

An Application of Cognitive Diagnostic Assessment on TIMMS-2007 8th Grade Mathematics Items

Turker Toker, Kathy Green

Paper presented at the Annual Meeting of the American Educational Research Association (Vancouver, British Columbia, Canada, April 14-16, 2012)

ABSTRACT

The least squares distance method (LSDM) was used in a cognitive diagnostic analysis of TIMSS items administered to 4,498 8th-grade students from seven geographical regions of Turkey, extending analysis of attributes from content to process and skill attributes. Logit item positions were compared between data for Turkey and the full international data. The Pearson correlation between item logit positions was $r = .82$, though several items were not invariant across the datasets. Results indicated that the majority of the TIMSS items were well explained by a set of 20 attributes ($R^2 = .65$). This study provides an extension of work by Dimitrov (2007) in use of LSDM in cognitive diagnostic analysis.

Key words: cognitive assessment, mathematics, Least Squares Distance Method, validation, Rash model, psychometrics

Educational assessment has evolved from grading one's achievement level to being diagnostically useful at every step in education (Bolt, 2007). Assessment affects grades, placement, advancement, instruction, curriculum, and in some cases, funding. Assessment is critical to both evaluating the effects of educational programs and also to directing those programs. Outcomes of education need to meet globally accepted criteria. Students must be able to think wisely and critically, to examine in detail, and to make inferences (Gierl, 2007). Changes in the skills base and the knowledge our students need require new learning goals; these new learning goals change the relationship between assessment and instruction with teachers needing to be informed about both.

In this study, data were taken from one of the most respected international exams, the Trends in International Mathematics and Science Study (TIMSS-2007), used to assess cognitive abilities of 8th graders' mathematics achievement. The purpose of this study was to validate cognitive attributes on the released TIMSS-2007 mathematics test items with respect to the cognitive attributes developed by Tatsuoka and her associates and so to extend use of the Least Square Distance Method (LSDM) to skills and cognitive processes as well as knowledge attributes. The study was specific to attribute identification for Turkish students. To better understand the idea behind the model, a brief review of CDA is given below.

Cognitive Diagnostic Assessment (CDA)

One of the most important concerns in education systems is summative assessment. In order to have powerful, effective, and meaningful summative assessment, evaluation should also be formative, which means it has to support teaching and learning processes with results (DiBello & Stout, 2007). An ideal assessment would not only be able to meet precise psychometric standards, but would also be able to provide specific feedback about how students

learn and what attributes they need to achieve goals. Cognitive diagnostic assessment is used to examine the cognitive processes necessary to successful task completion, because it supplies specific information about each attribute students need to master instead of a single score result (McGlohen, 2004).

One of the first cognitive diagnostic models was Fischer's Linear Logistic Test Model (LLTM), created in 1972. The LLTM is an extension of the basic Rasch model to take into account the cognitive steps needed for correct item response. In the LLTM, Rasch item difficulty is computed as a sum of discrete cognitive attribute-based difficulties. The point of the LLTM that makes it appropriate for cognitive diagnostic assessment is that the item difficulty is the composite of the influences of the basic cognitive levels, or "factors," critical for properly solving an item (Fischer, 1973).

The rule space methodology (RSM) was developed by Tatsuoka and her associates (1983), and comprises two parts. The first part involves determining the relation between the items of a test and the attributes that they are assessing. The description of which attributes are essential for each item is shown in a Q-matrix. The Q-matrix is a $[K \times n]$ binary matrix, where K is the number of attributes to be measured and n is the number of items on the test. For a given element of the Q-matrix in the k th row and the i th column, a value of one indicates that item i does indeed measure attribute k and a zero indicates it does not. Figure 1 provides an example of a 3x3 Qmatrix:

		i1	i2	i3
Q =	k1	1	1	0
	k2	0	0	0
	k3	1	0	1

Figure 1. Sample Q-matrix.

The first item in the above Q-matrix shows that students who mastered the first and third attributes should have correct responses to item one. For a correct response to item two students should master the first attribute. A correct response to item three means they should have mastery of attribute three. Attributes which are not necessary for a correct response to the item are shown by zeros meaning that students do not need to master those attributes to correctly respond to that item. Expert consultation in the content area is a way to build the structure of the Q-matrix to determine if an item measures a specific attribute (see Appendix). Some other ways to build the structure of the Q-matrices can be borrowing from the test blueprint or subjectively evaluating each item to draw conclusions about which attributes are being measured (Tatsuoka, 1983). The second part of the rule space method deals with interpretation of the relationship between test items and the Q-matrix. The results of using the rule space method are to obtain diagnostic information on students' mastery levels of specified cognitive skills to improve interpretations of test scores on the item level for students, teachers, and administrators.

There are different cognitive assessment models in learning and teaching mathematics (Dogan & Tatsuoka, 2008; Duval, 2006; Kuchemann, 1981; Tatsuoka et al., 2004). Several studies have applied CDA to mathematics items. Their results are summarized in the following section.

Cognitive Diagnostic Assessment of Mathematic Items

Tatsuoka, Corter, and Tatsuoka (2004) examined TIMSS-R math items across 20 countries. Their results showed that high-achieving countries in the eighth-grade TIMSS- 99 mathematics assessment mostly had higher level thinking skills. Chen, Gorin, Thompson, and Tatsuoka (2006) conducted a cognitive diagnostic assessment of TIMSS-99 items. They used three analyses, including calculation of classification rates, multiple regression analyses, and

comparisons of attribute mastery probabilities across four booklets. In general, a list of cognitive attributes predicted the performance of Taiwanese eighth graders on the TIMSS-1999 mathematics tests very well.

Dogan and Tatsuoka (2008) examined Turkish students' mathematics performance on the TIMSS-R. Their study was conducted using the rule space method. They used a Q-matrix that included a set of 23 attributes. Results showed that when compared to American students, Turkish students were poor in mastering attributes such as P10 (quantitative reading), S4 (approximation/estimation), S6 (patterns and relationships), and S10 (solving open-ended problems).

Ma, Çetin, and Green (2009) took data from a 2005 Turkish national assessment of eighth grade students' performance in math to examine attribute function using the least squares distance method (Dimitrov, 2007). They found that many attributes predictive of Rasch model item difficulties were beyond the students' abilities. The given time for students to complete the 25- item test was short. Their results suggest that the students at lower ability levels might guess at some attributes. They also found item difficulties were well predicted from the set of attributes used.

Results of these studies are based on the subjective judgment of experts leading to the development of the Q-matrices. Generally speaking, studies that use the same data can yield different results because of use of different Q-matrices.

For the purpose of this study, a relatively new approach to cognitive diagnostic analysis called the Least Squares Distance Method (LSDM) was used to explore the validation of cognitive attributes on the released TIMSS-2007 mathematics test items. This approach has not been used previously with TIMSS items. The intent of this study was to extend use of the LSDM

to skills and cognitive processes as well as knowledge attributes, as identified by Turkish education experts.

