & Muniz, ۲۰۰۴; Fidalgo, Hashimoto, Bartram, & Muniz, ۲۰۰۷; Rogers & Swaminathan, ۱۹۹۳).

As DIF detection is method dependent, DIF researchers highly recommended to use more than one method in order to cross-validate DIF results (e.g., Ferne & Rupp, ۲۰۰۷). Logistic Regression (LR) is employed in this study as a second method to triangulate DIF detection methods. LR has recently been widely used for DIF detection in many disciplines including language testing (Alavi, Amirian & Rezaee, ۲۰۱۱; Fidalgo, Alavi, & Amirian, ۲۰۱۴; Breland, Lee, Najarian, & Muraki, ۲۰۰۴; Kim, ۲۰۰۱). LR analysis has been used mainly to study group effect for dichotomously scored test items (Swaminathan & Rogers, ۱۹۹۰), but French and Miller (۱۹۹۶) demonstrated that this procedure can be extended for polytomous items as well. Some other reasons for the popularity of LR are that it allows modeling of Uniform DIF (UNDIF) and Non-Uniform (NUDIF) and requires less complicated computing than IRT-based analyses.

Although gender DIF has been extensively researched in first language contexts (Aryadoust, Goh, & Kim, ۲۰۱۱; Li & Suen, ۲۰۱۳; Ryan & Bachman, ۱۹۹۲; Zhang, Dorans & Matthews-Lopez, ۲۰۰۵) few studies have dealt with DIF in EFL/ESL settings (Alavi, Rezaee, & Amirian, ۲۰۱۱, Pae, ۲۰۰۴; Park & French, ۲۰۱۳, Shimizu & Zumbo, ۲۰۰۵). In fact, the literature on DIF studies on tests that are developed and administered outside the United States is critically meager (Ferne & Rupp, ۲۰۰۷). To bridge this gap, and in an attempt to validate UTEPT, the present DIF study uses MH and LR to find out whether UTEPT as a test developed and administered in Iran shows substantial DIF in favor of a specific gender group. Moreover, by conducting a content analysis of DIF items, the possible underlying linguistic sources of DIF are examined. Therefore, the study addresses the following questions:

Q^۱. Do items of UTEPT show substantial DIF in favor of males or females after matching on ability?

Q^۲. To what extent are MH and LR gender DIF findings comparable?

Q^۳. Does content analysis of DIF items indicate bias in UTEPT?

۲. Previous Gender DIF Studies

Many studies have recently been conducted to detect DIF with language tests (e.g. Alavi, Rezaee, & Amirian, ۲۰۱۱, Geranpayeh & Kunnan, ۲۰۰۷, Pae & Park, ۲۰۰۶) and DIF investigation has become “a key component of validity studies in virtually all large-scale assessments” (Penfield & Camilli, ۲۰۰۷, p. ۱۲۵). Nevertheless, research into gender DIF has indicated conflicting results (Mielikainen, ۱۹۸۸; Tannen, ۱۹۹۰). Some of these studies are reviewed here.

In an early DIF study, Ryan and Bachman (۱۹۹۲) investigated differential performance on the TOEFL and the FCE. They found little evidence that males and females performed differently at the item level on either test. Wainer and Lukhele (۱۹۹۷) also reported that the reading comprehension testlets of TOEFL showed essentially no differential functioning by

gender. This is in contrast with the findings of Carlton and Harris (۱۹۹۲) who in their analysis of gender DIF using MH and LR found that overall the female group performed differentially better than a matched male group.

Many studies have also focused on DIF investigation at skills level. In their investigation of the comparability of TOEFL writing prompts, Breland et al. (۲۰۰۴) employed LR procedures to estimate prompt difficulty and gender effects. They found that open-ended questions generally favored female test takers, which they hypothesized might have been due to gender differences in reasoning and cognitive processes.

MELAB listening test was examined for gender DIF by Aryadoust, Goh, and Kim (۲۰۱۱) using Rasch measurement. The result of a *t*-test UDIF analysis showed that two test items displayed substantive DIF favoring different gender subgroups; and NUDIF analysis revealed several test items with significant DIF, many of which favored low-ability male test takers.

Among very few DIF studies with a non-U.S sample, Gafni (۱۹۹۱) used the Standardization and the MH procedures to examine gender DIF. It was found that only three out of a total of ۵۰ English items were identified as showing gender DIF. Similarly, Lin and Wu (۲۰۰۳) in their study of EPT in China found that out of a total of ۱۲۰ items, only two items revealed C-level or large DIF and in fact ۸۹% of the test items manifested “no” DIF whatsoever.

To sum up, it can be said that research into gender DIF has produced inconsistent results and the issue needs further investigation using more elaborate methodologies. To this end, as a DIF study in an EFL context, the current study focuses on DIF analysis of UTEPT. In a previous DIF investigation with UTEPT focusing on academic discipline differences, Alavi, Rezaee, and Amirian (۲۰۱۱) found that overall UTEPT shows no DIF for humanities vs. science and engineering groups. Two previous studies (Rezaee & Shabani, ۲۰۱۰, and Karami, ۲۰۱۱) also examined earlier versions of UTEPT for gender DIF. In the first study, Rezaee and Shabani (۲۰۱۰) studied ۶,۵۵۵ examinees and found that ۳۹ of the ۱۰۰ items in the test displayed significant gender differences. However, these group differences were viewed as “negligible”. This finding was confirmed by Karami (۲۰۱۱) who utilized Rasch model to investigate gender DIF in UTEPT. The results of his study indicated that ۱۹ items were functioning differentially for the two groups. Only ۳ items, however, displayed DIF with practical significance. These two studies, however, did not investigate the content of DIF items to uncover underlying causes of DIF. To fill this gap, the present study investigates DIF items in all grammar, vocabulary and reading subtests of UTEPT both quantitatively and qualitatively in order to look into the linguistic sources of DIF.

