

The comparison of differential item functioning predicted through experts and statistical techniques

Sinan Yavuz*, Measurement and Assessment on Education, University of Wisconsin-Madison, Madison, Wisconsin 53706, UK.

Nuri Dogan, Department of Educational Measurement and Evaluation, Hacettepe University, Ankara 06100, Turkey.

Ronald K. Hambleton, Department of Educational Policy, University of Massachusetts Amherst, N104 Furcolo, Amherst, MA 01003.

Meltem Yurtcu, Faculty of Education, Department of Educational Sciences, Artvin Coruh University, 08000 Seyitler, Artvin, Turkey.

Suggested Citation:

Yavuz, S., Dogan, N., Hambleton, R. K. & Yurtcu, M. (2018). The comparison of differential item functioning predicted through experts and statistical techniques. *Cypriot Journal of Educational Science*. 13(2), 375–384.

Received date October 19, 2017; revised date February 07, 2018; accepted date June 05, 2018.

Selection and peer review under responsibility of Prof Dr. Huseyin Uzunboylu Near East University.

©2018 Academic World Education & Research Center. All rights reserved.

Abstract

Validity is one of the psychometric properties of the achievement tests. To determine the validity, one of the examination is item bias studies, which are based on differential item functioning (DIF) analyses and field experts' opinion. In this study, field experts were asked to estimate the DIF levels of the items to compare the estimations obtained from different statistical techniques. First, the experts were asked to examine the questions and make the DIF level estimations according to the gender variable for the DIF estimation; the agreement of the experts was examined. Second, DIF levels were calculated by using the logistic regression and Mantel-Haenszel test. Third, the experts' estimations and the statistical analyses results were compared. As a conclusion, it was observed that the experts and the statistical techniques were in agreement among themselves, and they were partially different from each other for the Sciences and equal for the Social Sciences tests.

Keywords: Item bias, differential item functioning (DIF), expert estimation.

* ADDRESS FOR CORRESPONDENCE: **Sinan Yavuz**, Measurement and Assessment on Education, University of Wisconsin-Madison, Madison, Wisconsin 53706, UK. **E-mail address:** yavuzsinan@gmail.com / **Tel.:** 608-262-3811

1. Introduction

Impartiality is a sign for validity (Camilli & Shepard, 1994). So, one of the validation studies, the application of which has become a routine in recent years, is the item bias studies. Item bias studies mostly cover a review of sensitivity and differential item functioning (DIF) (Educational Testing Service, 2009; Hambleton, 2006; Sireci & Mullane, 1994). The sensitivity review and DIF studies contribute to the achievement test scores in a proper manner (Zieky, 2002). DIF, for the achievement tests, is defined as the differentiation in the probabilities of giving the correct answer to an item on the part of the individuals with the same competency level who belong to different groups of the same population (Hambleton & Rogers, 1996; Zumbo, 1999). The fact that the probability of giving the correct answer to the item for different ability level students is expected to be different. Moreover, the probability of giving the correct answer to the item is supposed to be equal for the individuals with the same ability level who exist within the same population, even though they belong to different groups. If the students' levels, who are within the same population but belong to separate groups, have a different probability of giving the correct answer to the item, it is considered to be the item bias (Zumbo, 1999). Items with bias indication interfuse different variables to the assessment process other than examiners ability (Cromwell, 2002).

The total test scores, which contain biased items, will be non-objective, and the decision made according to the overall scores will be faulty and misjudged. As a result, the validity of the assessment will fail. At the beginning of the statistical analyses, the important point is to determine which variable the DIF analyses will be performed on (Hambleton & Rodgers, 1996; Hambleton, Ying & Klauck, 2001). DIF analyses can be carried out by taking different variables into consideration, such as gender, ethnic group, social class, subculture and beliefs. Another important respect regarding these analyses is the stage of determining the test hypothesis and statistical techniques. DIF analyses can be performed by applying statistical techniques based on the classical test theory (CTT) and item response theory (IRT). Even if the CTT hypothesis is sample dependent, CTT techniques still more practical than IRT methods (Budgell, Raju & Quartetti, 1995; Hambleton, 2006; Hambleton & Rogers, 1989; Jones & Hambleton, 1992). In this research, two statistical techniques based on the CTT were applied.

