**Psychometric Analyses of the 2006 MCAS High School Introductory Physics Test[1,2]**

**Nina Deng and Ronald K. Hambleton**
**University of Massachusetts Amherst**

**February 4, 2008**

**1. Goal of the Psychometric Analyses**

The primary goal of our work has been to provide readers with a number of worthwhile psychometric analyses of the 2006 MCAS high school Introductory Physics Test. These analyses provide more details on the Introductory Physics Test than it was possible to provide in the summary report prepared by Hambleton, Zhao, Smith, Lam, and Deng (2008). These analyses include (1) an item analysis, (2) descriptive statistics on the test scores including break-outs for several subgroups of students, (3) classical reliability analyses for the test scores organized by item format, and for the total test, (4) two investigations of test dimensionality, (5) item response theory (IRT) item calibrations obtained from fitting the three-parameter logistic model to binary-scored items and the graded response model to polytomously-scored items, (6) various item and test level model fit findings, (7) test information and conditional standard errors, and (8) the identification of differentially functioning test items.

**2. Description of the Introductory Physics Test**

The MCAS 2006 Grade 9/10 Introductory Physics Test consists of 45 items assessing six standards (sometimes called "curriculum strands"): Motion and Forces, Conservation of Energy and Momentum, Heat and Heat Transfer, Waves, Electromagnetism, and Electromagnetic Radiation, based on learning standards in the Physics content strand of the Massachusetts Science and Technology/Engineering Curriculum Framework (2006). The test was administered in a 2-day session in May of 2006, the first session consisted of the first 26 items on the test; and the second session consisted of the remaining 19 items. Each session included multiple-choice and open-response questions. More information about the curriculum and the test items can be found at **www.doe.mass.edu**.

Table 2.1 presents the number of items, by item type, and the total number of items and score points for the MCAS 2006 Grade 9/10 Introductory Physics Test. There are 40

multiple choice items (each with four choices) and five polytomously-scored performance items (or sometimes called "constructed response items").  Multiple choice items were scored dichotomously; a score of 1 for a correct answer, 0 otherwise. Performance items were scored polytomously, with possible scores ranging from 0 to 4.

**Table 2.1  Number of Items by Item Type on the Introductory Physics Test**

| Item Type | Points | Number of Items |
|---|---|---|
| Multiple Choice | 1 or 0 | 40 |
| Performance | 0 to 4 | 5 |
| Total | 60 | 45 |

## 3. Item Analyses

In total, 16,619 students were administered the Physics Test.  However, an exclusion criterion was implemented so as to reduce the distortion of findings due to the use of student responses that would introduce systematic errors into the data analyses.  Students who had a total test score of 0 were excluded.  Clearly, these students had not taken the test seriously, or perhaps were not even present for the test administration.  After applying the exclusion rule, there were 15,762 students left in the dataset.  Therefore, about 5% of the examinee data were excluded.  These students served no useful purpose for our psychometric analyses of the items and the test and so they were deleted.
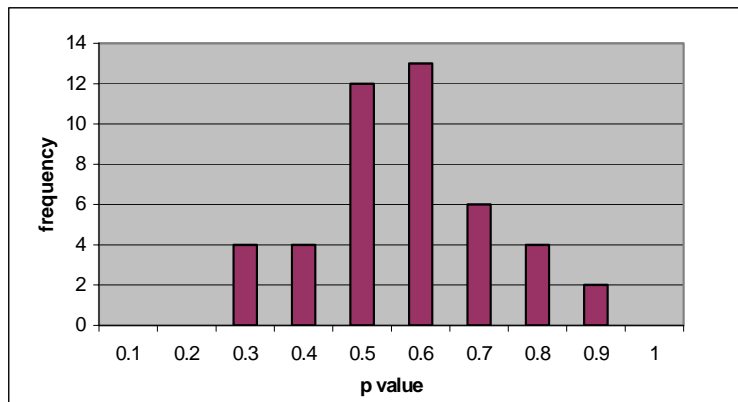
*Item Difficulty*

Difficulty indices (correct proportion for binary-scored data or item means for polytomously-scored data) were measured by averaging the proportion of points for correct answers received by all the students who answered the test items. As a result, for students who were absent for the test or who showed up with none of the items answered, their

responses were not included in these analyses. Multiple choice questions (MC) were scored dichotomously (0 and 1). For these items, difficulty indices were simply the proportion of students who got the correct answers. Performance items were scored from 0 to 4. For those items, difficulty indices were obtained by dividing the average proportions by 4, so that all the difficulty indices were on the same scale and within the range from 0 to 1, regardless of item type. Indices of larger values indicate easier items. For example, an index of 1 means everyone got the correct answer, while an index of 0 means no one received the point. Summary statistics and histogram of item difficulty index distribution are provided in Table 3.1 and Figure 3.1.

**Table 3.1 Summary of Classical Item Difficulty Indices**

| Number of Item | Mean | SD | Range | Minimum | Maximum |
|:--:|:--:|:--:|:--:|:--:|:--:|
| 45 | 0.53 | 0.15 | 0.62 | 0.25 | 0.87 |

**Figure 3.1 Histogram Showing the Distribution of Classical Item Difficulty Indices**



The figure showing that difficulty indices have a wide range:  From about the level of chance (.25) to about 0.9, with the mean around 0.5.
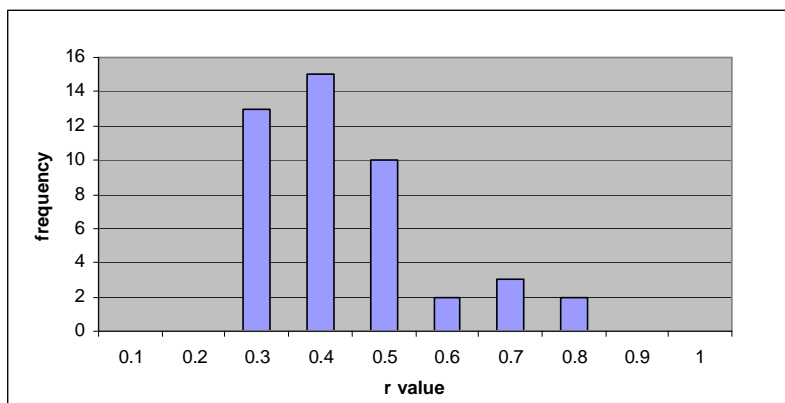
*Item-Test Correlations*

Item-test correlations are called item discrimination indices.  They indicate the degree to which test items distinguish between the performance of higher proficient and lower

4

proficient students. The discrimination index was calculated using Pearson correlations and the range is within -1.0 to 1.0. The typical range of discrimination indices for operational multiple-choice items is from 0.2 to 0.6. Summary statistics and histogram of discrimination indices distribution are provided in Table 3.2 and Figure 3.2.

**Table 3.2 Summary of the Classical Item Discrimination Indices**

| Number of Item | Mean | S.D. | Range | Minimum | Maximum |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 45 | 0.40 | 0.14 | 0.56 | 0.23 | 0.79 |

**Figure 3.2 Histogram Showing the Distribution of**

**Classical Item Discrimination Indices**



A comparison of difficulty and discrimination indices between multiple choice items and performance items are displayed in Table 3.3.

**Table 3.3  Average Item Difficulty and Discrimination Indices across Item Types**

| | Item Type | | |
|---|:---:|:---:|:---:|
| **Average Statistics** | **All** | **MCQ** | **Performance Items (PI)** |
| Difficulty (p) | 0.527 | 0.539 | 0.427 |
| Discrimination  (r) | 0.399 | 0.359 | 0.718 |
| Number of Items | 45 | 40 | 5 |

Because multiple choice items can be answered correctly by guessing, they generally have higher difficulty indices (i.e., are easier items) than many performance items. Besides, for the similar reason, their discrimination indices are usually lower than the performance items since the multiple choice questions could have chance scores and therefore decrease the possible score range and introduce error. In addition, candidates spend much more time in completing the performance items. Therefore the information provided by performance items is more reliable and more information is coming in. Hence their r values are higher. The indices of the 45 items are provided in Table 3.4.

**Table 3.4   Classical Item Statistics**

| Item | Item Type | Difficulty | Discrimination |
|---|---|---|---|
| 1 | MC | 0.847 | 0.253 |
| 2 | MC | 0.434 | 0.371 |
| 3 | MC | 0.532 | 0.43 |
| 4 | MC | 0.719 | 0.397 |
| 5 | MC | 0.733 | 0.418 |
| 6 | MC | 0.776 | 0.328 |
| 7 | MC | 0.389 | 0.297 |
| 8 | MC | 0.493 | 0.341 |
| 9 | MC | 0.758 | 0.4 |
| 10 | MC | 0.501 | 0.337 |
| 11 | PI | 0.660 | 0.681 |
| 12 | MC | 0.545 | 0.355 |
| 13 | MC | 0.598 | 0.231 |
| 14 | MC | 0.470 | 0.375 |
| 15 | MC | 0.700 | 0.444 |
| 16 | MC | 0.571 | 0.278 |
| 17 | MC | 0.501 | 0.282 |
| 18 | MC | 0.483 | 0.293 |
| 19 | MC | 0.582 | 0.337 |
| 20 | MC | 0.459 | 0.318 |
| 21 | MC | 0.520 | 0.274 |
| 22 | MC | 0.367 | 0.304 |
| 23 | MC | 0.593 | 0.368 |
| 24 | MC | 0.543 | 0.387 |
| 25 | PI | 0.435 | 0.787 |
| 26 | PI | 0.428 | 0.781 |
| 27 | MC | 0.874 | 0.415 |
| 28 | MC | 0.477 | 0.322 |
| 29 | MC | 0.646 | 0.467 |
| 30 | MC | 0.637 | 0.535 |
| 31 | MC | 0.611 | 0.478 |
| 32 | PI | 0.345 | 0.663 |
| 33 | MC | 0.378 | 0.256 |
| 34 | MC | 0.555 | 0.431 |
| 35 | MC | 0.485 | 0.362 |
| 36 | MC | 0.251 | 0.277 |
| 37 | MC | 0.278 | 0.235 |
| 38 | MC | 0.486 | 0.455 |
| 39 | PI | 0.268 | 0.677 |
| 40 | MC | 0.426 | 0.284 |
| 41 | MC | 0.428 | 0.357 |
| 42 | MC | 0.527 | 0.449 |
| 43 | MC | 0.628 | 0.505 |
| 44 | MC | 0.259 | 0.261 |
| 45 | MC | 0.511 | 0.465 |