۲.۱. Linguistic Causes of DIF

There is a scarcity of research into content-based reasons of DIF in DIF literature. In one of the early DIF studies investigating causes of DIF, Kim (۲۰۰۱), examined DIF across a sample of ۱۰۳۸ Asian and European participants in a speaking test using Likelihood ratio test and the LR methods. The results showed that ‘grammar’ and ‘pronunciation’ functioned differentially across the two groups. A content analysis of the DIF items suggested that the types and the numbers of scoring scales might influence the test validity.

In another study, Uiterwijk and Vallen (۲۰۰۵) investigated linguistic sources of item bias for second generation immigrants in Dutch tests employing IRT and MH. An investigation of sources of DIF was carried out which was followed by an in-depth analysis

of DIF items by the researchers and by external experts. The findings indicated that it was not easy for judges to decide if the DIF items were biased or not, but the description of the domains the items claim to measure turned out to be very informative and helpful in getting a final decision. While ۱۷.۴% of the items showed DIF, the researchers eventually concluded that only ۴% of all the items were biased.

In a similar study, Geranpayeh and Kunnan (۲۰۰۷) concentrated on the classification of linguistic sources of DIF of immigrant students. They investigated whether the test items on the listening section of the Certificate in Advanced English examination functioned differently for test takers from three different age groups. The main results showed that although a few items were detected based on statistical and content analyses procedures, expert judges could not clearly identify the sources of differential item functioning for the items. As it is essential to identify source of DIF, we employ a DIF questionnaire in the content analysis phase of the study to elicit expert panels' judgments (Ferne & Rupp, ۲۰۰۷) to uncover underlying sources of linguistic bias in UTEPT.

۲. Methodology

The data for the present study were gathered from ۱۵۵۰ test takers who took UTEPT in ۲۰۱۰. The sample was divided into a reference group of ۸۹۹ male and a focal group of ۶۵۱ female test takers with various age ranges. The participants were all PhD candidates seeking to sit for PhD exams of the University of Tehran. After generating the data, the descriptive statistics was computed using SPSS. The results indicated a slight difference between the mean score of male ($M=۵۰.۹۶$, $SD=۱۲.۷۰$) and female ($M=۵۰.۶۱$, $SD=۱۳.۱۴$) examinees. The effect size of mean difference is estimated small according to Cohen's (۱۹۸۸) test ($d=۰.۰۲۷$). The reliability of the data was also estimated ۰.۸۸ using Cronbach's alpha which indicates that the test enjoys high reliability.

MH statistic can be calculated using easily accessible general statistical softwares (SPSS, SAS) or more specific packages (MHDIF: Fidalgo, ۱۹۹۴; DIFAS: Penfield, ۲۰۰۵, ۲۰۰۹). In the present study, DIFAS ۵.۰ (Penfield, ۲۰۰۹) was used for MH DIF detection. In addition to reporting MH χ^2 (Holland and Thayer, ۱۹۸۸), DIFAS has the advantage of reporting a test of effect size based on ETS classification scheme (Zieky, ۱۹۹۳).

In order to run LR as a second method of DIF analysis, the Nagelkerke's SPSS syntax for nominal data (Zumbo, ۱۹۹۹) was utilized. The LR procedure uses the item response (۰ for incorrect response or ۱ for correct response) as the dependent variable, with grouping variable (dummy coded as ۱=male, ۲=female), total score (characterized as variable TOT) and a group by TOT interaction as independent variables. This appears in the following equation: $Y = b_0 + b_1 \text{TOT} (\theta) + b_2 \text{Group} (g) + b_3 \text{TOT} * \text{Group} (\theta * g)$.

After computing the two-degree-of-freedom χ^2 value as a test of DIF significance, Zumbo-Thomas' (۱۹۹۷) effect size measure was calculated. Jodoin and Gierl's (۲۰۰۱) more conservative criteria were also employed as a second test of DIF magnitude. Items which showed DIF were subjected to further content analysis using a DIF questionnaire developed by Geranpayeh and Kunnan (۲۰۰۷). Accordingly, two content experts were asked to rate the suitability of the test items for each gender group using a questionnaire with a ۵-point scale. It was expected that if evidence exists that the items detected with DIF were advantaging a particular group of test takers, then the items may be biased.

۴. Results

۴.۱. Mantel-Haenszel Results

The aim of this study was to assess the possibility of existence of DIF in UTEPT, to check the comparability of MH and LR DIF findings, and to examine the content of DIF items for potential sources of linguistic bias.

The result of MH analysis for males and females is presented in Table ۲. For the sake of convenience, only items that are flagged as DIF items are shown in Table ۳. Out of a total of ۱۰۰ UTEPT items, ۳۱ items (۳۱%) are flagged with MH gender DIF. According to ETS classification scheme, however, ۲۴ items displayed category A (negligible) effect size and only ۷ items manifested category B (Moderate) size DIF while no item was categorized as category C (large) DIF magnitude.

Table ۲. MH Gender DIF Results

| Item | Section | Favored | MH CHI | MH LOR | ETS |
|------|---------|---------|--------|--------|-----|
| ۱ | G | F | ۱۰.۷۰۹ | -۰.۴۳۷ | B |
| ۲ | G | F | ۶.۰۵۳ | -۰.۲۸۳ | A |
| ۳ | G | F | ۱۰.۶۸۲ | -۰.۳۶۸ | A |
| ۸ | G | F | ۶.۱۲۳ | -۰.۲۷۶ | A |
| ۹ | G | F | ۹.۸۶۶ | -۰.۳۸۷ | A |
| ۱۰ | G | F | ۷.۶۰۳ | -۰.۳۷۱ | A |
| ۱۳ | G | F | ۷.۹۰۸ | -۰.۳۲۷ | A |
| ۲۰ | G | F | ۱۶.۷۲۳ | -۰.۴۷۳ | B |
| ۲۲ | G | F | ۲۵.۸۸۷ | -۰.۶۲۵ | B |
| ۲۴ | G | F | ۷.۳۹۹ | -۰.۳۴۲ | A |
| ۲۵ | G | F | ۷.۹۴۲ | -۰.۳۵۹ | A |
| ۲۶ | G | F | ۸.۶۵۴ | -۰.۵۲۲ | B |
| ۲۸ | G | F | ۶.۶۸۰ | -۰.۳۳۷ | A |
| ۲۹ | G | F | ۰.۱۴۰ | -۰.۰۴۷ | A |
| ۳۵ | V | M | ۴.۸۵۳ | ۰.۲۴۸ | A |
| ۳۸ | V | M | ۱۲.۵۲۱ | ۰.۴۹۵ | B |
| ۴۱ | V | M | ۹.۹۵۹ | ۰.۴۶۴ | B |
| ۴۴ | V | M | ۷.۱۷۰ | ۰.۳۱۰ | A |
| ۴۸ | V | M | ۲.۵۹۸ | ۰.۲۲۱ | A |
| ۴۹ | V | M | ۵.۴۵۸ | ۰.۲۹۴ | A |
| ۵۰ | V | M | ۹.۹۴۸ | ۰.۴۰۰ | A |
| ۵۱ | V | M | ۰.۴۸۵ | ۰.۰۸۲ | A |
| ۵۲ | V | M | ۱۶.۶۰ | ۰.۵۰۵ | B |
| ۵۴ | V | M | ۶.۷۵۳ | ۰.۳۰۹ | A |
| ۵۹ | V | M | ۸.۱۷۱ | ۰.۳۲۶ | A |
| ۸۱ | R | F | ۵.۸۰۷ | -۰.۳۷۰ | A |
| ۹۰ | R | F | ۰.۳۹۲ | -۰.۰۸ | A |

| | | | | | |
|----|---|---|-------|--------|---|
| ۹۱ | R | M | ۶.۶۷۱ | ۰.۲۹۳ | A |
| ۹۲ | R | F | ۲.۷۱۱ | -۰.۲۰۳ | A |
| ۹۵ | R | M | ۵.۲۸۰ | ۰.۳۲۴ | A |
| ۹۷ | R | M | ۰.۴۰۲ | ۰.۰۷۶ | A |

Notes. * $p < .05$; G= Grammar; V= Vocabulary; R= Reading; F = Female; M = Male; A = negligible DIF; B = moderate DIF

MH CHI = Mantel-Haenszel Chi-Square statistic which is distributed as chi-square with one degree of freedom.

MH LOR= Mantel-Haenszel Common Log-Odds Ratio. Positive values indicate DIF in favor of males and negative values indicate DIF in favor of females.

DIF items come from different sections of UTEPT. Out of ۳۱ MH DIF items, ۱۴ items belong to the grammar section, ۱۱ items to the vocabulary section and ۶ items to the reading comprehension section. This shows that the grammar section displays more gender DIF items than vocabulary and reading sections. In fact, both vocabulary and reading sections contain ۱۷ DIF items while grammar section alone incorporates ۱۴ DIF items.

In terms of the direction of DIF, ۱۷ items favored females and ۱۴ items favored males. It was also indicated that the grammar section primarily worked to the advantage of females, vocabulary section to the advantage of males and reading comprehension section neither to the advantage of males nor females.

۴.۲. Logistic Regression Results

The findings of LR DIF are summarized in Table ۳. It should be noted that only items with significant DIF values at ۰.۰۵ level of significance are included in the table. It appears that the two-degree-of-freedom chi-squared (χ^2) test of the significance for DIF is significant for ۲۹ items. Of these ۲۹ items, ۱۴ items are detected in the grammar section, nine items in the vocabulary section and only six items in the reading comprehension section.

The obtained R^2 values reveal that gender DIF is predominantly of uniform nature on UTEPT. Of a total of ۲۹ DIF items, ۲۵ items displayed UDIF and only four items displayed NUDIF. ۱۹ of UDIF items favored the female group and only ۶ items favored the male group. Additionally, based on Copella and Sireci's (۲۰۰۹) guidelines, three NUDIF items (۵۱, ۹۰, and ۹۲) were found to favor females and only one item (۴۸) appeared to favor males. All in all, out of ۲۹ DIF items, ۱۷ items favored females and ۱۲ items favored males. Concerning DIF effect size, based on Jodoin and Gierl's (۲۰۰۱) conservative classification scheme, it is observed that all obtained R^2 values are category A that is smaller than ۰.۳۵. This means that all LR DIF items manifest a negligible DIF magnitude based on guidelines proposed by both Zubmo and Thomas (۱۹۹۷) and Jodoin and Gierl (۲۰۰۱).

Table ۳. LR Gender DIF Results

| Item | Section | Favored | UDIF | NUDIF | DIF size | χ^2 | P* | Category |
|------|---------|---------|------|-------|----------|----------|------|----------|
| ۱ | G | F | .۰۱۳ | .۰۰۲ | .۰۱۵ | ۱۷.۷۴۰ | .۰۰۰ | A |
| ۲ | G | F | .۰۰۷ | .۰۰۰ | .۰۰۷ | ۹.۲۴۲ | .۰۱۰ | A |
| ۳ | G | F | .۰۰۸ | .۰۰۰ | .۰۰۸ | ۱۰.۱۶۴ | .۰۰۶ | A |
| ۸ | G | F | .۰۰۴ | .۰۰۲ | .۰۰۶ | ۷.۱۵۲ | .۰۲۸ | A |
| ۹ | G | F | .۰۰۵ | .۰۰۱ | .۰۰۶ | ۸.۹۲۵ | .۰۱۲ | A |
| ۱۰ | G | F | .۰۱۰ | .۰۰۰ | .۰۱۰ | ۹.۵۶۱ | .۰۰۸ | A |
| ۱۳ | G | F | .۰۰۷ | .۰۰۰ | .۰۰۷ | ۷.۶۹۸ | .۰۲۱ | A |
| ۲۰ | G | F | .۰۱۵ | .۰۰۱ | .۰۱۵ | ۲۳.۴۷۶ | .۰۰۰ | A |
| ۲۲ | G | F | .۰۱۹ | .۰۰۲ | .۰۲۱ | ۲۸.۲۸۰ | .۰۰۰ | A |
| ۲۴ | G | F | .۰۰۶ | .۰۰۰ | .۰۰۶ | ۷.۸۹۳ | .۰۱۹ | A |
| ۲۵ | G | F | .۰۰۵ | .۰۰۱ | .۰۰۶ | ۶.۷۹۰ | .۰۳۴ | A |
| ۲۶ | G | F | .۰۰۹ | .۰۰۰ | .۰۰۹ | ۱۰.۰۷۶ | .۰۰۶ | A |
| ۲۸ | G | F | .۰۰۷ | .۰۰۱ | .۰۰۸ | ۱۰.۲۸۱ | .۰۰۶ | A |
| ۳۵ | V | M | .۰۰۴ | .۰۰۳ | .۰۰۷ | ۸.۵۰۹ | .۰۱۴ | A |
| ۳۸ | V | M | .۰۱۲ | .۰۰۱ | .۰۱۳ | ۱۳.۵۰۴ | .۰۰۱ | A |
| ۴۴ | V | M | .۰۰۵ | .۰۰۰ | .۰۰۵ | ۶.۱۳۹ | .۰۴۶ | A |
| ۴۸ | V | M | .۰۰۳ | .۰۰۷ | .۰۱۰ | ۹.۶۷۴ | .۰۰۸ | A |
| ۴۹ | V | M | .۰۰۴ | .۰۰۲ | .۰۰۶ | ۸.۲۲۶ | .۰۱۶ | A |
| ۵۰ | V | M | .۰۰۶ | .۰۰۰ | .۰۰۶ | ۷.۹۶۵ | .۰۱۹ | A |
| ۵۱ | V | F | .۰۰۰ | .۰۰۵ | .۰۰۵ | ۶.۳۲۲ | .۰۴۲ | A |
| ۵۲ | V | M | .۰۱۱ | .۰۰۱ | .۰۱۲ | ۱۵.۸۶۱ | .۰۰۰ | A |
| ۵۴ | V | M | .۰۰۵ | .۰۰۳ | .۰۰۶ | ۷.۰۰۵ | .۰۳۰ | A |
| ۵۹ | V | M | .۰۰۶ | .۰۰۱ | .۰۰۹ | ۱۱.۸۴۱ | .۰۰۳ | A |
| ۶۷ | R | M | .۰۰۶ | .۰۰۰ | .۰۰۷ | ۷.۹۶۱ | .۰۱۹ | A |
| ۸۱ | R | F | .۰۰۵ | .۰۰۵ | .۰۰۵ | ۶.۰۱۱ | .۰۵۰ | A |
| ۹۰ | R | F | .۰۰۱ | .۰۰۰ | .۰۰۶ | ۶.۶۵۱ | .۰۳۶ | A |
| ۹۱ | R | M | .۰۰۶ | .۰۰۶ | .۰۰۶ | ۸.۰۶۰ | .۰۱۸ | A |
| ۹۲ | R | F | .۰۰۱ | .۰۰۴ | .۰۰۷ | ۱۰.۱۲۱ | .۰۰۶ | A |
| ۹۵ | R | M | .۰۰۴ | .۰۰۱ | .۰۰۸ | ۸.۸۲۷ | .۰۱۲ | A |

Notes. * $p < .۰۵$; G= Grammar; V= Vocabulary; R= Reading; F = Female; M = Male; A = negligible DIF

As far as the comparability of MH and LR DIF findings is concerned, it was found that MH detected two more gender DIF items (۳۱) than LR (۲۹) (Table ۴). This finding contrasts with the dominant view in DIF literature that LR detects more DIF items in comparison to MH due to its capability of detecting both UDIF and NUDIF (Hidalgo, & López-Pina, ۲۰۰۴; Rogers & Swaminathan, ۱۹۹۳). Nevertheless, Phi correlation for nominal data was run to compare the performance of MH and LR methods in flagging similar DIF items in all ۱۰۰ items. The flagged items by either method were dummy coded as ۱ and the unflagged ones were dummy coded as ۰. Phi correlation value of .۹۰ indicated a highly significant correlation between the findings of the two methods.

Table ۴. Comparison of Gender DIF Results in Sections of UTEPT

| items methods | Number of items | Number of items | Number of |
|-----------------------|------------------------|------------------------|---------------------|
| | detected by MH only | detected by LR only | detected by both |
| Grammar | ۱۴ | ۱۳ | ۱۳ |
| Vocabulary | ۱۱ | ۱۰ | ۱۰ |
| Reading Comprehension | ۶ | ۶ | ۵ |
| Total | ۳۱ | ۲۹ | ۲۸ |

In terms of specific sections of DIF items, MH detected ۱۵ gender DIF items in grammar, ۱۰ items in vocabulary, and six items in reading comprehension sections while LR flagged ۱۴ DIF items in grammar, nine items in vocabulary and six items in the reading comprehension sections. This shows a great consistency between MH and LR findings even in their DIF detection in sections of UTEPT. As far as the magnitude of DIF is concerned, however, it was found that all DIF items detected by LR method displayed a negligible or type-A effect size while MH detected ۷ items with moderate or type-B effect size.

۴.۳. Content Analysis of DIF Items

۲۸ DIF items detected by both LR MH and LR methods were examined based on elicited responses from two expert judges. The results of content analysis for gender DIF are reported below.

Grammar Items. It is interesting that grammar section of UTEPT contains greater number of DIF items than vocabulary and reading comprehension sections. Among ۱۴ items flagged with DIF in the grammar section, seven items (۱, ۲, ۳, ۹, ۱۰, and ۱۳) came from the *fill in the blank structure subsection* four items (۲۰, ۲۲, ۲۴, and ۲۵) came from the *written expression subsection* and three items (۲۶, ۲۸, and ۳۵) came from the *grammar in context (cloze test) subsection* of the test. For many DIF items in this section, neither of the content specialists could hypothesize the sources of DIF. Item ۱ represents such items.

Item ۱. *Having been the prize, the professor continued working hard on his project.*

- A. awarded
- B. award
- C. awarding
- D. the award

As it is shown in the Item Characteristic Curve (ICC) of this item (Figure ۱), line ۱ and ۲ do not cross which means that this item shows UDIF favoring females at all levels of proficiency. This item shows moderate DIF magnitude in favor of female test takers. Content experts believed there is no clue in such a short grammar item as to why this item is working to the advantage of females.

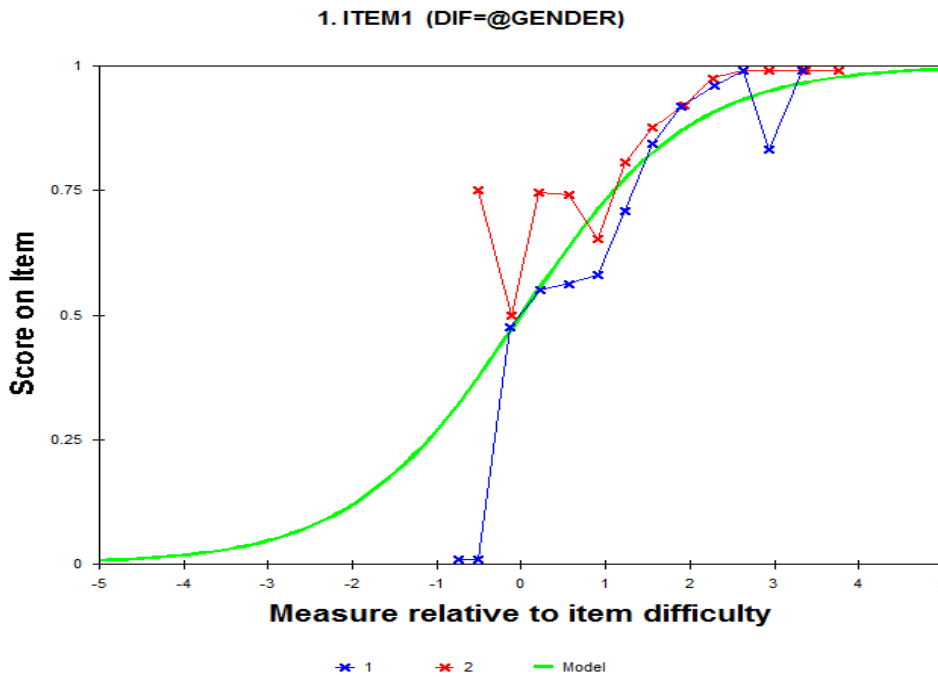


Figure ۱. ICC for item ۱

Note. ۱ = the performance of males; ۲ = the performance of females

Item ۲۶, which belongs to the *grammar in context (cloze test)* subsection, is an example of very few grammar items that content specialists unanimously rated as favoring females.

Item ۲۶. *Without regular supplies of some hormones, our capacity to behave would be seriously impaired; without others we would soon die. Tiny amount of some hormones can..... ۲۶.....our moods and our actions.*

- A. Modification
- B. Modifying
- C. Modify
- D. Modified

This item was shown to advantage examinees in the female group by both MH and LR methods. The expert judges believed that since females show more interest in topics such as *human biology*, they are systematically favored by this item. One of the judges also pointed out that words such as "*behave and mood*" in the item indicate social interactions which is a topic of women's interest.

Vocabulary Items. Out of ۳۰ DIF items, ۹ items in the vocabulary section were flagged with DIF by both methods. The underlying cause of many of DIF items in vocabulary section was attributed to the texts in which these words were more likely to occur. For example, one of experts commented that in item ۳۸ the word “*fundamental*” (the correct choice) is more likely to be found in scientific texts which are of more interest to males.

Item ۳۸. *Analytic tools enable one to get at the most fundamental logic of any discipline.*

- A. *enforced*
- B. *essential*
- C. *established*
- D. *escorted*

Reading Comprehension Items. Out of ۳۵ items in this section, only six items (۱۷% of UTEPT items) were found to manifest DIF (items ۸۱, ۹۰, ۹۱, ۹۲, ۹۵ and ۹۷). Among these six DIF items, item ۸۱ advantages females and comes from passage three which is on *history*. Content specialist commented that the content of this passage is responsible for differential performance because females score higher on humanities-oriented passages. Items ۹۰, ۹۱, ۹۲ and ۹۵ (favoring females, males, females and males respectively) come from passage six on *painting art*. Item ۹۲ which is a “*scanning details*” question is presented below. The ICC for this item also appears in Figure ۲.

Item ۹۲. *According to the passage, which of the following was one of the distinguishing characteristics of Impressionist painting?*

- A. *The emphasis on people rather than nature scenes*
- B. *The way the subjects were presented from multiple angles*
- C. *The focus on small solid objects*
- D. *The depiction of the effects of light and color*

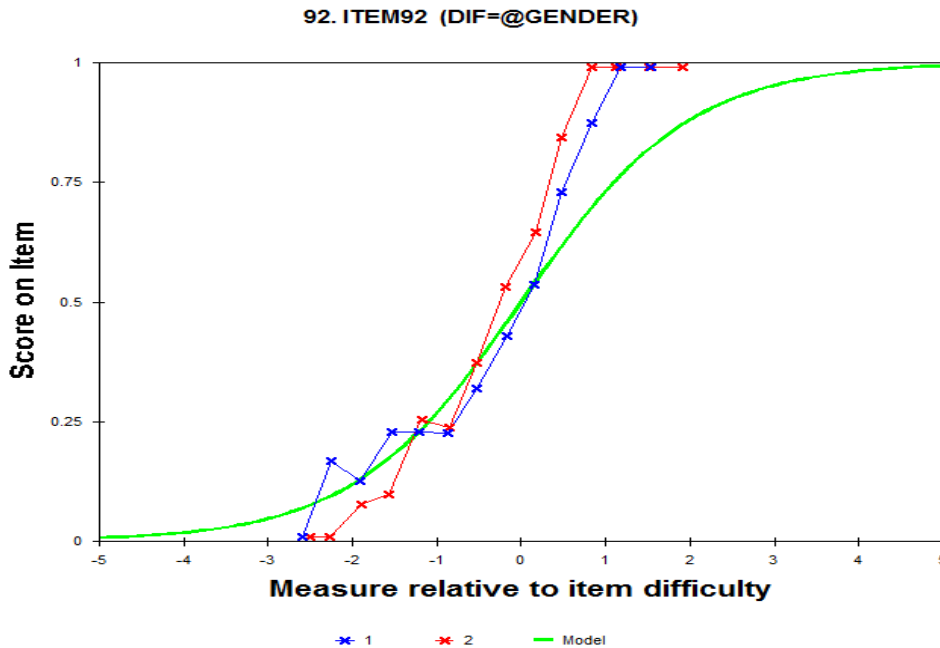


Figure ۲. ICC for item ۹۲

Note. ۱ = the performance of males; ۲ = the performance of females

The ICC for this NUDIF item indicates that line ۱ and ۲ cross at some points on the curve. It means that there is an interaction between proficiency level and grouping factor: the male group has a greater probability of answering an item correctly within the lower score range and the female group has a higher probability of answering an item correctly in the higher score range. However, overall the item favors females.

Content experts found *history* topic more humanities-oriented and as a result rated this item along with items ۹۰ (*main idea question*), ۹۱ (*vocabulary question*), and ۹۵ (*scanning details question*) in favor of females. However, according to DIF findings items ۹۱ and ۹۵ favored males.

۵. Discussion

The results of the study suggested that; altogether, MH flagged ۳۱ DIF items and LR flagged ۲۹ items. Moreover, ۲۸ items were flagged with gender DIF by both methods in grammar (۱۴ items), vocabulary (۹ items), and reading comprehension sections (۵ items). This indicates that MH and LR produce consistent result which is confirmed by a high correlation index of .۹۰. This is in line with the findings of previous DIF studies that various DIF detection methods show a close correlation (e.g. Rogers & Swaminathan, ۱۹۹۳).

The percentage of flagged DIF items in the present study (۲۸%) is smaller than Rezaee and Shabani's (۲۰۱۰) finding (۳۹%) and larger than Karami's (۲۰۱۱) finding (۱۹%) with earlier editions of UTEPT. The highest number of DIF items were detected in the grammar section while reading comprehension section incorporates the lowest number of DIF items which is against our expectations as the reading section of many tests is found to be responsible for differential performance of them and this skill is the one that has been most extensively studied for DIF (e.g. Jiao & Chen, ۲۰۱۴; Li & Suen, ۲۰۱۳; Pae, ۲۰۰۴).

A closer look at the direction of DIF items reveals that surprisingly most of the items that favor females belong to the grammar section (۱۴ items) and only three items belong to the reading section while no item in vocabulary section is in favor of females. In addition, most items that favor males come from the vocabulary section (۱۰ items) while only one grammar item and three reading comprehension items favor males. Thus, it is concluded that overall the grammar section of UTEPT is in favor of females, the vocabulary section is in favor of males and the reading section is neither in favor of females nor males. This finding indicates that the developers or users of UTEPT should be aware of the test-takers for whom the test is intended. For example, high scores on the grammar subtest may not provide sufficient information to permit inferences about a female test taker's overall English ability since it showed DIF in advantage of females.

The findings of the content analysis phase of the study indicated that determining underlying causes of gender DIF was very hard for content specialists especially with short grammar and vocabulary items. Conoley (۲۰۰۳) also pointed out that it is not always clear which element in an item causes DIF. Nonetheless, for ۷% of the items it was the content of the item that was deemed responsible for differential performance. Experts considered humanities-oriented topics in favor of females and science-oriented topics in favor of males giving support

to Doolittle and Welch (۱۹۸۹) who reported that females scored higher than males with humanities-oriented passages, but lower than males with science-oriented passages. This can be attributed to the educational system in Iran in which up to the recent years women showed more interest in studying in humanities majors and men showed greater interest in science and engineering majors. Of course, this pattern has changed lately and over ۶۰٪ of classes at universities are currently occupied by girls in all majors including science and engineering.

The findings of this study provide some implications for DIF researchers. The test of effect size in this study revealed that out of initially ۲۸ flagged items, MH only detected ۷ items with moderate magnitude while LR detected no sizable DIF magnitude whatsoever. This finding supports the results of Rezaee and Shabani's (۲۰۱۰) who did not detect any sizable DIF item with UTEPT. It also highlights the significance of reporting the effect size in DIF studies (Cohen, ۱۹۸۸; Kirk, ۱۹۹۶). Although considering the direction of DIF can give us a lot of information about the behavior of DIF items, if the magnitude of DIF along with the direction of DIF is considered, a better picture of misbehaving items is obtained.

UTEPT developers and administrators may also benefit from the findings of this study. First, considering the large number of applicants who take this test each year, it is obvious that a DIF study needs to be a necessary part of the validation process of this test. Second, in order to promote gender equity, test developers should make sure the passage topics appeal to both genders and cover a wide range of academic subjects. Ultimately, although the magnitudes of DIF items are relatively small, the present pool of UTEPT items contains ۲۸ items that produce statistically significant gender differences. Although the observed differences are not large by accepted statistical standards, UTEPT developers should adopt a policy by specifying what levels of difference should result in items being revised or dropped from active administration.

۵.۱. Conclusion and Suggestions for Further Research

This study made effort to make a contribution to the DIF literature by providing information about DIF with an Iranian sample. Since this is a test developed and administered in an EFL context, it may have some specific features that are observed only in Iranian contexts and is not common across nationalities. Even though initially about ۲۸٪ were consistently flagged as displaying DIF by MH and LR due to significant uniform or non-uniform group effects, their effect sizes were far too small for most of them to render the test unfair. That is, the item score differences between males and females compared in this study primarily seem to be because of item impact rather than group difference attributable to a construct irrelevant factor inherent in items or item bias. This fact reveals the importance of follow-up content analysis phase in DIF studies as not every DIF item is necessarily biased.

The present study was limited to the study of DIF at item level. Future research could address the differential functioning of the whole UTEPT test or sections of the test (bundles), particularly the grammar and vocabulary sections of UTEPT that were found to contain the largest number of DIF items in the current study.

Moreover, the current study analyzed the content of items only for linguistic sources of DIF. It is quite possible, however, that the DIF source is of a nonlinguistic nature. Therefore, other studies could be carried out looking for nonlinguistic sources of variance that contribute to differential performance of tests including item type, examinees' sociocultural background and

examinees' age. Finally, in the present study the researcher only consulted expert judges in order to identify the sources of DIF. Since identifying the potential sources of DIF in a test is a multidimensional process, interviews with examinees could be a promising line of inquiry into the causes of DIF.

Acknowledgments

Dr. Fidalgo's work in this article was supported by the Spanish Ministry of Economy and Competitiveness (grant number PSI۲۰۰۹-۰۸۵۲۹).

References

- Alavi, S. M., Rezaee, A. & Amirian, S. M. R. (۲۰۱۱). Academic discipline DIF in an English language proficiency test. *Journal of English Language Teaching and Learning*, ۵(۷): ۳۹-۶۵.
- Angoff, W. H. (۱۹۹۳). Perspectives on differential item functioning methodology. In P. W. Holland and H. Wainer (Eds). *Differential Item Functioning* (pp.۳-۲۵). New York: Routledge
- Aryadoust, V., Goh, C. C. M. & Kim, L. O. (۲۰۱۱). An investigation of differential item functioning in the MELAB listening test. *Language Assessment Quarterly*, ۸(۴): ۳۶۱-۳۸۵.
- Breland, H., Lee, Y., Najarian, M. & Muraki, E. (۲۰۰۴). *An analysis of TOEFL CBT writing prompt difficulty and comparability for different gender groups* (ETS-RR-۰۴-۰۵, Report ۷۶). Princeton, NJ: Educational Testing Service.
- Carlton, S. T., & Harris, A. M. (۱۹۹۲). *Characteristics associated with differential item functioning on the Scholastic Aptitude Test: Gender and majority/minority group comparisons*. Princeton, NJ: Educational Testing Service.
- Cohen. J. (۱۹۸۸). *Statistical power analysis for behavioral sciences*. Hillsdale, NJ: Lawrence Erlbaum Association.
- Conoley, C. A. (۲۰۰۳). *Differential item functioning in the Peabody Picture Vocabulary Test – Third Edition: Partial correlation versus Expert judgment*. PhD Thesis, Texas A and M University, TX.
- Copella, J., & Sireci, S. G. (۲۰۰۹). *Interpreting non-uniform DIF*. Poster for NCME ۲۰۰۹ Graduate Student Poster Session. University of Massachusetts Amherst.
- Doolittle, A., Welch, C. (۱۹۸۹). Gender differences in performance on a college level achievement test. Iowa City, IA.: American College Testing Program.
- Dorans, N.J., Kulick, E. (۲۰۰۶). Differential item functioning on the Mini-Mental State Examination: An application of the Mantel-Haenszel and standardization procedures. *Medical Care*, ۴۴(۱۱, Suppl.۳), S۱۰۷-S۱۱۴.
- Ferne, T., Rupp, A. A. (۲۰۰۷). A synthesis of ۱۵ years of research on DIF in language testing: Methodological advances, challenges, and recommendations. *Language Assessment Quarterly*, ۴(۲): ۱-۳۶.
- Fidalgo, A.M. (۱۹۹۴). MHDIF: A computer program for detecting uniform and non-uniform differential item functioning with the Mantel-Haenszel procedure. *Applied Psychological*

- Measurement*, ۱۸, ۳۰۰-۳۱۵.
- Fidalgo, A. M., Alavi, S. M. & Amirian, S. M. R. (۲۰۱۴). Strategies for testing statistical and practical significance in detecting DIF with logistic regression models. *Language Testing*, first published on April ۱۱, ۲۰۱۴ as doi:۱۰.۱۱۷۷/۰.۲۶۵۵۳۲۲۱۴۵۲۶۷۴۸ pp. ۱-۱۹
- Fidalgo, A. M., Hashimoto, K., Bartram, D. & Muñoz J. (۲۰۰۷) Empirical Bayes versus standard Mantel-Haenszel statistics for detecting differential item functioning under small sample conditions. *Journal of Experimental Education*, ۷۵, ۲۹۳-۳۱۴.
- Fidalgo, A. M., Ferreres, D. & Muniz, J. (۲۰۰۴). Liberal and conservative differential item functioning detection using Mantel-Haenszel and SIBTEST: Implications for type I and type II rates. *The Journal of Experimental Education*, ۷۳(۱):۲۳-۳۹
- French, A.W., & Miller, T.R. (۱۹۹۶). Logistic regression and its use in detecting differential item functioning in polytomous items. *Journal of Educational Measurement*, ۳۳(۳): ۳۱۵-۳۳۲.
- Gafni, N. (۱۹۹۱). *Differential Item Functioning: performance by sex on reading comprehension tests*. ERIC Document ED ۳۳۱۸۴۴. Rockville, MD: Educational Resources Information Center.
- Geranpayeh, A., & Kunnan, A.J. (۲۰۰۷) Differential item functioning in terms of age in the certificate in advanced English examination. *Language Assessment Quarterly*, ۴(۲): ۱۹۰-۲۲۲.
- Hambleton, R.K., & Rogers, H.J. (۱۹۸۹). Detecting potentially biased test items: Comparison of IRT area and Mantel-Haenszel methods. *Applied Measurement in Education*, ۲, ۳۱۳-۳۳۴.
- Hidalgo, M. H., López-Pina, J. A. (۲۰۰۴). Differential item functioning detection and effect size: A comparison between logistic regression and Mantel-Haenszel procedures. *Educational and Psychological Measurement*, ۶۴, ۹۰۳-۹۱۵.
- Holland, P. W., & Thayer, D.T. (۱۹۸۸). Differential item functioning detection and the Mantel-Haenszel procedure. In: Wainer, H. and Braun, H.I (eds) *Test Validity*. Hillsdale, NJ: Elbaum, pp. ۱۲۹-۱۴۵.
- Jiao, H., & Chen, Y. (۲۰۱۴) Differential item and testlet functioning analysis. In: Kunnan AJ (ed) *The Companion to Language Assessment* (1st ed.). John Wiley and Sons.
- Jodoin, M. C., & Gierl, M.J. (۲۰۰۱). Evaluating type I error and power rates using an effect size measure with logistic regression procedure for DIF detection. *Applied Measurement in Education*, ۱۴, ۳۲۹-۳۴۹.
- Karami, H. (۲۰۱۱). Detecting gender bias in a language proficiency test. *International Journal of Language Studies*, ۵(۲): ۲۷-۳۸.
- Kim, M. (۲۰۰۱). Detecting DIF across the different language groups in a speaking test. *Language Testing*, ۱۸(۱):۸۹-۱۱۴.
- Kirk, R.E. (۱۹۹۶). Practical significance: A concept whose time has come. *Educational and Psychological Measurement*, ۵۶, ۷۴۶-۷۵۹.
- Kunnan, A. J. (۲۰۰۷). Test fairness, test bias and DIF. *Language Assessment Quarterly*, ۴(۲): ۱۰۹-۱۱۲.
- Kunnan, A. J. (۲۰۱۰) Fairness matters and Toulmin's argument structures. *Language*

- Testing*, ۲۴(۲): ۱۸۳-۱۸۹.
- Lin, J., & Wu, F. (۲۰۰۳). Differential performance by gender in foreign language testing. Poster for the ۲۰۰۳ annual meeting of NCME in Chicago.
- Li, H., & Suen, H. K. (۲۰۱۳). Detecting native language group differences at the subskills level of reading: A differential skill functioning approach. *Language Testing*, ۳۰(۲): ۲۷۳-۲۹۷.
- Mielikainen, A. (۱۹۸۸). Women as preservers and innovators of spoken language. In: Laitinen L(ed), *Isosuinennainen*. Helsinki: Yliopistopainos, pp.۹۲-۱۱۱.
- Pae, T. (۲۰۰۴). Gender effect on reading comprehension with Korean EFL learners. *System*, ۳۲(۲): ۲۶۵-۲۸۱
- Pae, T., & Park, G.P. (۲۰۰۶). Examining the relationship between differential item functioning and differential test functioning. *Language Testing*, ۲۳(۴): ۴۷۵-۴۹۶.
- Park, G. P., & French, B. F. (۲۰۱۳). Gender differences in the Foreign Language Classroom Anxiety Scale. *System* ۴۱(۲): ۴۶۲-۴۷۱
- Penfield, R.D. (۲۰۰۵). DIFAS: Differential item functioning analysis system. *Applied Psychological Measurement*, ۲۹, ۱۵۰-۱.
- Penfield, R.D. (۲۰۰۹). DIFAS ۵.۰. User's manual. Manuscript in preparation
- Penfield, R.D. & Camilli, G. (۲۰۰۷). Differential item functioning and item bias. In: Sinharay, S. and Rao, C.R. (eds.), *Handbook of Statistics, Volume ۲۶: Psychometrics*. New York: Elsevier, pp.۱۲۵-۱۶۷.
- Rezaee, A., & Shabani, E. (۲۰۱۰). Gender differential item functioning analysis of the University of Tehran English Proficiency Test. *Pazhuhesh-e Zabanha-ye Khareji*, ۵۶, ۸۹-۱۰۸.
- Roever, C. (۲۰۰۷). DIF in the assessment of second language pragmatics. *Language Assessment Quarterly*, ۴(۲): ۱۶۵-۱۸۹.
- Rogers, H. J., & Swaminathan, H. (۱۹۹۳). A comparison of logistic regression and Mantel-Haenszel procedures for detecting differential item functioning. *Applied Psychological Measurement*, ۱۷, ۱۰۵-۱۱۶.
- Ryan, K., & Bachman, L. F. (۱۹۹۲). Differential item functioning on two tests of EFL proficiency. *Language Testing*, ۹(۱): ۱۲-۲۹.
- Shimizu, Y., Zumbo, B.D. (۲۰۰۵). Logistic regression for differential item functioning: A primer. *Japan Language Testing Association Journal*, ۷, ۱۱۰-۱۲۴.
- Swaminathan, H., & Rogers, H.J. (۱۹۹۰). Detecting differential item functioning using logistic regression procedures. *Journal of Educational Measurement*, ۲۷, ۳۶۱-۳۷۰.
- Tannen. D. (۱۹۹۰). *You just don't understand: women and men in conversation*. New York: William Morrow.
- Uiterwijk, H., & Vallen, T. (۲۰۰۵). Linguistic sources of item bias for second generation immigrants in Dutch test. *Language Testing* ۲۲(۲): ۲۱۱-۲۳۴
- Wainer, H., & Lukhele, R. (۱۹۹۷). How reliable are TOEFL scores? *Educational and Psychological Measurement*, ۵۷(۵): ۷۴۱-۷۵۹.
- Xi, X. (۲۰۱۰). How do we go about investigating test fairness? *Language Testing*, ۲۷(۲): ۱۴۷-۱۷۰.
- Zhang, Y., Dorans, N. J. & Matthews-Lopez, J. L. (۲۰۰۳). Using DIF dissection method to assess effects of item deletion. *Research Report No. ۲۰۰۵-۱۰*. College Board

- Zieky, M. (۱۹۹۳). Practical questions in the use of DIF statistics in test development. In: Holland P. & Wand
Wainer, H. (eds) *Differential Item Functioning*. Hillsdale, NJ: Lawrence Erlbaum, pp. ۳۲۱-۳۳۶.
Zumbo, B. D. (۱۹۹۹). *A handbook on the theory and methods of differential item functioning (DIF): Logistic regression modeling as a unitary framework for binary and Likert-type (ordinal) item scores*. Ottawa, ON: Directorate of Human Resource Research and Evaluation, Department of National Defense.
Zumbo, B. D., & Thomas, D.R. (۱۹۹۷). *A measure of effect size for a model-based approach for studying DIF*.
Working paper of the Edgeworth Laboratory for Quantitative Behavioral Sciences, University of Northern British Columbia: Prince George, B. C.