The decision on the superiority or importance of DIF values is made according to the results obtained from the statistical analysis. Since the statistical significance is influenced by the sample size of the test, calculation of the effect size is also becoming more and more common (Benito, Hidalgo & Guilera, 2010; Hambleton, 2006). While DIF is an item analysis methodology that describes the sample as whole, ignoring how the psychometric properties of the scale may vary as a function of variation within the sample (Zumbo, 1999), it uses a variety of techniques. The techniques, such as Mantel-Haenszel (MH) and Logistic Regression, indicate the effect size quantity for the calculated DIF value (Hambleton, 2006). Hence, the items in MH and Logistic Regression techniques can be grouped as those yielding weak (A), moderate (B) and high (C) DIF levels by making use of the effect size quantity. The final step of the analytical processes is about determining in favour or in disfavour of which group the DIF values prove to be. At this stage, the DIF type of the items is identified as uniform or non-uniform (Mellenbergh, 1982). In the uniform DIF-yielding items, the item functions favourable for either the reference or the focal group on all ability levels. However, the non-uniform DIF-yielding items, the item advantageous for the reference group on some ability levels, while it is favourable for the focal group on different ability levels.

It is rather challenging to make an interpretation about item bias and fairness of the test through statistical techniques (Zieky, 2002). An item might have yielded DIF for some reason other than item bias (Camilli & Shepard, 1994). Apart from the fact that a study on DIF is considered to be one of the evidence-gathering ways for the validity of the test, since such evidence have no single correct answer, an expert or a referee advice is required to evaluate and interpret this evidence (Benito et al., 2010).

Hambleton and Rogers (1995) emphasised that the sensitivity reviews could be beneficial if it is done from before performing a statistical analysis to determine whether or not the items have a structure that is in favour or disfavour of an aggressive, controversial and particular group. The aim of the sensitivity review is to reveal the source of DIF in the items after statistical analyses.

To determine the items, which contain the expressions that are likely to cause bias, the experts try to examine whether these items bear stereotyped expressions or not. Whether the content is unfavourable for the experiences of a given group or not in the sub-groups, they have equal chances regarding learning the substance of the item (Benito et al., 2010; Hambleton et al., 2001). In this evaluation, education and training of DIF and item bias can be provided to avoid any difference among the experts and elevate their adaptability (Hambleton, 2006). It is hard to predict bias, even though differences among the experts as regard to their training might have been made up (Gierl, Rogers & Klinger, 1999; Jensen, 1977; Plake, 1980; Sandoval & Miille, 1980).

There are rather few studies found in the literature to determine the items with DIF and to understand the sources of DIF by depending on the experts' view (Gierl et al., 1999; Roth, Oliveri, Sandilands, Lyons-Thomas & Ercikan, 2013). The studies regarding the experts' predictions without having any knowledge of statistic results and the results of the statistical DIF-determining techniques are rather insufficient.

In the same way, no study has been found in the literature as to the comparison of the predictions of the field experts, who lack sufficient level of theoretical knowledge even if they may have received training on DIF.

This research was designed for the purpose of eliminating such imperfections. The primary goal of the study is to compare the results of the experts' predictions on DIF and the statistical results. In line with this purpose, the explanations to the following questions were sought for:

1. How are the DIF level predictions of the field experts for the Sciences and Social Sciences tests?
2. How are the DIF predictions made through statistical techniques for the Sciences and Social Sciences tests?
3. How are the comparison results between the DIF level-predictions addressed by the field experts and statistical methods results for the Science and Social Sciences tests?

2. Method

2.1. Sample and data

This research has been designed as a descriptive study since it aimed at putting forward the compliance between the DIF predictions of the field experts and the DIF results calculated through the statistical techniques. The research population comprises 1,055,508 eighth-grade students who entered the placement test (PT, Turkish acronym is SBS) performed by the Ministry of National Education (MNE) of Turkey. 130,564 students selected from this population through the unbiased methods were incorporated into the sampling within the scope of the research. During the sampling process, 13% of the population were selected randomly (draw technique); yet, those who left with their gender unwritten and those with PT (SBS) scores proved to be zero were excluded from the selected sampling. As a result, the analyses were conducted with 130,564 students, 65,505 of whom were male, and 65,058 of whom were female. In the first stage of the research, the items and test scores within the PT Sciences and Social Sciences tests performed on the eighth-grade students by MNE in 2011 were practiced. These data were provided from MNE. The PT carried out for the eighth-grade students consists of five tests as Turkish, Mathematics, Sciences, Social Sciences and Foreign Languages. The entire test performed for the eighth-grade students is composed of a total of 100 questions comprising 23 questions in Turkish test, 20 in Math test, 20 in Sciences test, 20 questions in Social Sciences test and 17 questions in Foreign Languages test. The exam duration is

120 minutes. The study was implemented according to the data obtained from the Sciences and Social Sciences tests.

2.2. DIF-determining techniques (DIF procedures)

To determine DIF levels, PT Sciences and Social Sciences tests items statistical techniques and the techniques regarding the decisions made by an expert (Judgemental Technique) were used. To ascertain the DIF levels of the items with a statistical approach, the techniques referred to as 'MH' and 'indices of conditional p -value differences' were applied. MH is a chi-square statistic, and the obtained results can be interpreted as DIF in favour of the reference group if $MH > 1$; DIF in favour of the focal group if $MH < 1$; and no DIF if $MH \cong 1$. A logarithmic transformation is performed to be able to interpret the MH statistics more easily. The logarithmic transformation formula is as follows: $\Delta\Omega = \Delta MH = -(4/1.7) * \ln MH = -2.35 * \text{logit}$. The results obtained by the logarithmic conversion of the formula are interpreted as DIF in favour of the focal group if $\Delta MH > 0$; DIF in favour of the reference group if $\Delta MH < 0$; and no DIF if $\Delta MH \cong 0$ (Holland & Thayer, 1986).

Separately, the DIF level can also be interpreted according to the size/greatness of MH. It is stated that if $|\Delta MH| < 1$, then A indicates an insignificant DIF level; if $1 \leq |\Delta MH| < 1.5$, then B indicates a moderate DIF level, and if $|\Delta MH| \geq 1.5$, then C indicates a high DIF level (Dorans & Holland, 1993). One of the weakest aspects of this technique is that it cannot distinguish the uniform DIF and non-uniform DIF from one another. The second statistical technique referred to as the 'indices of conditional p -value differences', which is used in defining the item bias within the Sciences and Social Sciences tests, is also termed as the standardised differences, or merely, standardisation differences. In this technique, the test takers within the reference and focal groups were equalised according to the total test scores and score categories in the first place. Afterwards, the difference among the percentages of correct items of the focal and reference groups was taken and standardised for each equalised score. The positive values obtained from the analyses suggest that the DIF is in favour of the reference group, whereas the negative values indicate that the DIF is in favour of the focal group. If the DIF statistics obtained is within the range of ± 0.05 , then DIF is regarded as insignificant, and if it falls out of this range, then DIF is considered to be significant. Signed and Unsigned DIF (SDIF and UDIF) statistics can be calculated through the standardisation method. If the difference between the two statistics is small, then the existence of a uniform DIF is considered; whereas, if it is large, then a non-uniform DIF is mentioned. On the other hand, in the process of consulting the experts' decisions in determining DIF; DIF predictions of the field experts, which were made according to gender for Sciences and Social Sciences tests, were also collected. The volunteerism of all the experts for their participation in this research was taken as the cornerstone. The field experts consisted of the teachers who had completed the involved undergraduate program as well as those experienced in at least a 2-year teaching process. Seven field experts of Sciences (three males, four females) and five field experts of social sciences (two males, three females) have performed their DIF predictions.

First of all, during face-to-face interviews, the experts of Sciences and Social Sciences were informed as to what DIF was and how it could be identified, and the sample questions with weak, moderate and high DIF levels determined in different studies were shown to them. The field experts were given tests regarding their area, and then they were asked to fill in the form given to them by marking the DIF level of the items within the test as weak (A), moderate (B) and high (C).

The field experts were asked to use a list of nine items adapted from Hambleton and Rogers (1996) and Hambleton et al. (2001) regarding providing assistance for the DIF predictions to be performed according to gender. These items are as follows: 'There is some content likely to arouse different emotions or cause fallacy according to gender', 'It contains structure and language bearing vulgar and insulting characteristics according to gender'; 'There is some content showing difference according to gender'; 'It contains a structure/structures that the individuals may take advantage of in their lives according to their sexual identities'; 'It contains the information male/female students can benefit

from'; 'It contains words, structures or situations causing differences in meaning for female/male students'; 'For cultural reasons, the distractor(s) differ(s) according to gender'; 'The explanations given within the question may cause confusion in students' minds according to gender'; and 'It contains a hint/clue in the way that it will be of use to female/male students'.

The predictions of the field experts for the DIF level of each question (A-weak level or none; B-moderate level; C-high level) were tabulated within the Excel file. These predictions included in the classification scale were evaluated following the mode value. The research problems were tried to be solved by comparing the DIF predictions obtained from the field experts through statistical techniques.

3. Findings

In this study, 65,505 male and 65,088 female students answered the tests. When it is examined descriptive statistics in Sciences test, female's mean test scores are 9.04 and male's mean test scores are 8.18. For Social Sciences test, female's mean test scores are 9.74 and male's mean test scores are 8.68. It can be noted that female students are more successful than male students for both tests. When these means are compared in the Sciences test, the difference between female and male mean test scores is significantly important on behalf of the female students ($t = 30.827, p < 0.001$). The same results were also found for Social Sciences test ($t = 31.408, p < 0.001$). The estimations made by the field experts using a list of DIF indicators to evaluate the items in Sciences and Social Sciences test were collected as a reply to the first question of the study. The DIF estimations of the field experts on the items of the Sciences test are given in Table 1.

Table 1. The DIF estimations of the field experts on sciences test items

		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Field Experts	A1	1	1	1	1	1	1	2	1	1	2	1	1	1	1	1	2	1	1	1	1
	A2	1	1	1	2	1	1	1	1	1	2	1	1	1	1	1	2	2	1	1	1
	A3	1	1	1	2	2	1	2	1	1	1	2	1	1	1	1	2	2	1	1	1
	A4	1	1	2	2	2	1	2	1	2	1	2	1	1	1	2	2	1	1	1	2
	A5	1	1	2	2	2	1	2	1	1	1	2	1	1	1	2	2	1	1	2	1
	A6	1	1	2	2	1	1	1	1	1	1	1	1	1	1	2	2	2	1	2	2
	A7	1	1	2	2	2	1	1	1	1	2	2	1	1	1	2	2	1	1	2	2
1	%	100	100	43	14	43	100	43	100	86	57	43	100	100	100	43	0	57	100	57	57
2	%	0	0	57	86	57	0	57	0	14	43	57	0	0	0	57	100	43	0	43	43
3	%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Field experts' answers; 1: Weak level (A), 2: Moderate level (B) and 3: High level (C).

It is observed in Table 1 that DIF level estimations of the field experts are highly close on many items for Sciences test. Percentage values show that number of items showing DIF at a weak level is much more than other DIF levels. For the 16th item, all of the field experts stated that it showed DIF at a medium level. The field experts did not find any of the items as being at a higher-level DIF. The field experts expressed that the finding of DIF in Sciences test items was due to the following reasons: 'It contains information that may be benefited by male/female students'; 'The distractor/s show differences in terms of cultural reasons according to gender'; and 'It contains structure/s that may be advantageous in their lives according to the gender identities of the individuals'. The agreement of the replies given by the seven field experts was checked with percentage agreement. The measured value was found as 0.7142.

The DIF estimations of the field experts on the items of the Social Sciences test are given in Table 2.

Table 2. The DIF estimations of the experts of the field on social sciences test items

		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Field Experts	S1	1	1	1	1	1	1	1	2	1	1	1	1	1	1	1	1	1	1	1	1
	S2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	S3	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	S4	1	1	1	1	1	1	1	1	1	2	1	1	1	1	1	1	1	1	1	1
	S5	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
1	%	100	100	100	100	100	100	100	80	100	80	100	100	100	100	100	100	100	100	100	100
2	%	0	0	0	0	0	0	0	20	0	20	0	0	0	0	0	0	0	0	0	0
3	%	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Field experts' answers; 1: Weak level (A), 2: Moderate level (B) and 3: High level (C).

It is observed in Table 2 that DIF level estimations of the field experts, on the Social Sciences test, two of all items have a weak level of DIF. Only two experts expressed medium-level DIF for two different questions. According to the estimations made by experts using a list with DIF indicators, the experts stated that these two items were considered under the item of 'The distractor(s) show(s) differences regarding the cultural reasons according to gender'. The agreement of the replies given by the five experts was checked with percentage agreement. The measured value was found as 0.96.

The DIF levels were calculated for the Sciences and Social Sciences test items by using the MH and standardisation technique for the second question of the study. The DIF results calculated according to the mentioned methods and results are given in Table 3.

According to the results of the MH chi-square statistics computed for the Sciences test in Table 3, all the other items were found to be significant at 0.05 level, except for the 6, 11, 12 and 19. When the MH D-DIF results that were corrected by considering the size of the sampling are examined, it is observed that none of the items gave DIF at a significant level. According to the MH technique, all items have DIF at A-level, i.e., at a weak level.

The SDIF and UDIF values that were calculated with the standardisation technique are highly close to each other, and these values vary between -0.041 and 0.046. When the SDIF and UDIF values that were computed with the standardisation technique are examined, it is noted that all the items showed weak (A) level of DIF.

Table 3. The MH and standardisation DIF statistics calculated for gender for sciences and social sciences tests

Item	Alpha	Sciences				Social sciences				
		MH DIFL	MH D-DIF	Standardisation SDIF	Standardisation UDIF	Alpha	MH DIFL	MH D-DIF	Standardisation SDIF	Standardisation UDIF
1	1.154	A	-0.337	0.027	0.028	0.923	A	0.188	-0.014	-0.022
2	1.16	A	-0.349	0.027	0.028	0.971	A	0.07	-0.003	-0.01
3	0.919	A	0.199	-0.012	-0.015	1.027	A	-0.064	0.003	0.023
4	1.064	A	-0.146	0.01	0.012	1.022	A	-0.051	0.005	0.013
5	0.913	A	0.213	-0.016	-0.018	1.015	A	-0.035	0.002	0.016
6	1.001	A	-0.002	0.001	0.007	1.31	A	-0.635	0.04	0.04
7	1.07	A	-0.158	0.011	0.015	1.292	A	-0.602	0.038	0.038
8	1.306	A	-0.628	0.046	0.046	1.02	A	-0.047	0.003	0.011
9	0.939	A	0.147	-0.011	-0.021	0.904	A	0.238	-0.011	-0.02
10	0.74	A	0.707	-0.036	-0.041	0.947	A	0.128	-0.006	-0.022
11	0.976	A	0.056	-0.004	-0.012	0.907	A	0.228	-0.012	-0.015
12	0.979	A	0.05	-0.003	-0.012	0.995	A	0.011	-0.001	-0.011
13	1.081	A	-0.184	0.014	0.017	0.883	A	0.291	-0.022	-0.023
14	1.127	A	-0.28	0.02	0.024	1.107	A	-0.24	0.017	0.019
15	0.866	A	0.338	-0.02	-0.023	1.127	A	-0.28	0.017	0.017

16	0.814	A	0.483	-0.04	-0.04	0.963	A	0.089	-0.006	-0.019
17	0.949	A	0.124	-0.01	-0.017	1.055	A	-0.125	0.007	0.018
18	0.941	A	0.142	-0.011	-0.023	0.832	A	0.432	-0.018	-0.024
19	0.997	A	0.008	0	0.007	0.851	A	0.379	-0.017	-0.026
20	1.045	A	-0.102	0.008	0.013	0.854	A	0.372	-0.024	-0.028

DIFL = DIF level.