*Distracter Analysis*

The proportions of students choosing each option of the multiple choice items were calculated by **Test Analysis Program** (TAP), a software program coming from Ohio University. The detailed results for each item are provided in Appendix A.

In Appendix A, the frequencies and percentages of students who chose each of the four choices are displayed for each item. Besides, the statistics and the differences were calculated for high and low groups (respectively top and low 27% of total scores). The correct answer should have a distinct positive difference (more of the high group than of the low group), and negative differences should be displayed across all the other incorrect answers. The results show that all the correct answers of 40 multiple choice items have positive differences between the high and low groups, and they enjoy satisfactory large proportions in the high group. Besides, most of the incorrect answers show negative differences, except for two options, option 3 of item 33 and option 3 of item 44, with difference indices of 3% and 10% respectively. The indices indicate that the two incorrect options attracted more students in high group than in low group and may need some improvement, however, the differences are not significant. Overall, most incorrect answers, called distracters, functioned well in distracting the students from the low group with an average percentage of around 20%.

In sum, the item statistics are highly supportive of an excellent test. Of special importance for tests like the MCAS tests, are items with high discriminating powers, and a range of p values to support consistent and accurate performance classifications at three widely spaced performance standards along the proficiency scale.

## 4. Basic Test Score Descriptive Statistics and Reliability Analyses

*Basic Test Score Descriptive Statistics*

There were 16,619 students in total enrolled in the 2006 Introductory Physics Test. The distributions of attendance, gender and ethnicity are displayed in Table 4.1.

**Table 4.1 Demographic Statistics**

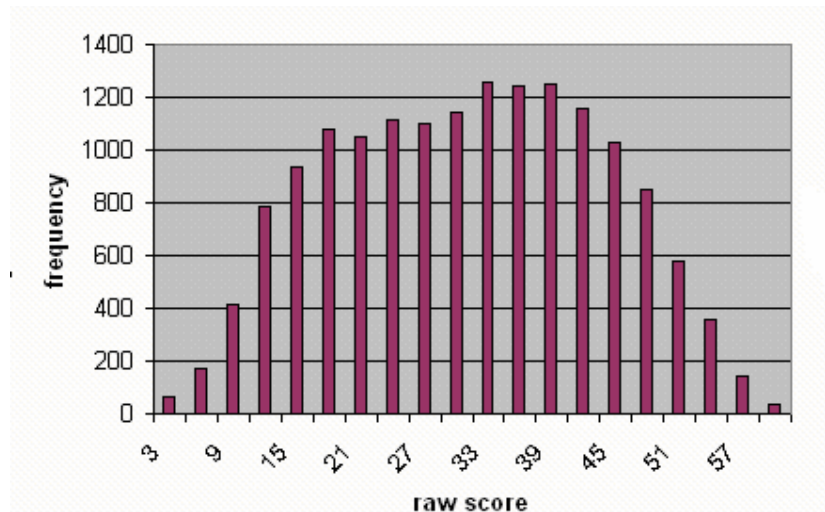|  | Subgroup | Number | Percentage (%) |
|---|---|---|---|
| **Attendance** | Valid | 15,762 | 94.8 |
|  | Absent | 622 | 3.7 |
|  | Scored "0" | 235 | 1.4 |
| **Gender** | Male | 8,219 | 49.5 |
|  | Female | 7,925 | 47.7 |
|  | Unspecified | 475 | 2.9 |
| **Ethnicity** | White | 10,939 | 65.8 |
|  | Black | 2,141 | 12.9 |
|  | Hispanic | 1,997 | 12.0 |
|  | Asian | 1,007 | 6.1 |
|  | Native | 50 | 0.3 |
|  | Unspecified | 485 | 2.9 |
| **Enrolled in Total** | | 16,619 | 100 |

There are 15,762 valid total scores out of the enrolled students. The summary statistics of the valid total scores are shown in Table 4.2, and its histogram of distribution is shown in Figure 4.1. The distribution might be described as platykurtic.

**Table 4.2 Descriptive Statistics of Total Scores**

| Mean | SD | Median | Mode | Minimum | Maximum | Number of Students * |
|-------|-------|--------|------|---------|---------|----------------------|
| 30.12 | 12.55 | 31 | 35 | 1 | 60 | 15762 |

* For the purpose of test analysis only, the number of students excludes the number of absent students and students who showed up with zero scores.

**Figure 4.1   Histogram Showing the Distribution of Total Scores**



*Reliability*

Internal consistency reliability was estimated using *coefficient alpha*, which measures the degree to which all of the items measure a common characteristic of the person and it depends on the consistency of the individual's performance from item to item. The reliability statistics of all items, the multiple-choice questions and the performance items are shown in Table 4.3.

**Table 4.3  Test Score Reliability Statistics**

| Statistic | All Items | MCQ | Performance Items |
|-----------|-----------|-------|-------------------|
| Coefficient Alpha | 0.901 | 0.874 | 0.83 |

In addition, the correlation coefficients between MCQ scores, performance item scores and total scores are displayed in Table 4.5.

**Table 4.5 Correlations Among MCQ Scores, Performance Scores, and Total Scores**

| Score | MCQ Score | Performance Item Score |
|---|---|---|
| Performance Items Score | 0.801 | |
| Total Score | 0.967 | 0.927 |

## 5. Test Dimensionality

Before any IRT-related analyses, one important assumption to be confirmed is unidimensionality. The assumption of unidimensionality means that all the items are measuring a single dominant trait; in this case it refers to physics proficiency. Good model fit requires a reasonable good approximation of the unidimensionality assumption. Even the data in Table 4.5 provides an initial estimate of unidimensionality since the MCQ scores and performance item scores are very highly correlated (0.80) without any adjustments for the unreliability of each score. If any multidimensionality in the data were to be present, it might be expected to show up in items assessed by different formats and measuring different learning standards, and then the correlation would be expected to be considerably lower than .80.

Table 5.1 shows the eigenvalues of the 45x45 correlation matrix. Figure 5.1 shows the scree plot. The largest engenvalue accounted for about 30% of the total variance. Also, the first eigenvalue is over 7 times larger than the second eigenvalue. Based on the conventional standards, that is, the first dominant factor accounts for at least 20% of the total variance and is four or five times larger than the second factor, it demonstrates that the test is strongly unidimensional.

**Table 5.1  Eigenvalues and Variances Explained**

| Component | Eigenvalue | Percentage of Variance Explained (%) | Cumulative Percentage of Variance Explained (%) |
|---|---|---|---|
| 1 | 13.5 | 30.05 | 30 |
| 2 | 1.65 | 3.67 | 34 |
| 3 | 1.41 | 3.14 | 37 |
| 4 | 1.32 | 2.94 | 40 |
| 5 | 1.10 | 2.44 | 42 |
| 6 | 1.07 | 2.38 | 45 |
| 7 | 1.01 | 2.24 | 47 |
| 8 | 0.99 | 2.21 | 49 |
| 9 | 0.97 | 2.16 | 51 |
| 10 | 0.93 | 2.06 | 53 |
| 11 | 0.92 | 2.05 | 55 |
| 12 | 0.90 | 2.00 | 57 |
| 13 | 0.87 | 1.94 | 59 |
| 14 | 0.83 | 1.84 | 61 |
| 15 | 0.81 | 1.80 | 63 |
| 16 | 0.80 | 1.78 | 65 |
| 17 | 0.79 | 1.76 | 66 |
| 18 | 0.77 | 1.72 | 68 |
| 19 | 0.76 | 1.69 | 70 |
| 20 | 0.73 | 1.62 | 71 |
| 21 | 0.72 | 1.59 | 73 |
| 22 | 0.70 | 1.56 | 75 |
| 23 | 0.69 | 1.54 | 76 |
| 24 | 0.69 | 1.53 | 78 |
| 25 | 0.67 | 1.49 | 79 |
| 26 | 0.63 | 1.41 | 81 |
| 27 | 0.61 | 1.36 | 82 |
| 28 | 0.61 | 1.35 | 83 |
| 29 | 0.59 | 1.31 | 85 |
| 30 | 0.58 | 1.28 | 86 |
| 31 | 0.56 | 1.24 | 87 |
| 32 | 0.54 | 1.20 | 88 |
| 33 | 0.52 | 1.15 | 89 |
| 34 | 0.51 | 1.13 | 91 |
| 35 | 0.50 | 1.11 | 92 |
| 36 | 0.47 | 1.04 | 93 |
| 37 | 0.45 | 1.01 | 94 |
| 38 | 0.44 | 0.99 | 95 |
| 39 | 0.43 | 0.95 | 96 |
| 40 | 0.38 | 0.85 | 97 |
| 41 | 0.37 | 0.83 | 97 |
| 42 | 0.35 | 0.77 | 98 |
| 43 | 0.30 | 0.68 | 99 |
| 44 | 0.27 | 0.59 | 99 |
| 45 | 0.25 | 0.55 | 100 |

**Figure 5.1  Plot of the 45 Eigenvalues**

**Eigenvalue Plot**



The students' responses to the 45 items were further analyzed using Confirmatory

Factor Analysis (CFA) with Structural Equation Modeling (SEM) by the software package

LISREL. Table 5.2 shows the factor loadings of a 1-factor model on the 45 items (variables).

**Table 5.2 Estimated Factor Loadings in a One-Factor Solution**

| Item | Factor Loading |
|------|----------------|
| 1 | .40 |
| 2 | .56 |
| 3 | .66 |
| 4 | .64 |
| 5 | .64 |
| 6 | .51 |
| 7 | .43 |
| 8 | .52 |
| 9 | .62 |
| 10 | .57 |
| 11 | .74 |
| 12 | .54 |
| 13 | .35 |
| 14 | .55 |
| 15 | .65 |
| 16 | .43 |
| 17 | .43 |
| 18 | .49 |
| 19 | .53 |
| 20 | .49 |
| 21 | .43 |
| 22 | .47 |
| 23 | .56 |
| 24 | .61 |
| 25 | .83 |
| 26 | .83 |
| 27 | .71 |
| 28 | .46 |
| 29 | .67 |
| 30 | .80 |
| 31 | .71 |
| 32 | .73 |
| 33 | .36 |
| 34 | .70 |
| 35 | .61 |
| 36 | .47 |
| 37 | .37 |
| 38 | .66 |
| 39 | .75 |
| 40 | .44 |
| 41 | .54 |
| 42 | .70 |
| 43 | .81 |
| 44 | .44 |
| 45 | .69 |

All of the loadings are high, and so these results combined with the other analyses we carried out are strongly supportive of a strong first factor, a perquisite for a good fitting unidimensional IRT model.

## 6. Item Calibrations and Model Fit

Forty dichotomous items and five polytomous items were calibrated using Parscale, with the 3P logistic model fitted to dichotomous items, and the Graded Response Model (GRM) fitted to the polytomous items. The estimated discrimination (slope), difficulty (location), and guessing parameter estimates of the 45 items and their summary statistics are shown in Table 6.1 and Table 6.2. Figure 6.1 shows the score category curves for all the items.
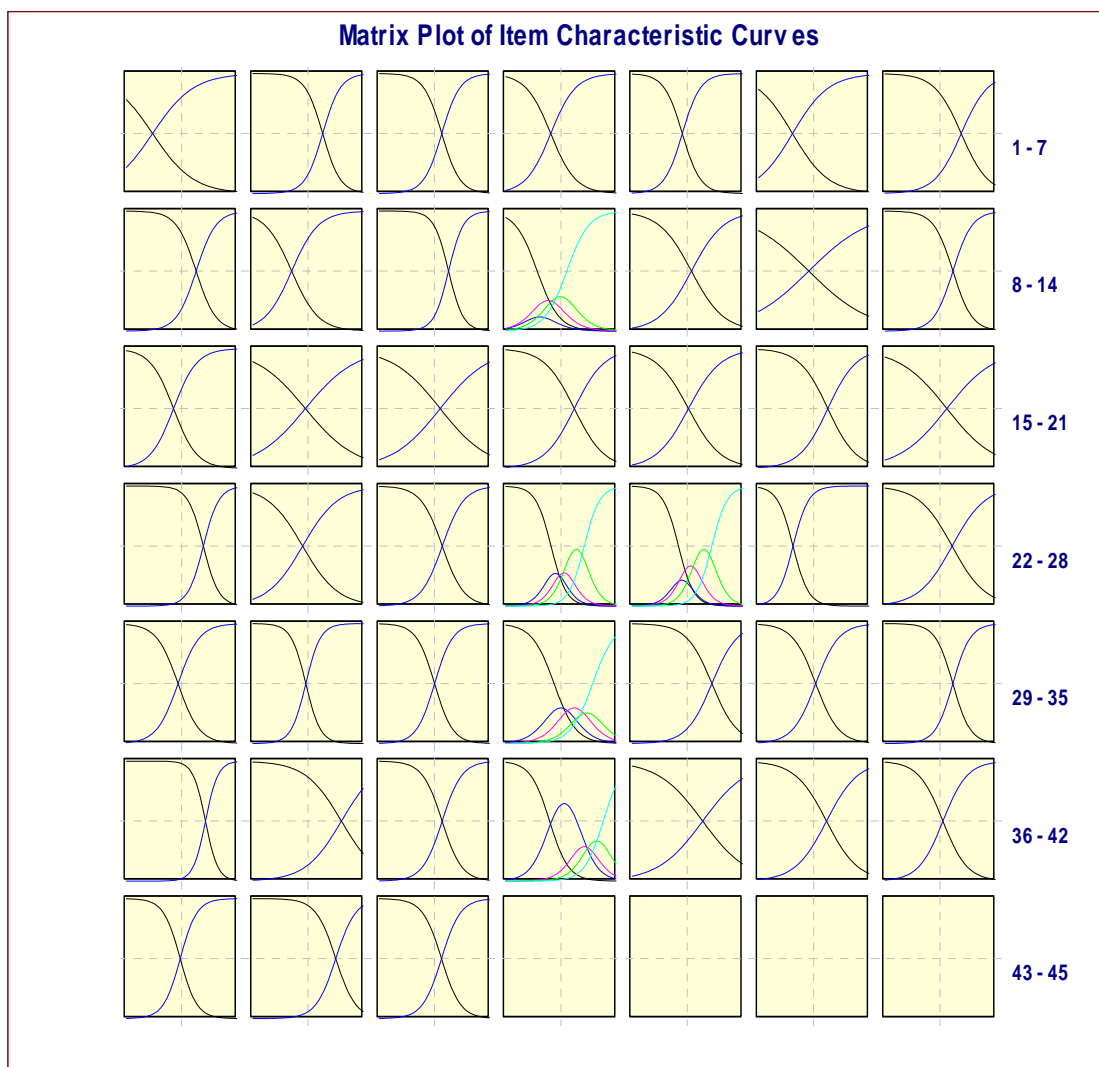
**Table 6.1 Item Parameter Estimates**

| Item | Slope(a) | SE of a | Location(b) | SE of b | Guessing(c) | SE of c |
|------|----------|---------|-------------|---------|-------------|---------|
| 1 | 0.53 | 0.04 | -1.58 | 0.29 | 0.32 | 0.09 |
| 2 | 1.23 | 0.08 | 0.80 | 0.03 | 0.22 | 0.01 |
| 3 | 1.20 | 0.07 | 0.40 | 0.04 | 0.23 | 0.02 |
| 4 | 0.87 | 0.05 | -0.55 | 0.08 | 0.20 | 0.03 |
| 5 | 1.19 | 0.07 | -0.27 | 0.06 | 0.34 | 0.02 |
| 6 | 0.60 | 0.04 | -1.15 | 0.15 | 0.17 | 0.06 |
| 7 | 0.81 | 0.07 | 1.15 | 0.05 | 0.21 | 0.02 |
| 8 | 1.10 | 0.08 | 0.79 | 0.04 | 0.30 | 0.02 |
| 9 | 0.81 | 0.04 | -0.89 | 0.09 | 0.15 | 0.04 |
| 10 | 1.47 | 0.10 | 0.77 | 0.03 | 0.31 | 0.01 |
| 11 | 0.98 | 0.02 | -0.59 | 0.02 | 0.00 | 0.00 |
| 12 | 0.67 | 0.05 | 0.22 | 0.08 | 0.17 | 0.03 |
| 13 | 0.35 | 0.03 | -0.27 | 0.22 | 0.13 | 0.05 |
| 14 | 1.13 | 0.07 | 0.70 | 0.04 | 0.23 | 0.01 |
| 15 | 0.93 | 0.05 | -0.44 | 0.07 | 0.20 | 0.03 |
| 16 | 0.43 | 0.03 | -0.14 | 0.13 | 0.09 | 0.04 |
| 17 | 0.45 | 0.03 | 0.31 | 0.11 | 0.09 | 0.03 |
| 18 | 0.72 | 0.06 | 0.73 | 0.07 | 0.23 | 0.02 |
| 19 | 0.67 | 0.05 | 0.06 | 0.09 | 0.19 | 0.03 |
| 20 | 0.76 | 0.06 | 0.77 | 0.06 | 0.21 | 0.02 |
| 21 | 0.45 | 0.04 | 0.37 | 0.14 | 0.13 | 0.04 |
| 22 | 1.40 | 0.11 | 1.18 | 0.03 | 0.24 | 0.01 |
| 23 | 0.61 | 0.03 | -0.29 | 0.07 | 0.07 | 0.02 |
| 24 | 1.00 | 0.06 | 0.43 | 0.05 | 0.26 | 0.02 |
| 25 | 1.41 | 0.02 | 0.27 | 0.01 | 0.00 | 0.00 |
| 26 | 1.41 | 0.02 | 0.30 | 0.01 | 0.00 | 0.00 |
| 27 | 1.39 | 0.08 | -1.12 | 0.07 | 0.31 | 0.04 |
| 28 | 0.65 | 0.05 | 0.65 | 0.08 | 0.19 | 0.03 |
| 29 | 1.01 | 0.05 | -0.20 | 0.05 | 0.19 | 0.02 |
| 30 | 1.47 | 0.07 | -0.12 | 0.03 | 0.18 | 0.02 |
| 31 | 1.19 | 0.06 | 0.01 | 0.04 | 0.20 | 0.02 |
| 32 | 0.98 | 0.02 | 0.70 | 0.02 | 0.00 | 0.00 |
| 33 | 0.85 | 0.08 | 1.34 | 0.06 | 0.23 | 0.02 |
| 34 | 0.94 | 0.05 | 0.11 | 0.05 | 0.15 | 0.02 |
| 35 | 1.24 | 0.08 | 0.69 | 0.04 | 0.28 | 0.01 |
| 36 | 1.75 | 0.12 | 1.30 | 0.03 | 0.15 | 0.01 |
| 37 | 0.60 | 0.07 | 1.81 | 0.09 | 0.12 | 0.02 |
| 38 | 1.10 | 0.06 | 0.42 | 0.03 | 0.15 | 0.01 |
| 39 | 1.14 | 0.02 | 1.06 | 0.01 | 0.00 | 0.00 |
| 40 | 0.48 | 0.04 | 0.84 | 0.10 | 0.11 | 0.03 |
| 41 | 0.69 | 0.05 | 0.70 | 0.06 | 0.13 | 0.02 |
| 42 | 0.89 | 0.05 | 0.15 | 0.05 | 0.12 | 0.02 |
| 43 | 1.29 | 0.06 | -0.08 | 0.04 | 0.19 | 0.02 |
| 44 | 1.09 | 0.09 | 1.51 | 0.05 | 0.15 | 0.01 |
| 45 | 1.14 | 0.06 | 0.38 | 0.04 | 0.19 | 0.02 |

**Table 6.2    Summary Statistics of Item Parameter Estimates**

| Parameter | Mean | SD | Number |
|-----------|------|-----|--------|
| Slope | 0.958 | 0.336 | 45 |
| Location | 0.294 | 0.727 | 45 |
| Guessing | 0.193 | 0.067 | 40 |

**Figure 6.1 Item Score Category Function**



Matrix Plot of Item Characteristic Curves

*Model Fit*

Item response theory possesses many advantages over classical test theory in analyzing the measurement data of latent traits. However, the advantages will be greatly undermined if the model used to analyze does not fit the observed data. Therefore the assessment of model fit should always be carried out as an integrated part of IRT analyses.

Most IRT calibration programs offer the fit statistics at the item level. Table 6.3 provides the Chi-square item fit statistics and probabilities from Parscale, though there is not much confidence in these statistics because they are very dependent on sample size. With very big samples as we used in this study, it will appear that none or only a few items will actually fit the data. This is a common finding and it is little value.

**Table 6.3  Item Fit Statistics**

| Item | Chi-square | D.F. | Probability |
|------|-----------|------|-------------|
| 1 | 49.77 | 28 | 0.007 |
| 2 | 33.59 | 30 | 0.297 |
| 3 | 43.78 | 30 | 0.050 |
| 4 | 40.93 | 27 | 0.042 |
| 5 | 21.89 | 26 | 0.695 |
| 6 | 53.29 | 30 | 0.006 |
| 7 | 22.17 | 30 | 0.848 |
| 8 | 25.42 | 30 | 0.705 |
| 9 | 62.76 | 27 | 0.000 |
| 10 | 36.86 | 30 | 0.181 |
| 11 | 239.08 | 99 | 0.000 |
| 12 | 51.12 | 30 | 0.010 |
| 13 | 66.30 | 30 | 0.000 |
| 14 | 26.31 | 30 | 0.660 |
| 15 | 38.93 | 27 | 0.064 |
| 16 | 49.52 | 30 | 0.014 |
| 17 | 58.99 | 30 | 0.001 |
| 18 | 21.59 | 30 | 0.869 |
| 19 | 32.07 | 30 | 0.364 |
| 20 | 27.16 | 30 | 0.615 |
| 21 | 40.91 | 30 | 0.088 |
| 22 | 41.22 | 30 | 0.083 |
| 23 | 68.07 | 30 | 0.000 |
| 24 | 59.43 | 30 | 0.001 |
| 25 | 137.45 | 90 | 0.001 |
| 26 | 126.98 | 90 | 0.006 |
| 27 | 30.89 | 20 | 0.057 |
| 28 | 35.44 | 30 | 0.227 |
| 29 | 24.90 | 28 | 0.634 |
| 30 | 25.76 | 26 | 0.477 |
| 31 | 33.20 | 27 | 0.190 |
| 32 | 376.24 | 97 | 0.000 |
| 33 | 26.01 | 30 | 0.675 |
| 34 | 62.83 | 30 | 0.000 |
| 35 | 63.55 | 30 | 0.000 |
| 36 | 39.26 | 30 | 0.120 |
| 37 | 34.48 | 30 | 0.262 |
| 38 | 37.44 | 30 | 0.164 |
| 39 | 108.53 | 88 | 0.068 |
| 40 | 53.02 | 30 | 0.006 |
| 41 | 28.14 | 30 | 0.563 |
| 42 | 26.33 | 30 | 0.658 |
| 43 | 33.88 | 27 | 0.169 |
| 44 | 60.36 | 30 | 0.001 |
| 45 | 40.95 | 30 | 0.088 |
| Total | 2616.83 | 1627 | 0.000 |

In addition to fit statistics, standardized residuals (SR) were calculated for each item using the FIT output file produced by Parscale. This type of analysis is much more insightful. Figure 6.2 shows the distribution of SR across all items. Table 6.4 displays the percentage of SR in each interval. The results suggest large proportions of SRs around the value of zero and suggest excellent model fit.

**Figure 6.2   Distribution of Standardized Residuals (SR)**
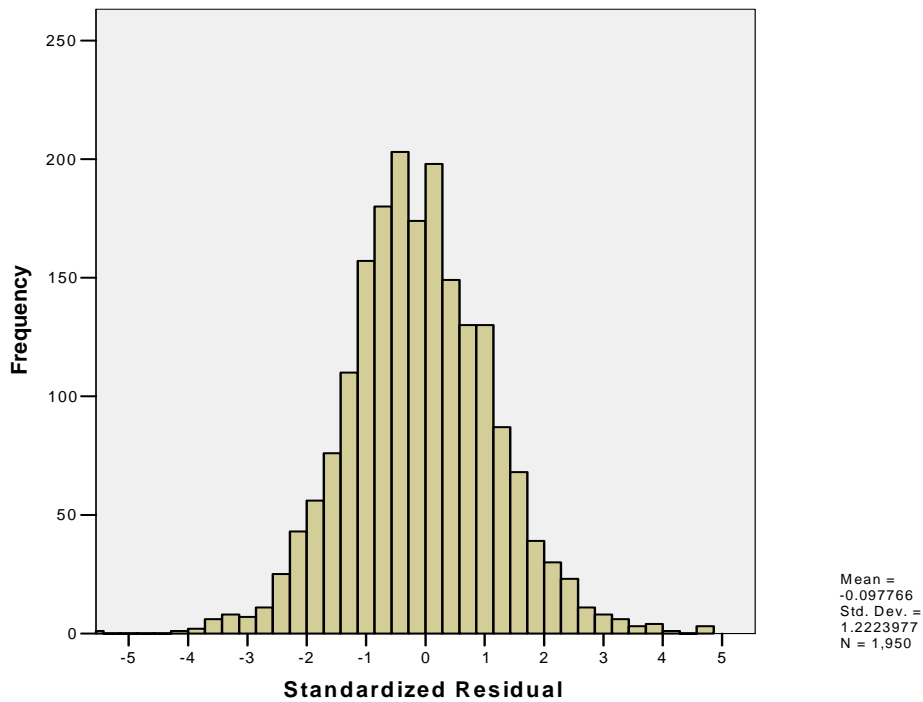


Mean =
-0.097766
Std. Dev. =
1.2223977
N = 1,950

**Table 6.4  Percentage of Standardized Residuals (SR), Organized in Intervals**

**along the Proficiency Continuum**

| SR Interval | Percentage |
|---|---|
| < -3 | 1% |
| (-3, -2) | 4% |
| (-2, -1) | 17% |
| (-1, 1) | 61% |
| (1, 2) | 13% |
| (2, 3) | 4% |
| >3 | 1% |

An item would show good fit with the model if its SR falls into the range (-2, 2), suggesting a 95% confidence when the model predicts the proportion of candidates who answer the item correctly. Table 6.4 shows that, out of a total of 1950 SRs (30 theta intervals for each dichotomous item, 150 for each polytomous item), there are 90% falling into the interval (-2, 2), and 98% falling into (-3, 3). As a result, the model shows very good fit for nearly all of the test items.  More information of SR can be found in Appendix C.

As it is widely recognized that fit statistics are easily affected by sample size, residual plots for the 45 items are provided in Appendix B to help illustrate the fit between model and data. There is one plot for each dichotomous item, which shows the correct response category. And five plots for each polytomous item, which show the response categories from 0 to 4. