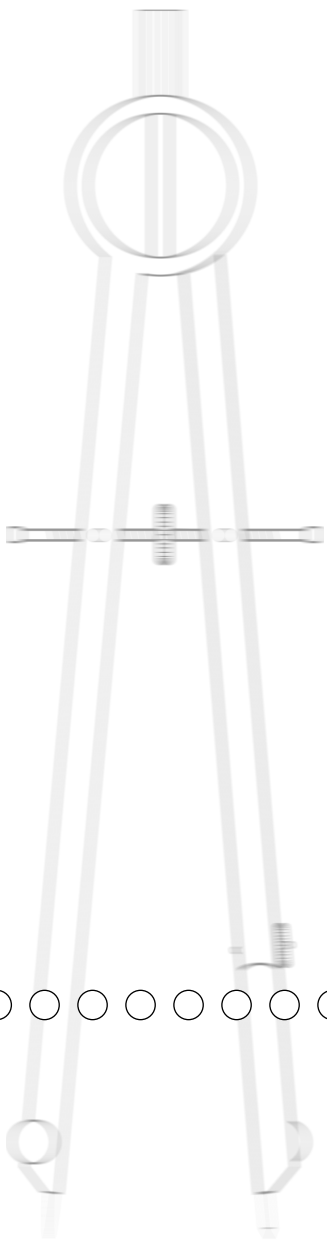


Item Analysis and Review

Ina V.S. Mullis
Michael O. Martin







13

Item Analysis and Review

Ina V.S. Mullis
Michael O. Martin

13.1 Overview

In order to assess the psychometric properties of the TIMSS 1999 achievement items before proceeding with item response theory (IRT) scaling,¹ TIMSS computed a series of diagnostic statistics for each item in each country. As part of the TIMSS quality assurance process, these statistics were carefully checked for any evidence of unusual item behavior. If an item was uncharacteristically easy or difficult for a particular country, or had unusually low discriminating power, this sometimes suggested a translation or printing problem. On the relatively few occasions that such items were found, the test booklets were examined for flaws, and where necessary the national research coordinator was consulted. Any item that was discovered to have a flaw in a particular country was removed from the database for that country.

13.2 Statistics for Item Analysis

The basic statistics for the item review were calculated at the IEA Data Processing Center and summarized in graphical form for review at the International Study Center. Item statistics were computed for each of the 38 TIMSS 1999 countries, and where countries tested in more than one language, for each of the languages tested. For each item, the basic item-analysis display presents the number of students that responded in each country, the difficulty level (the percentage of students that answered the item correctly), and the discrimination index (the point-biserial correlation between success on the item and a total score).² For multiple-choice items (see Exhibit 13.1 for an example), the display presents the percentage of students that chose each option, including the percentage that omitted or did not reach the item, and the point-biserial correlation between each option and the total score. For free-response items (which could have more than one score level – see Exhibit 13.2 for an example), the display

○○○

1. The TIMSS 1999 IRT scaling is described in Chapter 14.
2. For the purpose of computing the discrimination index, the total score was the percentage of items a student answered correctly.

presents the difficulty and discrimination of each score level. As a prelude to the main IRT scaling, it shows some statistics from a preliminary Rasch analysis, including the Rasch item difficulty for each item and the standard error of this difficulty estimate.

The item-analysis display presents the difficulty level of each item separately for male and female students. As a guide to the overall statistical properties of the item, it also shows the international item difficulty (the mean of the item difficulties across countries) and the international item discrimination (the mean of the item discriminations).