Least Squares Distance Method

The least squares distance method (Dimitrov, 2007) uses the Rasch item position parameters of binary test items to (a) validate cognitive attributes that underlie item responses and (b) assess the probability of correct processing of such attributes across levels of the scale continuum. The LSDM is a conjunctive model in which a correct answer on a test item requires mastery of all cognitive attributes associated with that item. The cognitive attributes for all test items are outlined in a Q-Matrix, where a “1” shows an attribute is needed for an item and a “0” means an attribute is not needed. The basic assumption with LSDM is that, theoretically, the probability of correct item response is equal to the probability that all required attributes are correctly applied; which is,

$$P_{ij} = \prod_{k=1}^K [P(A_k = 1|\theta_i)]^{q_{jk}} \quad (1)$$

where P_{ij} is the probability of correct response on item j at ability level θ_i (item probability),

$P(A_k = 1|\theta_i)$ is the probability of correct response on attribute A_k at ability level θ_i (attribute probability),

and q_{jk} is a 0 or 1 element of the attribute matrix for item j and attribute A_k (Q-matrix).

Formula 2 is generated by taking the natural logarithm of both sides of Formula 1:

$$\ln P_{ij} = \sum_{k=1}^K q_{ik} \ln P(A_k = 1|\theta_i) \quad (2)$$

Then, Formula 2 is simplified to:

$$L = QX \quad (3)$$

where L is the vector with known elements $\ln P_{ij}$,

Q is the Q -matrix,

and X is the vector with unknown elements $\ln P(A_k = 1 | \theta_i)$.

According to Dimitrov, Equation 3 does not have an exact solution since it is

“overdetermined” —the number of equations $[i*j]$ is greater than the number of unknowns $[k*i]$.” To solve this problem, the LSDM is used to minimize the Euclidean norm of the vector $\|QX - L\|$. For a participant with ability level θ_i , the probability of a correct answer on attribute A_k is $P(A_k = 1 | \theta_i) = \exp(X_k)$. Item probabilities are recovered from the attribute probabilities $P(A_k = 1 | \theta_i)$ across ability levels. The graphical image of this probability across ability levels shows the probability curve for cognitive attribute A_k . The Rasch item characteristic curve (ICC) is represented by the recovered item probabilities (*Prec*), or the recovery curve. The ICC recovery compared to LSDM provides information about how well the required attributes describe the item across ability levels. The mean absolute difference (MAD) between the LSDM curve and the ICC provides for validation of attributes for each item across ability levels. A MAD equal to 0.0 would indicate perfect ICC recovery. According to Dimitrov (2007), a classification for level of ICC recovery was developed: “(a) very good ($0.00 \leq \text{MAD} < 0.02$), (b) good ($0.02 \leq \text{MAD} < 0.05$), (c) somewhat good ($0.05 \leq \text{MAD} < 0.10$), (d) somewhat poor ($0.10 \leq \text{MAD} < 0.15$), (e) poor ($0.15 \leq \text{MAD} < 0.20$), and (f) very poor ($\text{MAD} \geq 0.20$).” (p. 373). These criteria were applied in this study.

Dimitrov (2007) pointed out an interpretation of LSDM results with respect to heuristic criteria for validation of cognitive attributes: “(1) The smaller the LSD..., the better the

cognitive attributes hold together (jointly for all items) at this ability level; (2) The attribute probability curves (APCs) should exhibit logical and substantively meaningful behavior in terms of monotonicity, relative difficulty, and discrimination; (3) The better the ICC recovery for an item, the better the required attributes explain the item” (pp. 372-373).

Method

Participants

Released items and Turkish students’ responses to TIMSS-2007 in mathematics administered in 2007 were used in this study. There were 4,498 8th-grade students from seven geographical regions of Turkey. There were 14 booklets which were randomly assigned to students with between 314 and 331 students receiving each booklet. Booklets contained at least one released item. Since there are limited numbers of released items which were repeated in different booklets, all booklets were used in this study. The study used data from 2,093 female students with a mean age of 13.96. There were 2,405 male students with a mean age of 14.09.

A Q-matrix was developed by two Turkish speaking mathematics teachers and the first author of this paper. All of the experts had bachelor’s degrees in teaching and were male. Two of them worked as mathematics teachers in the U.S. The ages of experts were 27, 30 and 35. The first author of this paper had three years of teaching experience in Turkish schools. The other two experts had three and five years of teaching experience.

Instrument

The TIMSS-2007 for 8th grade consisted of 179 questions which included 96 multiple-choice questions. There were only 51 multiple-choice items released. No information was found about why TIMSS administrators released those items. Two of the items were dropped since they did not provide any variation at all in student responses (i.e., all students answered correctly or

all answered incorrectly). The final dataset was based on 49 items. (See Appendix for sample items.) For the Q-matrix, only released multiple-choice items were used. This test covered content domains of numbers, geometric shapes and measures, and data display. The cognitive domains included in the exam were knowing, reasoning, and applying. Two examples of items from this test can be seen in Figures 2 and 3.

Analysis

A Q-matrix of cognitive attributes (see Appendix and Table 1) of each item was developed based on Tatsuoka and her associates classification of cognitive attributes (K. Tatsuoka, Corter, & C. Tatsuoka, 2004). The Q-matrix includes 27 attributes divided into the three categories of content, process, and skill. Cognitive attributes which are represented by all items or not represented by all items were not used since they would provide no variation. The final version of the Q-matrix included 20 attributes. To ensure the consistency of the subjective judgment of attributes, the cognitive attribute matrix was developed based on the independent identifications of three experts.

A correct answer on an item means that a student has mastered all attributes required, within a margin of error. A linear regression analysis was conducted to see if the Q-matrix could explain item difficulty. Using the responses of Turkish students to the released math items of TIMSS-2007, an IRT analysis was conducted to both compare results with international item parameters and to use for LSDM analysis. Item parameters were correlated and a scatterplot constructed to identify items that had distinctly different positions for the Turkish and international data. Logit item difficulties for those items were compared using a significance test.

The probability of correct response was calculated via the WINSTEPS program for each item for Turkish students (Linacre, 2007). The individual attribute probabilities and the average

LSDs for the 49 items across ability levels were estimated using the MATLAB computer program (The MathWorks, Inc., 2005). Finally, the mean absolute difference between the ICC and the LSDM recovery curve for each item was calculated to validate the cognitive attributes.

Results

Item Parameter Comparison

It was hypothesized that the item parameters would be invariant across Turkish students and the international data. To test this hypothesis for the released 49 items, WINSTEPS results for Turkish students and item parameters from the official website of TIMSS were used. There were three items which were statistically significantly different at $p < .01$ in logit position: items 1 ($t = 5.69$), 25 ($t = 7.35$), and 37 ($t = 8.45$). These three items assessed basic concepts and operations in fractions and decimals (C2), computational applications of knowledge in arithmetic and geometry (P2), basic concepts and operations in elementary algebra (C3).

A Pearson correlation between these two sets of item logit positions produced a correlation of $r = .82$, and so strong consistency in item logit position. Both overall item parameters and Turkish students' item parameters are given in Table 1.

Cognitive Attributes Matrix

The overall agreement level among the three experts on the Q-matrix elements was 68%. After a final meeting, the judges agreed on the final version of the Q-matrix. Seven attributes (P5, P6, S4, S5, S8, and S10) which were in the original classification of cognitive attributes were deleted because no items required the attributes. One more attribute (S11) was excluded because it was present in all items. The last version of the cognitive attribute matrix had 6 content, 8 process, and 6 skill attributes, or 20 total attributes (see Table 1).

There were some statistically significant relationships between attributes. The relationship between attributes C1 (Basic concepts and operations in whole numbers and integers) and P9 (Management of data and procedures) was statistically significant ($r = .43, p < .05$). Additionally, the correlation between C1 (Basic concepts and operations in whole numbers and integers) and C4 (Basic concepts and operations in two-dimensional geometry), and this correlation was also significant was $r = .62, p < .05$. For attributes P7 (Generating, visualizing, and reading figures and graphs) and S3 (Using figures, tables, charts, and graphs) the correlation was high and statistically significant ($r = .89, p < .05$). From these results, we can conclude that some attributes overlap. We can also say in order to master one attribute, one needs to master another attribute.

A multiple regression analysis showed that most of the variance in item difficulties was accounted for by the identified cognitive attributes; however none of the individual attributes was statistically significant. The estimates of R^2 and adjusted R^2 were .65 and .42, respectively.

LSDM Analysis

The LSDM was conducted across 17 ability levels in the interval from -4.0 to 4.0 , with increments of 0.5 on the logit ability scale. The 20 attributes were applied more accurately and consistently by higher ability examinees. The LSDM estimates of the probabilities of correct performance on each attribute across ability levels are presented in Table 2. The attribute probability curves (APCs) monotonically increase across the ability levels and provide information about the relative difficulty and discrimination of the 20 attributes. C4 (basic concepts and operations in two-dimensional geometry) was the most difficult attribute because its APC was consistently below the other APCs across all ability levels. Other attributes can be

put in increasing difficulty order as follows: C1, C5, P1, P3, P8, P10, S1, S9, C3, C2, S2, P9, C6, S7, P2, P7, S6, P4, S3, and C4.

The APCs obtained with these probabilities are graphed in Figure 4. C4 was the most difficult attribute (the lowest curve), followed in decreasing difficulty by S3, P4, P2, P7, S6, S7, C6, P9, S2, and C2. All the other attributes have similar difficulty values. For example, the probability of correct performance on C4 at the ability level $y=0$ was .654 (see Table 2). That is, the likelihood of examinees with ability at the origin of the logit scale to correctly process geometrical operations in two dimensional geometry (C4) was .654. For the same examinees, the likelihood to correctly process algebra operations in elementary algebra (C3) was higher (.743).

For each item, the mean absolute differences for the ICC recovery across ability levels are given in Table 3. Graphically, the ICC recovery is presented (Figure 5) for four examples of items (47, 40, 19, and 17) with differing MAD levels. According to Dimitrov's criteria, item 47 is in the category of very good with a MAD of .017, item 40 is in the category of good with a MAD of 0.036, item 19 is in the category of somewhat good with a MAD of .060 and item 17 is in the category of somewhat poor with a MAD of .14. With the conventional rule for degree of ICC recovery described earlier in the LSDM section of this paper, the examination of all 49 graphs for ICC recovery and their MAD values revealed that the ICC recovery was very good for four items (20, 36, 41, and 47), good for 16 items (1, 2, 5, 11, 13, 16, 24, 25, 28, 31, 34, 35, 39, 40, 45, and 46), somewhat good for 15 items (4, 7, 10, 18, 19, 21, 22, 29, 30, 33, 37, 38, 42, 48, and 49), somewhat poor for eight items (6, 8, 9, 12, 17, 26, and 27), poor for four items (3, 14, 43, and 44) and very poor for two items (15 and 23). Generally speaking, such diagnostic information on ICC recovery can be particularly useful in validating math sub-skills for students.

Overall, the findings indicate that the 20 attributes relate to difficulties in mathematical skills of students (Mean MAD = .075). The mean MAD value suggests that overall item recovery is somewhat good based on the Dimitrov's criteria. For this reason, the APCs of the 20 attributes provide valuable information in terms of their difficulty (see Figure 4). But the results for ICC recovery suggest that there is room for improvement regarding the set of attributes and their links to items in the Q-matrix. Indeed, using Dimitrov's criteria compared to the MAD values of items in Table 3, 35 items have very good, good, or somewhat good ICC recovery while 14 items do not. According to these results, it might be said that the Q-matrix can be improved to get better results since 14 items were not well-recovered; for example, see the location of Items 47 and 17 in Figure 5.

Discussion

Cognitive diagnostic assessment is a useful way to examine validation of test items (Tatsuoka, Corter, & Tatsuoka, 2004). This study validated cognitive attributes using item position parameters from the Rasch model and attributes identified by three experts. A benefit of this approach is that student performance on individual attributes is obtained. With information about student performance on an individual attribute in hand, instruction can be tailored to the individual attribute level and then to overall success on item solution. With the LSDM, the present study investigated the validation of cognitive attributes on TIMSS-2007 items for Turkish students. First, the cognitive validation of the items was evaluated using the 27 initially-identified attributes. Once independent responses of experts were collected, a lack of variability was found resulting in deletion of seven of the attributes. In addition to this, two of the items were deleted due to a lack of response variability. The degree of the LSDM recovery of ICC was then assessed for the 49 items and item parameters for both international and Turkish students

were analyzed in order to assess the correlation between those parameters and so the generalizability of results.

The results showed the validation of cognitive attributes for the test with respect to the 20 revised attributes. The monotonic decrease of LSDs across ability levels demonstrated that it was unlikely that different cognitive strategies had been employed by students with different ability levels, and that the higher ability students applied the attributes more accurately and consistently. The APCs generally indicated attribute difficulties relative to each other and clear discriminations. The LSDM recovery of ICCs showed that on average, 71% of the items were recovered well by the attribute probabilities, revealing that the most of the 49 items could be substantially explained by the identified attributes. But the LSDM recovery also suggested that there is need for modifying attributes for some items such as Items 3, 6, 8, 9, 12, 14, 15, 17, 23, 26, 27, 43, and 44, because they were not explained very well by the attributes.

With respect to the three categories of attributes, the present research generated some useful results. It was found that identifying the content attributes was easier than identifying the cognitive process and skill attributes. The MAD values also showed that the content attributes accounted better for the items than the other two kinds of attributes and so, the LSDM recovery of ICCs was more accurate for content attributes. Although all attributes together contributed to the correct response for an item, their overlap might lead to unclear results. Some of the attributes include the same mastery areas. For example P7 and S3 overlapped on usage of figures and graphs. This shows that it might not be reasonable to combine different types of attributes in one analysis. The most difficult attribute for students to master was C4 (basic concepts and operations in two-dimensional geometry). This indicates that Turkish students' level of mastery is not sufficient in geometry.

Item logit positions did not differ significantly between the international and Turkish indices, with the exception of three items (items 1, 25, and 37).