According to the results of the MH chi-square statistics computed for the Social Sciences test in Table 3, all the other items were found to be significant at 0.05 level, except for the 3, 4, 5, 8, and 12 items. When the MH D-DIF results that were adjusted by considering the size of the sampling are examined, it is seen that none of the items gave DIF at a significant level. According to the MH Technique, all items have DIF at A-level, i.e., at a weak level.

The SDIF and UDIF values that were calculated with the standardisation technique are highly close to each other, and these values vary between -0.028 and 0.040. When the SDIF and UDIF values that were measured with the standardisation technique are tested, it is noted that all the items showed weak (A) level of DIF.

The DIF level estimations made for the Sciences and Social Sciences test items to answer the third question of the study, and the DIF level estimations obtained from the statistical techniques are compared.

All the items that were calculated for the Sciences test and estimated according to the MH Statistical Method showed DIF at A-Level. Meanwhile, according to the UDIF values calculated with the Standardisation Technique for Sciences test, all the items also showed DIF at A-Level. When the experts of the field were asked to make estimations according to a list in which the DIF indicators are given, it was remarked that the majority of the DIF level of the items in the Sciences test was weak (at A-Level) for the most part. It is noted that the 16th item in the Sciences test shows Medium (B) level DIF according to all of the experts. According to 86% of the experts for the fourth item, 57% of the experts for the 3rd, 5th, 7th, 11th and 15th items also estimated as a Medium (B) level of DIF. According to these findings, the number of the items with DIF, which was discovered by statistical techniques, and the number of the items with DIF, which was based on the estimations of the experts of the field, are different. It is possible to claim that the accordance is at an acceptable level as a percentage.

All of the items that were calculated for the Social Sciences test and estimated according to the MH Statistical Method showed DIF at A-level. Meanwhile, according to the UDIF values calculated with the Standardisation Technique for Social Sciences test, all the items had also DIF at A-level. When the experts of the field were asked to make estimations according to a list in which the DIF indicators are given, it was observed that the majority of the DIF level of the items in the Social Sciences test was weak (at A-level) for the most part. It is apparent that only two items in this test showed DIF at A-level according to the 80% of the experts, and all the other items show DIF at A-level. According to these results, the number of the items with DIF, which was determined with statistical techniques, and the number of the items that were based on the estimations of the experts of the field are the same, and it is possible to claim that there is an acceptable agreement between them. On the other hand, it is possible to suggest that the agreement between the experts in Social Sciences test and the agreement between the experts and the statistical technique results is higher than that of the Sciences test.

4. Results

The estimations made by the experts were studied by using a list having DIF indicators in nine items to evaluate the Sciences and Social Sciences tests. It is noted that the majority of the items have DIF at a weak level, and the rest of the items have medium-level DIF. Meanwhile, according to the experts,

none of these items have DIF at a high level. The field experts stated that the finding of DIF in seven Sciences test items as medium-level DIF was due to the following reasons: 'It contains information that may be benefited by male/female students'; 'The distractor/s show differences in terms of cultural reasons according to gender' and 'It contains structure/s that may be advantageous in their lives according to the gender identities of the individuals'. Statistical values like Krippendorff Alpha and Fleiss Kappa were not obtained, and the percentage agreement was calculated instead of these values because the estimations of the experts were not very different. The experts did not prefer the expression 'It contains rude or insulting structures or language' as a DIF indicator in the list. The agreement of the replies given by the seven experts to the Sciences test items was reviewed with percentage agreement. The measured value was calculated as 0.7142.

According to the experts, all of the Social Sciences test items have a weak level of DIF. Twenty percent of the experts, who made estimations by using a DIF indicators list, affirmed that only two items had 'The distractor/s show differences in terms of cultural reasons according to gender'. Social Sciences class, which is a verbal subject, may be defined as a discipline that increases the readiness of the individuals in society. In this situation, it can be claimed that there were no items according to gender with DIF in Social Sciences Test in 2011 PT. The agreement of the responses given by the five experts of the field to the Social Sciences test items was checked with percentage agreement. The value was found as 0.96. This value shows that there is a remarkably high agreement between the raters.

The DIF levels were calculated for the Sciences and Social Sciences test items by using the MH and standardisation technique. The D-DIF results according to the MH Technique for Sciences and Social Sciences show that the items in these tests have DIF at a weak level. In a similar large-scale examination conducted in Turkey, there are items with DIF in Sciences test, and this has been demonstrated in the literature (Kalaycioglu & Kelecioğlu, 2011; Yurdugul & Askar, 2004). Similarly, when the calculations were made for both tests by using standardisation technique, the SDIF and UDIF values indicate that all the items show DIF at a weak level.

The DIF level estimations made by the experts for the items in Sciences and Social Sciences test and the DIF level estimations obtained from the statistical techniques were compared. All the items estimated according to the MH statistical method and standardisation technique for Sciences test showed DIF at a low level. When the experts were asked to make estimations according to a list, in which the DIF indicators were given, it was observed that more than 50% of the experts stated that the items number 3, 4, 5, 7, 11, 15 and 16 in the Sciences test, which showed DIF at a medium (B) level. It was asserted that the DIF level of the other items was weak. According to these results, although the number of the items with DIF, both regarding statistical techniques and experts' opinions are partly different, it is possible to claim that the agreement is at an acceptable level as percentage.

All the items estimated according to the MH statistical method and standardisation technique for Social Sciences test showed DIF at a weak level. When the experts were asked to make estimations according to the list in which the DIF indicators were given, it was observed that all of the items in Social Sciences test showed DIF at a weak (A) level. According to these results, the number of the items with DIF determined with statistical techniques and the estimations of the experts is equal. It is possible to suggest that there is a full agreement between them. It can be indicated that the agreement between the experts for the Social Sciences is higher than the agreement between the results of the experts' decisions and statistical techniques in Sciences test is greater.

Acknowledgements

This study has been realised in the process of 'Comparison of Classical Test Theory and Item Response Theory regarding Test Development Process', which is carried on by Nuri Dogan and is supported in the scope of TUBITAK BİDEB 2219 'Post-doctorate Research Scholarship Program Abroad' with the number 1059B191400868.

References

- Benito, J. G., Hidalgo, M. D. & Guilera, G. (2010). Bias in measurement instruments: fair tests. *Papeles del Psicologo*, 31(1), 75–84.
- Budgell, G. R., Raju, N. S. & Quartetti, D. A. (1995). Analysis of differential item functioning in translated assessment instruments. *Applied Psychological Measurement*, 1(4), 309–321.
- Camilli, G. & Shepard, L. A. (1994). *Methods for identifying biased test items*. Vol. 4, Thousand Oaks, CA: Sage Publications.
- Cromwell, S. (2002). *Primer on ways to explore item bias*. Paper presented at the 25th Annual Meeting of the Southwest Educational Research Association, 14–16 February, Austin, TX.
- Dorans, N. J. (2013). ETS contributions to the quantitative assessment of item, test, and score fairness. *Educational Testing Service*. Retrieved from <http://www.ets.org/Media/Research/pdf/RR-13-27.pdf>
- Dorans, N. J. & Holland, P. W. (1993). DIF detection and description: Mantel Haenszel and standardization. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). New Jersey: Educational Testing Service. (2009). Guidelines for fairness review of assessment. Retrieved from http://www.ets.org/Media/About_ETS/pdf/overview.pdf
- Gierl, M.J., Rogers, W.T. & Klinger, D.A. (1999). Using statistical and judgmental reviews to identify and interpret translation differential item functioning. *The Alberta Journal of Educational Rese*, 45(4), 353–376.
- Hambleton, R. K. (2006). Good practices for identifying differential item functioning. *Medical Care*, 44(11), 182–188.
- Hambleton, R. K. & Jones, R. W. (1992). *Comparison of empirical and judgmental methods for detecting differential item functioning*. Paper presented at the Annual Meeting of the National Council on Measurement in Education, 21–23 April, San Francisco, CA.
- Hambleton, R. K. & Rogers, H. J. (1989). Detecting potentially biased test items: comparison of IRT area and Mantel-Haenszel methods. *Applied Measurement in Education*, 2(4), 313–333.
- Hambleton, R. K. & Rogers, H. J. (1995). Item bias review. *Practical Assessment, Research and Evaluation*. Retrieved from <http://pareonline.net/getvn.asp?v=4&n=6>
- Hambleton, R. K. & Rogers, H. J. (1996). Developing an item bias review form. Retrieved from <http://ericae.net/ft/tamu/biaspub2.htm>
- Hambleton, R. K., Ying, L. & Klauck, S. (2001). Pennsylvania department of education item sensitivity review procedure (Unpublished document).
- Holland, P.W. & Thayer, D.T. (1986). *Differential item performance and the Mantel-Haenszel procedure* (Technical Report No. 86–69). Princeton, NJ: Educational Testing Service.
- Jensen, A. R. (1977). An examination of cultural bias in the Wonderlic Personnel Test. *Intelligence*, 1, 51–64.
- Jones, R. W. & Hambleton, R. K. (1992). Recent advances in psychometric methods. *Laboratory of Psychometric and Evaluative Research Report No. 233*. Amherst, MA: University of Massachusetts, School of Education.
- Kalaycioglu, B. D. & Kelecioğlu, H. (2011). Ogrenci Secme Sinavinin made yanliligi acisindan incelenmesi. *Egitim ve Bilim*, 36(161), 3–12.
- Mellenbergh, G. (1982). Contingency table models for assessing item bias. *Journal of Educational Statistics*, 32(1), 92–109.
- Plake, B. S. (1980). A comparison of statistical and subjective procedures to ascertain item validity: one step in the test validation process. *Educational and Psychological Measurement*, 40, 397–404.
- Roth, W.-M., Oliveri, M. E., Sandilands, D., Lyons-Thomas, J. & Ercikan, K. (2013). Investigating sources of differential item functioning using expert think-aloud protocols. *International Journal of Science Education*, 35, 546–576.
- Sandoval, J. & Miille, W. P. W. (1980). Accuracy of judgments of WISC-R item difficulty for minority groups. *Journal of Consulting and Clinical Psychology*, 48, 249–253.
- Sireci, S. G. & Mullane, L. A. (1994). Evaluating test fairness in licensure testing: the sensitivity review process. *CLEAR Exam Review*, 5(2), 22–27.
- Yurdugul, H. & Askar P. (2004). Ortaogretim Kurumlari Ogrenci Secme ve Yerlestirme Sinavinin cinsiyete gore made yanliligi acisindan incelenmesi [The investigation of the student selection and placement

Yavuz, S., Dogan, N., Hambleton, R. K. & Yurtcu, M. (2018). The comparison of differential item functioning predicted through experts and statistical techniques. *Cypriot Journal of Educational Science*. 13(2), 375-384.

examination for secondary education with respect to gender in terms of item bias]. *Egitim Bilimleri ve Uygulama Dergisi*, 3(5), 3–20.

Zieky, M. (2002). Ensuring the fairness of licensing tests. *CLEAR Exam Review*, 12(1), 20–26.

Zumbo, B. D. (1999). *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, Canada: Directorate of Human Resources Research and Evaluation National Defense Headquarters.