Exhibit 13.2 International Item Statistics for a Free-Response Item

Country	N	Diff	Disc	Pct_L0	Pct_L1	Pct_L2	Pct_L3	Pct_OM	Pct_NR	Pct_O	Pct_B1	Pct_B2	Pct_B3	Pct_OM	RDIFF	Cases	Reliability Score	Code	Flags
Australia	494	23.30	0.39	62.30	23.30	.	14.40	0.00	0.00	-0.22	0.39	.	.	-0.16	2.10	123.00	96.70	92.70_H	----
Belgium (Flem)	656	67.20	0.40	27.70	67.20	.	5.00	0.00	0.00	-0.28	0.40	.	.	-0.39	0.35	90.00	100.00	98.90	-----
Bulgaria	411	68.60	0.52	20.00	68.60	.	11.40	0.20	0.20	-0.29	0.52	.	.	-0.39	-0.60	88.00	100.00	97.70	-----
Canada (Englis)	783	32.30	0.42	57.00	32.30	.	10.70	0.10	0.10	-0.33	0.42	.	.	-0.10	1.36	.	.	.	-----
Canada (Fren)	302	33.80	0.32	56.00	33.80	.	10.30	0.00	0.00	-0.26	0.32	.	.	-0.07	1.44	.	.	.	-----
Chile (7th Gr)	723	26.40	0.48	51.50	26.40	.	22.10	0.40	0.40	-0.20	0.48	.	.	-0.27	0.26	.	.	.	-----
Chile (7th Gr)	754	12.60	0.35	56.60	12.60	.	30.80	0.70	0.30	-0.03	0.35	.	.	-0.20	0.89	.	.	.	-----
Chinese Taipei	725	83.70	0.64	13.40	83.70	.	2.90	0.00	0.00	-0.53	0.64	.	.	-0.33	-0.83	180.00	100.00	100.00	-----
Cyprus	388	21.40	0.49	54.10	21.40	.	24.50	0.00	0.00	-0.14	0.49	.	.	-0.31	1.51	.	.	.	-----
Czech Republi	406	67.70	0.52	27.30	67.70	.	4.90	0.00	0.00	-0.48	0.52	.	.	-0.14	0.06	121.00	98.20	94.20	-----
Egypt	373	4.60	0.26	82.30	4.60	.	13.10	0.00	0.00	-0.03	0.26	.	.	-0.13	4.12	89.00	100.00	93.30	-----
Finland (Fin)	344	16.00	0.31	70.90	16.00	.	13.10	0.30	0.30	-0.16	0.31	.	.	-0.09	2.42	.	.	.	-----
Finland (Swe)	44	4.50	0.12	61.40	4.50	.	34.10	0.00	0.00	0.05	0.12	.	.	-0.10	3.77	.	.	.	-----
Hong Kong (S)	648	81.60	0.39	14.50	81.60	.	3.90	0.20	0.20	-0.23	0.39	.	.	-0.34	-0.44	171.00	98.80	97.70	-----
Hungary	400	68.00	0.46	25.00	68.00	.	7.00	0.00	0.00	-0.30	0.46	.	.	-0.13	-0.13	98.00	98.00	96.90	-----
Indonesia	730	24.10	0.59	62.20	24.10	.	13.70	0.00	0.00	-0.43	0.59	.	.	-0.13	0.61	173.00	98.30	94.80	-----
Iran, Islamic R	868	35.30	0.38	56.70	35.30	.	7.90	0.00	0.00	-0.30	0.38	.	.	-0.13	0.03	170.00	92.90	73.50	-----
Israel (Arabic)	99	16.20	0.50	64.60	16.20	.	19.20	0.00	0.00	-0.19	0.50	.	.	-0.24	1.09	.	.	.	-----
Israel (Hebrew)	409	24.00	0.39	46.50	24.00	.	29.60	0.70	0.10	-0.10	0.39	.	.	-0.25	1.36	.	.	.	-----
Italy	419	62.10	0.52	24.60	62.10	.	13.40	0.20	0.20	-0.35	0.52	.	.	-0.29	-0.61	115.00	99.10	96.50	-----
Japan	589	79.30	0.47	15.60	79.30	.	5.10	0.00	0.00	-0.34	0.47	.	.	-0.29	-0.32	140.00	100.00	99.30	-----
Jordan	635	32.00	0.59	57.00	32.00	.	11.00	0.00	0.00	-0.41	0.59	.	.	-0.23	0.31	153.00	100.00	96.10	-----
Korea, Rep. of	765	80.70	0.56	14.80	80.70	.	4.60	0.00	0.00	-0.40	0.56	.	.	-0.37	-0.28	123.00	97.60	95.10	-----
Latvia (USS)	355	67.60	0.44	22.50	67.60	.	9.90	0.60	0.60	-0.28	0.44	.	.	-0.29	-0.56	97.00	97.90	95.90	-----
Lithuania	287	65.50	0.43	26.80	65.50	.	7.70	0.30	0.30	-0.28	0.43	.	.	-0.31	-0.57	68.00	98.50	98.50	-----
Mexiconia (A	119	8.40	0.10	43.70	8.40	.	47.90	0.80	0.20	0.22	0.10	.	.	-0.27	1.56	.	.	.	-----
Mexiconia (M	375	28.00	0.45	44.50	28.00	.	27.50	0.50	0.50	-0.08	0.45	.	.	-0.36	0.93	.	.	.	-----
Malaysia	699	54.80	0.61	36.80	54.80	.	8.40	0.10	0.10	-0.45	0.61	.	.	-0.31	0.37	191.00	100.00	99.00	-----
Madova (Ror	378	55.60	0.53	33.30	55.60	.	11.10	0.30	0.30	-0.38	0.53	.	.	-0.26	-0.34	.	.	.	-----
Moldova (Rus	70	64.30	0.51	28.60	64.30	.	7.10	1.40	1.40	-0.36	0.51	.	.	-0.34	-0.29	.	.	.	-----
Morocco	659	4.60	0.09	61.90	4.60	.	33.50	0.90	0.90	0.06	0.09	.	.	-0.08	1.46	108.00	95.40	88.90	-----
Netherlands	370	13.20	0.35	76.40	13.20	.	8.40	0.30	0.30	-0.16	0.35	.	.	-0.16	3.21	38.00	100.00	92.10	-----
New Zealand	456	9.90	0.31	73.50	9.90	.	16.70	0.00	0.00	-0.04	0.31	.	.	-0.20	2.80	114.00	98.10	95.60	-----
Philippines (E	754	12.20	0.41	77.70	12.20	.	10.10	0.40	0.40	-0.24	0.41	.	.	-0.11	0.84	.	.	.	-----
Philippines (T	79	2.50	0.26	86.10	2.50	.	11.40	0.00	0.01	0.01	0.26	.	.	-0.14	2.30	.	.	.	-----
Romania (Hun	15	53.30	0.72	40.00	53.30	.	6.70	0.00	0.00	-0.63	0.69	.	.	-0.15	-0.78	.	.	.	-----
Romania (Ron	418	68.40	0.58	21.80	68.40	.	8.90	0.50	0.50	-0.35	0.58	.	.	-0.41	-0.93	.	.	.	-----
Russian Feder	527	73.60	0.50	19.70	73.60	.	6.60	0.00	0.00	-0.37	0.50	.	.	-0.29	-0.61	120.00	98.30	96.70	-----
Singapore	625	83.70	0.48	13.90	83.70	.	2.40	0.00	0.00	-0.41	0.48	.	.	-0.23	-0.42	100.00	98.00	98.00	-----
Slovak Repub	430	73.30	0.42	20.90	73.30	.	5.80	0.00	0.00	-0.32	0.42	.	.	-0.23	-0.15	101.00	99.00	99.00	-----
Slovenia	386	74.10	0.43	22.30	74.10	.	3.60	0.00	0.00	-0.33	0.43	.	.	-0.26	-0.67	104.00	100.00	95.20	-----
South Africa (170	5.90	0.42	85.30	5.90	.	8.80	0.00	0.00	-0.22	0.42	.	.	-0.07	2.24	.	.	.	-----
South Africa (851	4.20	0.32	85.50	4.20	.	7.30	0.20	0.20	-0.19	0.32	.	.	0.00	1.45	.	.	.	-----
Thailand	711	41.90	0.59	51.10	41.90	.	7.00	0.10	0.10	-0.52	0.59	.	.	-0.12	0.28	207.00	100.00	100.00	-----
Tunisia	625	26.60	0.36	46.30	26.60	.	13.00	0.20	0.20	-0.21	0.36	.	.	-0.14	0.65	157.00	98.40	96.80	-----
Turkey	984	36.60	0.56	46.20	36.60	.	23.00	0.20	0.20	-0.38	0.56	.	.	-0.26	-0.17	247.00	100.00	98.80	-----
United States	1128	37.40	0.52	55.30	37.40	.	7.30	0.10	0.10	-0.41	0.52	.	.	-0.18	0.77	118.00	98.20	97.50	-----
International Avg.:		48.80	0.46	39.60	48.80	.	10.60	0.10	0.10	-0.30	0.46	.	.	-0.24	0.42	.	.	.	-----

Keys: Diff = Percent obtaining maximum score; RDIFF = Difficulty (1-P₀); Pct_In = Invalid Responses; Pct_NR = Not Reached; Pct_OM = Omitted
 Flags: A = Ability not ordered; Attractive Distractor; C = Difficulty less than chance; D = Negative/low discrimination; E = Easier than average;
 F = Distractor chosen by less than 10%; H = Harder than average; R = Scoring reliability < 80%; V = Difficulty greater than 95.

Exhibits 13.1 and 13.2 contained the statistics described below.

N: This is the number of students to whom the item was administered. If an item was not reached by a student it was considered to be not administered for the purpose of the item analysis.³

Diff: The item difficulty was the percentage of students that provided a fully correct response to the item. In the case of free-response items worth more than one point this was the percentage of students achieving the maximum score on the item. When computing this statistic, not reached items were treated as not administered.

Disc: The item discrimination was computed as the correlation between correctly answering the item and total score on all of the items in the subject area in the test booklet.⁴ This correlation should be moderately positive for items with good measurement properties.

PCT_A, PCT_B, PCT_C, PCT_D and PCT_E: Used for multiple-choice items only (Exhibit 13.1), these represent the percentage of students choosing each response option for the item. Not reached items were excluded from the denominator for these calculations.

PCT_0, PCT_1, PCT_2 and PCT_3: Used for open-ended items only (Exhibit 13.2), these are the percentage of students scoring at each score level for the item. Not reached items were excluded from the denominator for these calculations.

PCT_IN: Used for multiple-choice items only, this was the percentage of students that provided an invalid response to a multiple-choice item. Invalid responses were generally the result of choosing more than one response option.

PCT_OM: This is the percentage of students that did not provide a response to the item even though the item was administered and they had reached it. Not reached items were excluded from the denominator when calculating this statistic.

○○○

3. In TIMSS, for the purposes of item analysis and item parameter estimation in scaling, items not reached by a student were treated as if they had not been administered. For purposes of estimating student proficiency, however, not reached items were treated as incorrectly answered.
4. For free-response items, the discrimination is the correlation between the number of score points and total score.

PCT_NR: This is the percentage of student that did not reach the item. An item was coded as not reached when there was no evidence of a response to any of the items following it in the booklet and the response to the item preceding it was omitted.

PB_A, PB_B, PB_C, PB_D and PB_E: Used for multiple-choice items only, these present the correlation between choosing each of the response options A, B, C, D, or E and the score on the test booklet. Items with good psychometric properties have zero or negative correlations for the distracter options (the incorrect options) and moderately positive correlations for the correct answer.

PB_0, PB_1, PB_2 and PB_3: Used for free-response items only, these present the correlation between the score levels on the item (0,1,2, or 3) and the score on the test booklet. For items with good measurement properties the correlation coefficients should change from negative to positive as the score on the item increases.

PB_OM: This is the correlation between a binary variable - indicating an omitted response to the item - and the score on the test booklet. This correlation should be negative or near zero.

PB_IN: Used for multiple-choice items only, this presents the correlation between an invalid response to the item (usually caused by selecting more than one response option) and the score on the test booklet. This correlation also should be negative or near zero.

RDIFF: This is an estimate of the difficulty item based on a Rasch one-parameter IRT model. The difficulty of the item is expressed in the logit metric (with a positive logit indicating a difficult item) and was scaled so that the average Rasch item difficulty was zero within each country.

Reliability - Cases: It was expected that the free-response items in approximately one-quarter of the test booklets would be scored by two independent scorers. This column indicates the number of times each item was double scored in a country.

Reliability - Score: This column contains the percentage of times the two independent scorers agreed on the score level for the item.

Reliability - Code: This column contains the percentage of times the two scorers agreed on the two-digit code (score and diagnostic code) for the item.