To assess the stability of results, it is suggested that in future study an LSDM analysis be run with a random Q-matrix and with Q-matrices of multiple experts run separately. In addition to this, an aggregate Q-matrix can be iteratively revised for poor items to decrease the MAD and increase the multiple correlation between attributes and item difficulties. If an item could be written to uniquely address a specific attribute, these results would provide concurrent validation evidence and could be the best scenario to see how well students are doing on certain attributes. To clearly see the relationship between items and attributes, simple attributes are better than complex ones. It is feasible that items might be constructed for content attributes and less so for process and skill attributes, which are conveyed via a content item. Thus process and skill attributes would be overlaid on content.

This study was well developed with a large sample of students and the findings represented most of the released items in the mathematics test for the eighth grade in Turkey. But some caveats should be considered. The first limitation is the selection of the attributes. With respect to the cognitive model, if the same attributes are shown as relevant for two different items, the item difficulties of the two items should be close. Therefore, large differences in item difficulties would signify the misspecification or inadequacy of the chosen attributes for the items. Because of this, items 2 (logit position = .42) and 3 (logit position = -1.08), items 23 (logit position = -1.89) and 28 (logit position = .20), and items 27 (logit position = -1.45) and 33 (logit position = 1.13) were flagged for potential problems in the specification of the attributes since similar attributes were identified but the logit difficulties were very different.

Moreover, no guessing or omissions were taken into account in the LSDM (Dimitrov, 2007). In the current study, since the results from the IRT model showed that most of the items were beyond the students' abilities, making guessing likely. These problems with the LSDM were not taken into account in this study. Future studies might investigate effects of guessing with the LSDM.

References

- Bolt, D. (2007). The present and future of IRT-based cognitive diagnostic models (ICDMs) and related methods. *Journal of Educational Measurement*, 44(4), 377-383.
- Chen, Y., Gorin, J., Thompson, M., & Tatsuoka, K. (2006, April). *Verification of cognitive attributes required to solve the TIMSS-1999 mathematics items for Taiwanese students*. Paper presented at the annual meeting of the American Educational Research Association, San Francisco.
- DiBello, L. V., & Stout, W. (2007). Guest editors' introduction and overview: IRT-based cognitive diagnostic models and related methods. *Journal of Educational Measurement*, 44, 285–291.
- Dimitrov, D. M. (2007). Least squares distance method of cognitive validation and analysis for binary items using their item response theory parameters. *Applied Psychological Measurement*, 31(5), 367-387.
- Dogan, E., & Tatsuoka, K. (2008). An international comparison using a diagnostic testing model: Turkish students' profile of mathematical skills on TIMSS-R. *Educational Studies in Mathematics*, 68(3), 263-272.
- Duval, R. (2006). A cognitive analysis of problems of comprehension in a learning of mathematics. *Educational Studies in Mathematics*, 61, 103–131.
- Fischer, G. (1973). Linear logistic test model as an instrument in educational research. *Psychologica*, 37, 359-374.

- Gierl, M. J. (2007). Making diagnostic inferences about cognitive attributes using the rule space model and attribute hierarchy method. *Journal of Educational Measurement*, 44, 325–340.
- Kuchemann, D. (1981). Cognitive demands of secondary school mathematics items. *Educational Studies in Mathematics*, 12, 301-316.
- Linacre, J. M. (2007). Winsteps (Version 3.64.2) [Computer Software]. Chicago: Winsteps.com.
- Ma, L., Cetin, E., & Green, K. E. (2009, April). *Cognitive assessment in mathematics with the least squares distance method*. Paper presented at the annual meeting of the American Educational Research Association, San Diego, US.
- The MathWorks, Inc. (2005). MATLAB (Version 7.1) [Computer software]. Natick, MA: The MathWorks, Inc.
- McGlohen, M. (2004). *The application of cognitive diagnosis and computerized adaptive testing to a large-scale assessment*. Unpublished Doctoral Dissertation, University of Texas at Austin, TX.
- Tatsuoka, K. K. (1983). Rule space: An approach for dealing with misconceptions based on item response theory. *Journal of Educational Measurement* 20(4), 901-926.
- Tatsuoka, K. K., Corter, J. E., & Tatsuoka, C. (2004). Patterns of diagnosed mathematical content and process skills in TIMSS-R across a sample of 20 countries. *American Educational Research Journal*, 41(4), 901-926.

Wiggins, G. (1998). *Educative assessment: Designing assessments to inform and improve student performance*. San Francisco: Jossey-Bass.

Table 1
Item Difficulties and the 20 Initial Attributes for the 49 Items

Item	Item Difficulty International	Item Difficulty Turkey	Content Attributes						Cognitive Process Attributes							Skill Attributes						
			C1	C2	C3	C4	C5	C6	P1	P2	P3	P4	P7	P8	P9	P10	S1	S2	S3	S6	S7	S9
1*	0.12	-0.45	1	1	0	0	0	1	0	1	1	0	0	0	1	1	0	0	1	1	1	1
2	0.95	0.42	1	1	1	0	0	0	1	1	1	1	1	1	0	1	0	1	1	0	0	1
3	-0.35	-1.08	1	1	1	0	0	0	1	1	1	0	1	1	0	0	1	1	1	1	0	1
4	0.55	-0.17	1	1	1	0	1	0	1	1	0	1	1	1	0	0	1	0	0	0	0	1
5	0.63	0.63	0	0	1	1	0	1	0	1	1	0	0	0	1	1	1	1	0	0	1	1
6	0.097	-0.4	1	1	1	0	0	0	1	1	1	1	1	1	0	1	0	1	1	0	0	0
7	-	-0.95	1	0	1	0	0	0	0	0	0	0	1	0	0	1	0	0	0	1	0	0
8	-	-0.56	0	1	1	0	0	0	0	1	0	1	0	0	0	0	0	1	0	1	0	0
9	-	0.99	0	1	1	0	0	0	1	1	0	0	1	1	0	0	0	1	1	0	0	1
10	-	0.70	1	1	1	0	0	0	0	1	0	0	1	1	0	0	0	0	1	0	0	1
11	-	-0.59	0	0	1	0	0	0	0	1	0	0	0	0	0	1	0	0	1	0	0	0
12	-	1.09	0	0	0	1	0	0	0	0	1	0	1	0	1	0	0	0	1	0	1	0
13	-	0.14	0	0	1	1	0	1	1	1	1	0	0	0	1	1	0	1	0	0	1	1
14	-	-0.86	0	0	0	1	0	0	0	0	1	1	1	0	1	0	0	1	1	0	1	0
15	-	-1.89	1	0	1	0	1	0	0	1	1	0	1	0	1	0	0	1	0	1	1	0
16	-0.19	0.26	0	1	1	0	1	0	0	0	1	1	1	0	0	1	1	0	0	0	1	0
17	-0.50	-1.11	1	0	1	0	1	0	0	0	1	0	1	0	0	1	1	1	1	0	0	1
18	-0.49	0.01	0	1	1	0	0	0	0	1	1	1	1	0	0	0	0	1	0	1	0	0
19	0.70	0.68	0	0	0	1	0	1	0	0	1	0	1	1	1	0	1	0	1	1	1	0
20	0.20	0.30	0	1	1	1	0	1	1	1	0	0	1	1	1	0	0	0	0	1	1	1
21	-0.68	-0.90	0	0	0	0	1	1	0	1	0	0	1	1	1	1	1	1	0	1	0	0
22	1.13	0.91	0	1	1	0	0	0	0	1	0	0	1	0	0	0	0	0	1	1	0	1
23	-0.93	-1.89	1	0	1	0	0	0	0	1	0	0	1	0	0	0	0	1	1	0	1	0
24	-0.01	-0.68	1	1	1	0	0	1	1	0	1	0	1	1	0	0	1	0	0	1	0	1
25*	0.31	-0.48	1	0	1	0	0	0	0	1	1	0	0	0	0	1	0	1	0	1	0	0
26	-0.15	-0.63	0	0	1	1	0	1	0	1	1	0	0	0	1	1	1	0	0	0	1	1
27	-	-1.45	1	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	1	0
28	-	0.20	1	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	1	0	1	0
29	-	0.82	0	0	1	0	0	0	0	1	0	1	0	0	0	0	0	0	0	1	0	0
30	-	0.28	0	0	1	0	0	1	0	0	1	0	1	0	1	0	1	0	0	1	0	0
31	-	0.39	1	0	1	0	0	0	0	1	0	1	0	1	0	0	0	0	1	0	0	0
32	-	1.53	0	0	1	0	0	0	0	1	0	1	0	0	0	0	0	0	1	0	1	0
33	-	1.13	1	0	1	0	0	0	0	1	0	1	0	0	0	0	0	0	1	0	1	0
34	-	0.11	0	0	0	1	0	0	0	0	1	0	1	0	1	0	1	0	1	0	1	0
35	-	0.64	0	0	1	1	0	1	1	1	1	0	0	0	1	1	0	1	0	0	1	1
36	-	0.03	1	0	1	0	1	0	0	0	1	0	1	0	0	1	1	1	1	0	0	1
37*	1.05	-0.43	0	1	1	0	0	0	0	1	1	0	0	1	1	0	1	0	1	1	0	1
38	0.64	0.79	0	0	1	0	0	0	0	1	0	1	0	0	0	0	0	0	1	0	1	0
39	0.03	0.37	0	0	1	1	0	1	0	1	1	0	1	0	1	0	1	1	1	0	1	0
40	0.64	0.64	0	1	1	0	1	0	0	0	1	1	1	0	0	1	1	0	1	0	1	0
41	0.89	0.25	0	1	1	0	0	0	0	1	1	1	0	0	0	1	1	1	0	1	0	0
42	1.23	1.09	0	1	1	0	0	0	0	1	0	0	1	0	0	0	0	0	1	1	0	1
43	-0.23	-1.17	0	0	1	0	0	0	1	1	1	1	0	0	1	1	1	0	0	1	1	1
44	0.05	-1.11	0	1	1	0	0	0	0	1	0	1	0	0	0	0	0	0	1	0	1	0
45	0.53	-0.09	0	0	1	0	0	0	1	1	1	1	0	0	1	0	1	0	0	0	0	1
46	1.30	0.43	0	0	1	1	0	1	1	1	0	1	0	1	0	1	1	0	1	0	1	0
47	0.51	0.81	0	0	1	1	0	1	1	1	0	0	1	0	1	0	0	1	1	0	1	1
48	-0.15	0.03	0	0	1	0	0	0	0	1	1	1	0	1	1	0	0	1	1	0	1	0
49	0.91	1.22	1	1	1	0	0	0	0	1	0	1	1	0	0	0	1	1	1	1	0	1

*Item logit positions differ at $p < .05$ between Turkey/international.

Table 2
Estimates of Probability for Correct Performance on 20 Abilities across 17 Ability Levels

Ability Level	Content Attributes						Cognitive Process Attributes								Skill Attributes					
	C1	C2	C3	C4	C5	C6	P1	P2	P3	P4	P7	P8	P9	P10	S1	S2	S3	S6	S7	S9
-4.00	1.000	1.000	1.000	0.168	1.000	0.804	1.000	0.421	1.000	0.265	0.393	1.000	0.869	1.000	1.000	0.907	0.198	0.391	0.596	1.000
-3.50	1.000	1.000	1.000	0.180	1.000	0.795	1.000	0.443	1.000	0.293	0.446	1.000	0.988	1.000	1.000	0.952	0.233	0.446	0.663	1.000
-3.00	1.000	0.977	1.000	0.199	1.000	0.800	1.000	0.479	1.000	0.321	0.505	1.000	1.000	1.000	1.000	0.976	0.278	0.564	0.775	1.000
-2.50	1.000	0.941	1.000	0.221	1.000	0.816	1.000	0.521	1.000	0.357	0.571	1.000	1.000	1.000	1.000	0.909	0.336	0.676	0.896	1.000
-2.00	1.000	0.924	1.000	0.253	1.000	0.850	1.000	0.564	1.000	0.404	0.634	1.000	1.000	1.000	1.000	0.996	0.409	0.789	1.000	1.000
-1.50	1.000	0.925	1.000	0.313	1.000	0.919	1.000	0.617	1.000	0.470	0.698	1.000	1.000	1.000	1.000	0.983	0.500	0.882	1.000	1.000
-1.00	1.000	0.948	1.000	0.392	1.000	0.945	1.000	0.673	1.000	0.533	0.751	1.000	1.000	1.000	1.000	0.966	0.602	0.943	1.000	1.000
-0.50	1.000	0.980	0.578	0.514	1.000	1.000	1.000	0.740	1.000	0.640	0.800	1.000	1.000	1.000	1.000	0.937	0.705	0.972	1.000	1.000
0.00	1.000	1.000	0.743	0.654	1.000	1.000	1.000	0.808	1.000	0.734	0.850	1.000	1.000	1.000	1.000	0.921	0.797	0.975	1.000	1.000
0.50	1.000	1.000	0.876	0.785	1.000	1.000	1.000	0.873	1.000	0.822	0.901	1.000	1.000	1.000	1.000	0.924	0.873	0.972	1.000	1.000
1.00	1.000	1.000	0.958	0.885	1.000	1.000	1.000	0.928	1.000	0.889	0.944	1.000	1.000	1.000	1.000	0.940	0.930	0.973	1.000	0.991
1.50	1.000	1.000	0.990	0.945	1.000	1.000	1.000	0.963	1.000	0.938	0.972	1.000	1.000	1.000	1.000	0.962	0.965	0.981	1.000	0.992
2.00	1.000	1.000	0.999	0.975	1.000	1.000	1.000	0.983	1.000	0.969	0.987	1.000	1.000	1.000	1.000	0.979	0.983	0.989	1.000	0.996
2.50	1.000	1.000	1.000	0.989	1.000	1.000	1.000	0.993	1.000	0.986	0.945	1.000	1.000	1.000	1.000	0.990	0.993	0.995	1.000	0.998
3.00	1.000	1.000	1.000	0.995	1.000	1.000	1.000	0.997	1.000	0.994	0.998	1.000	1.000	1.000	1.000	0.996	0.997	0.998	1.000	0.999
3.50	1.000	1.000	1.000	0.998	1.000	1.000	1.000	0.999	1.000	0.997	0.999	1.000	1.000	1.000	1.000	0.998	0.999	0.999	1.000	1.000
4.00	1.000	1.000	1.000	0.999	1.000	1.000	1.000	0.999	1.000	0.999	1.000	1.000	1.000	1.000	1.000	0.999	0.999	0.100	1.000	1.000

Table 3
Absolute Differences for Item Characteristic Curve Recovery with the Least Squares Distance Method

Item	Ability (logits)																	MAD
	-4.0	-3.5	-3.0	-2.5	-2.0	-1.5	-1.0	-0.5	0.0	0.5	1.0	1.5	2.0	2.5	3.0	3.5	4.0	
1	0.024	0.006	0.013	0.030	0.067	0.143	0.282	0.479	0.683	0.834	0.922	0.965	0.985	0.993	0.997	0.999	1.000	0.043
2	5.403	0.001	0.003	0.007	0.016	0.037	0.082	0.173	0.329	0.534	0.729	0.863	0.936	0.972	0.988	0.995	0.998	0.024
3	0.007	0.016	0.037	0.082	0.173	0.329	0.534	0.729	0.863	0.936	0.972	0.988	0.995	0.998	0.999	1.000	1.000	0.173
4	0.002	0.003	0.008	0.019	0.043	0.094	0.196	0.363	0.572	0.758	0.880	0.945	0.976	0.990	0.996	0.998	0.999	0.058
5	3.780	8.848	0.002	0.005	0.011	0.026	0.059	0.128	0.255	0.445	0.652	0.815	0.912	0.960	0.983	0.993	0.997	0.033
6	0.002	0.005	0.012	0.027	0.062	0.133	0.265	0.458	0.664	0.822	0.916	0.962	0.984	0.993	0.997	0.999	0.999	0.119
7	0.006	0.013	0.030	0.067	0.143	0.282	0.479	0.683	0.834	0.922	0.965	0.985	0.993	0.997	0.999	1.000	1.000	0.066
8	0.003	0.007	0.016	0.036	0.079	0.168	0.321	0.526	0.722	0.859	0.934	0.971	0.987	0.995	0.998	0.999	1.000	0.102
9	2.048	4.790	0.001	0.003	0.006	0.014	0.033	0.073	0.156	0.303	0.504	0.704	0.848	0.929	0.968	0.986	0.994	0.112
10	3.355	7.855	0.002	0.004	0.010	0.023	0.053	0.115	0.233	0.416	0.625	0.796	0.901	0.955	0.980	0.992	0.996	0.055
11	0.003	0.007	0.016	0.037	0.083	0.175	0.332	0.538	0.732	0.865	0.937	0.972	0.988	0.995	0.998	0.999	1.000	0.049
12	1.728	4.046	9.471	0.002	0.005	0.012	0.028	0.063	0.135	0.268	0.462	0.668	0.825	0.917	0.963	0.984	0.993	0.119
13	8.869	0.002	0.005	0.011	0.026	0.058	0.126	0.252	0.441	0.649	0.812	0.910	0.960	0.982	0.992	0.997	0.999	0.025
14	0.005	0.011	0.026	0.058	0.126	0.252	0.441	0.649	0.812	0.910	0.960	0.982	0.992	0.997	0.999	0.999	1.000	0.165
15	0.027	0.061	0.131	0.262	0.453	0.660	0.820	0.914	0.962	0.983	0.993	0.997	0.999	0.999	1.000	1.000	1.000	0.232
16	7.093	0.002	0.004	0.009	0.021	0.048	0.105	0.215	0.391	0.601	0.779	0.892	0.951	0.978	0.991	0.996	0.998	0.023
17	0.007	0.017	0.039	0.086	0.180	0.340	0.547	0.739	0.869	0.939	0.973	0.988	0.995	0.998	0.999	1.000	1.000	0.139
18	0.001	0.003	0.006	0.014	0.032	0.071	0.152	0.296	0.496	0.697	0.844	0.927	0.967	0.986	0.994	0.997	0.999	0.054
19	3.472	8.127	0.002	0.004	0.010	0.024	0.054	0.118	0.239	0.424	0.633	0.802	0.904	0.957	0.981	0.992	0.997	0.060
20	6.627	0.002	0.004	0.008	0.020	0.045	0.099	0.204	0.375	0.584	0.767	0.885	0.948	0.977	0.990	0.996	0.998	0.019
21	0.005	0.012	0.027	0.062	0.133	0.265	0.458	0.664	0.822	0.916	0.962	0.984	0.993	0.997	0.999	0.999	1.000	0.078
22	2.347	5.496	0.001	0.003	0.007	0.016	0.037	0.083	0.175	0.332	0.538	0.732	0.865	0.937	0.972	0.988	0.995	0.068
23	0.027	0.061	0.131	0.262	0.453	0.660	0.820	0.914	0.962	0.983	0.993	0.997	0.999	0.999	1.000	1.000	1.000	0.233
24	0.004	0.008	0.019	0.043	0.096	0.199	0.367	0.576	0.761	0.882	0.946	0.976	0.990	0.996	0.998	0.999	1.000	0.041
25	0.003	0.006	0.014	0.031	0.070	0.150	0.292	0.492	0.694	0.841	0.926	0.967	0.986	0.994	0.997	0.999	1.000	0.043
26	0.003	0.008	0.017	0.040	0.089	0.185	0.348	0.555	0.745	0.873	0.941	0.974	0.989	0.995	0.998	0.999	1.000	0.103
27	0.013	0.030	0.067	0.143	0.282	0.479	0.683	0.834	0.922	0.965	0.985	0.993	0.997	0.999	1.000	1.000	1.000	0.100
28	7.855	0.002	0.004	0.010	0.023	0.053	0.115	0.233	0.416	0.625	0.796	0.901	0.955	0.980	0.992	0.996	0.998	0.044
29	2.736	6.405	0.002	0.004	0.008	0.019	0.043	0.096	0.199	0.367	0.576	0.761	0.882	0.946	0.976	0.990	0.996	0.075
30	6.856	0.002	0.004	0.009	0.020	0.046	0.102	0.210	0.383	0.593	0.773	0.889	0.949	0.978	0.990	0.996	0.998	0.076
31	5.686	0.001	0.003	0.007	0.017	0.039	0.086	0.180	0.340	0.547	0.739	0.869	0.939	0.973	0.988	0.995	0.998	0.022
32	8.173	1.914	4.481	0.001	0.003	0.006	0.013	0.031	0.069	0.148	0.289	0.487	0.690	0.839	0.924	0.966	0.985	0.145
33	1.614	3.780	8.848	0.002	0.005	0.011	0.026	0.059	0.128	0.255	0.445	0.652	0.815	0.912	0.960	0.983	0.993	0.098
34	9.154	0.002	0.005	0.012	0.027	0.061	0.131	0.262	0.453	0.660	0.820	0.914	0.962	0.983	0.993	0.997	0.999	0.022
35	3.716	8.699	0.002	0.005	0.011	0.026	0.058	0.126	0.252	0.441	0.649	0.812	0.910	0.960	0.982	0.992	0.997	0.034
36	0.001	0.003	0.006	0.013	0.031	0.069	0.148	0.289	0.487	0.690	0.839	0.924	0.966	0.985	0.994	0.997	0.999	0.019
37	0.002	0.005	0.012	0.029	0.065	0.139	0.275	0.470	0.675	0.830	0.919	0.964	0.984	0.993	0.997	0.999	1.000	0.063
38	2.879	6.740	0.002	0.004	0.009	0.020	0.045	0.100	0.207	0.379	0.588	0.770	0.887	0.948	0.977	0.990	0.996	0.058
39	5.883	0.001	0.003	0.008	0.017	0.040	0.089	0.185	0.348	0.555	0.745	0.873	0.941	0.974	0.989	0.995	0.998	0.034
40	3.716	8.699	0.002	0.005	0.011	0.026	0.058	0.126	0.252	0.441	0.649	0.812	0.910	0.960	0.982	0.992	0.997	0.036
41	7.215	0.002	0.004	0.009	0.021	0.048	0.107	0.218	0.395	0.605	0.782	0.894	0.952	0.979	0.991	0.996	0.998	0.018
42	1.728	4.046	9.471	0.002	0.005	0.012	0.028	0.063	0.135	0.268	0.462	0.668	0.825	0.917	0.963	0.984	0.993	0.089
43	0.008	0.019	0.043	0.094	0.196	0.363	0.572	0.758	0.880	0.945	0.976	0.990	0.996	0.998	0.999	1.000	1.000	0.161
44	0.007	0.017	0.039	0.086	0.180	0.340	0.547	0.739	0.869	0.939	0.973	0.988	0.995	0.998	0.999	1.000	1.000	0.166
45	0.001	0.003	0.007	0.016	0.037	0.083	0.175	0.332	0.538	0.732	0.865	0.937	0.972	0.988	0.995	0.998	0.999	0.028
46	5.312	0.001	0.003	0.007	0.016	0.036	0.081	0.170	0.325	0.530	0.725	0.861	0.935	0.971	0.988	0.995	0.998	0.030
47	2.783	6.515	0.002	0.004	0.008	0.019	0.044	0.097	0.201	0.371	0.580	0.764	0.883	0.947	0.977	0.990	0.996	0.017
48	0.001	0.003	0.006	0.013	0.031	0.069	0.148	0.289	0.487	0.690	0.839	0.924	0.966	0.985	0.994	0.997	0.999	0.053
49	1.385	3.243	7.592	0.002	0.004	0.010	0.022	0.051	0.111	0.227	0.408	0.617	0.790	0.898	0.954	0.980	0.991	0.065

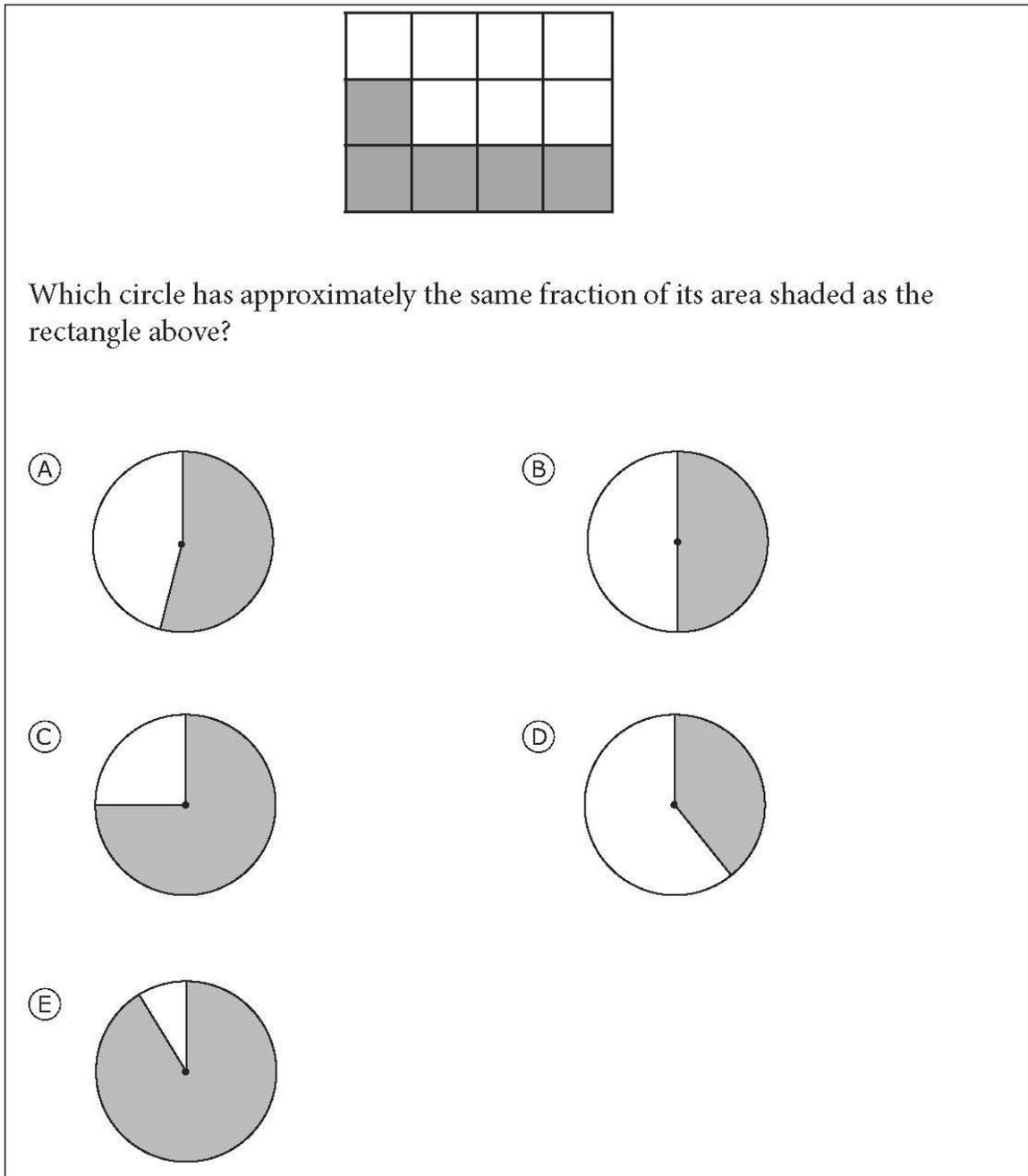


Figure 2. Sample item (TIMSS-2007 Technical Report, 2008)

A bowl contains 36 colored beads all of the same size: some blue, some green, some red, and the rest yellow. A bead is drawn from the bowl without looking. The probability that it is blue is $\frac{4}{9}$. How many blue beads are in the bowl?

(A) 4
 (B) 8
 (C) 16
 (D) 18
 (E) 20

Figure 3. Sample item (TIMSS-2007 Technical Report, 2008)

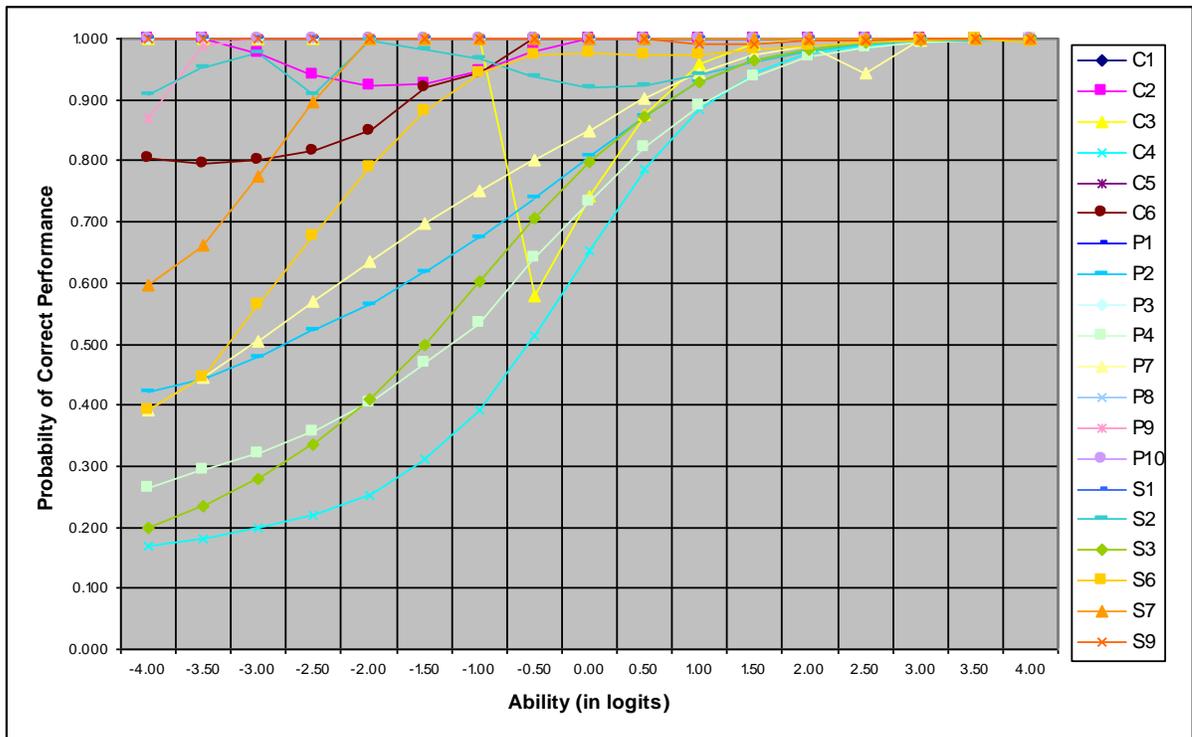


Figure 4. Attribute Probability Curves

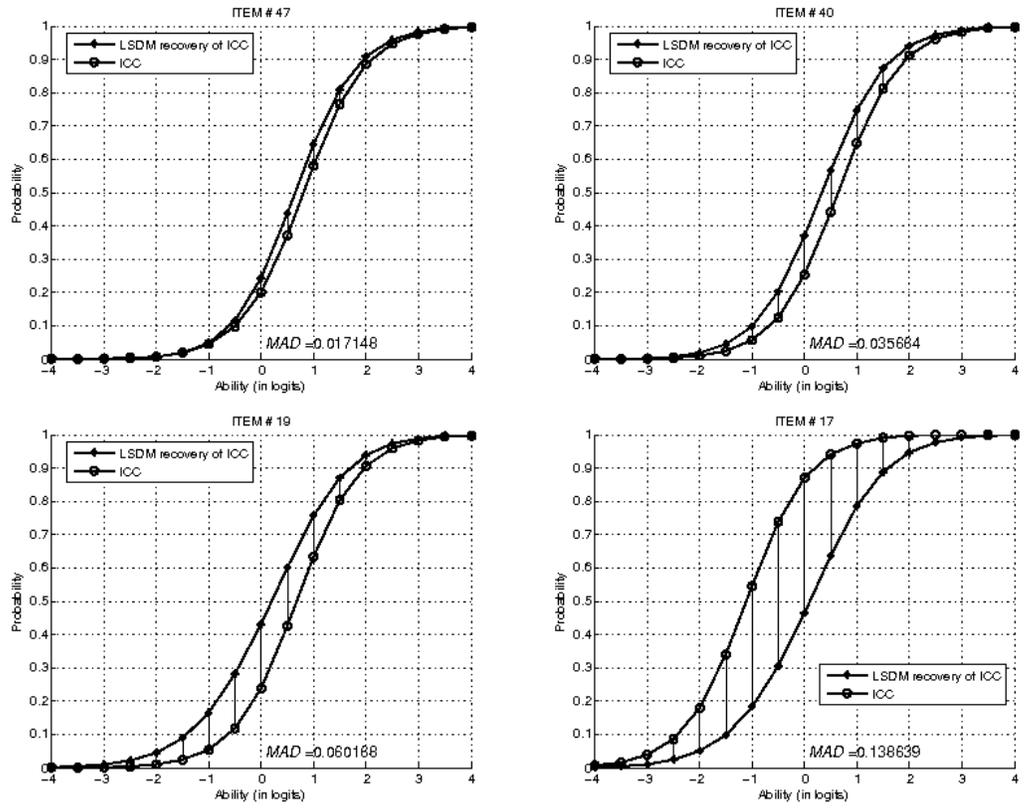


Figure 5. Item Characteristic Curve Recovery with the LSDM for Four Sample

Items

Appendix

The Content, Cognitive Process and Skill Attributes for the TIMSS-R (1999)

Content attributes

- C1 Basic concepts and operations in whole numbers and integers
- C2 Basic concepts and operations in fractions and decimals
- C3 Basic concepts and operations in elementary algebra
- C4 Basic concepts and operations in two-dimensional geometry
- C5 Data, probability, and basic statistics
- C6 Measuring or estimating: length, time, angle, temperature, etc.

Cognitive Process attributes

- P1 Translate/formulate equations and expressions to solve a problem
- P2 Computational applications of knowledge in arithmetic and geometry
- P3 Judgmental applications of knowledge in arithmetic and geometry
- P4 Applying rules in algebra
- P5 Logical reasoning-includes case reasoning, deductive thinking skills, if-then, necessary and sufficient, generalization skills
- P6 Problem search; analytic thinking, problem restructuring; inductive thinking
- P7 Generating, visualizing, and reading figures and graphs
- P8 Applying and evaluating mathematical correctness
- P9 Management of data and procedures
- P10 Quantitative and logical reading

Skill (item type) attributes

- S1 Unit conversion
- S2 Apply number properties and relationships; number sense/number line
- S3 Using figures, tables, charts, and graphs
- S4 Approximation/estimation
- S5 Evaluate/verify/check options
- S6 Patterns and relationships (inductive thinking skills)
- S7 Using proportional reasoning
- S8 Solving novel or unfamiliar problems
- S9 Comparison of two/or more entities
- S10 Open-ended items, in which an answer is not given
- S11 Understanding verbally posed questions