As an aid to reviewers, the item-analysis display includes a series of “flags” signaling the presence of one or more conditions that might indicate a problem with an item. The following conditions are flagged:

- Item difficulty exceeds 95% in the sample as a whole
- Item difficulty is less than 25% for 4-option multiple-choice items in the sample as a whole (20% for 5-option items)
- Item difficulty exceeds 95% or is less than 25% (20% for 5-option items)
- One or more of the distracter percentages is less than 5%
- One or more of the distracter percentages is greater than the percentage for the correct answer
- Point-biserial correlation for one or more of the distracters exceeds zero
- Item discrimination (i.e., the point-biserial for the correct answer) is less than 0.2
- Item discrimination does not increase with each score level (for an item with more than one score level)
- Rasch goodness-of-fit index is less than 0.88 or greater than 1.12
- Difficulty levels on the item differ significantly for males and females
- Difference in item difficulty levels between males and females diverge significantly from the average difference between males and females across all the items making up the total score

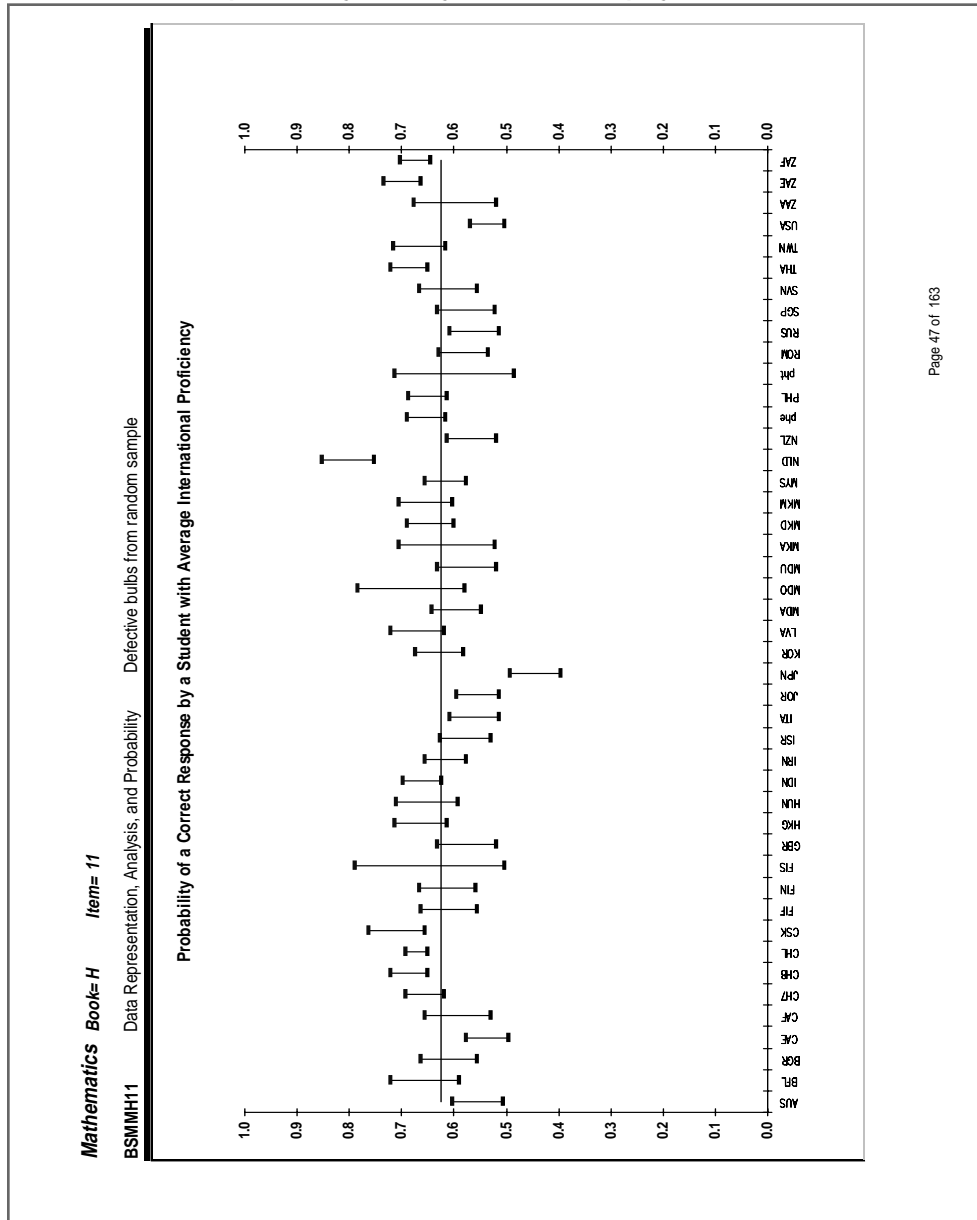
Although not all of these conditions necessarily indicate a problem, the flags are a useful way to draw attention to potential sources of concern. The IEA Data Processing Center also produced information about the inter-rater agreement for the free-response items.

13.2.1 Item-by-Country Interaction

Although there is room for variation across items, in general countries with high average performance on the achievement tests as a whole should perform relatively well on each of the items, and low-scoring countries should do less well on each of items. When this does not occur, i.e., when a high-scoring country has low performance on an item on which other countries are

doing well, there is said to be an item-by-country interaction. Since large item-by-country interactions can indicate an item that is flawed in some way, the item review also included this aspect of item performance.

Exhibit 13.3 Example Item-by-Country Interaction Display



To help examine item-by-country interactions, the International Study Center produced a graphical display for each item showing the average probability across all countries of a correct response for a student of average international proficiency, compared with the probability of a correct response by a student of average proficiency in each country (see Exhibit 13.3 for an example). The probability for each country is presented as a 95% confidence interval, which includes a built-in Bonferroni correction for multiple comparisons.

The limits for the confidence interval are computed as follows:

$$UpperLimit = \left(1 - \frac{e^{-\frac{RDIFF_{ik} + SE_{RDIFF_{ik}} \times Z_b}{1 + e^{-\frac{RDIFF_{ik} + SE_{RDIFF_{ik}} \times Z_b}}}}}{1 + e^{-\frac{RDIFF_{ik} + SE_{RDIFF_{ik}} \times Z_b}}}} \right)$$

$$LowerLimit = \left(1 - \frac{e^{-\frac{RDIFF_{ik} - SE_{RDIFF_{ik}} \times Z_b}{1 + e^{-\frac{RDIFF_{ik} - SE_{RDIFF_{ik}} \times Z_b}}}}}{1 + e^{-\frac{RDIFF_{ik} - SE_{RDIFF_{ik}} \times Z_b}}}} \right)$$

where $RDIFF_{ik}$ is the Rasch difficulty of item k within country i ; $SE_{RDIFF_{ik}}$ is the standard error of the difficulty of item k in country i ; and Z_b is the critical value from the Z distribution, corrected for multiple comparisons using the Bonferroni procedure.

13.3 Item Checking Procedures

Prior to the IRT scaling of the TIMSS 1999 achievement data by Educational Testing Service, the International Study Center thoroughly reviewed the item statistics for all participating countries to ensure that items were performing comparably across countries. Although only a small number of items were found to be inappropriate for international comparisons, throughout the series of item-checking steps a number of reasons were discovered for differences in items across countries. Most of these were inadvertent changes in the items during printing, such as omitting an item option or misprinting the graphics associated with an item. Differences attributable to translation problems, however, were found for an item or two in several countries.

In particular, items with the following problems were considered for possible deletion from the international database:

- Errors were detected during translation verification but were not corrected before test administration

- Data cleaning revealed more or fewer options than in the original version of the item
- The item-analysis information showed the item to have a negative biserial
- The item-by-country interaction results showed a very large negative interaction for a given country
- The item-fit statistic indicated that the item did not fit the model
- For free-response items, the within-country scoring reliability data showed an agreement of less than 70% for the score level. Also, performance in items with more than one score level was not ordered by score, or correct levels were associated with negative point-biserials.

When the item statistics indicated a problem with an item, the documentation from the translation verification⁵ was used as an aid in checking the test booklets and contacting National Research Coordinators (NRCs). If a problem could be detected by the International Study Center (such as a negative point-biserial for a correct answer or too few options for the multiple-choice questions), the item was deleted from the international scaling. If there was a question about potential translation or cultural issues, however, then the NRC was consulted before deciding how the item should be treated. Appendix D provides a list of deleted items as well as a list of recodes made to free-response item codes.

13.4 Summary

Considering that the checking involved more than 300 items for 38 countries (almost 12,000 item-country combinations), very few deviations from the international format were found. Appendix D summarizes the changes that were made to items in the international database before beginning the 1999 IRT scaling.

○○○

5. See Chapter 5 for a description of the translation and verification of the TIMSS data-collection instruments.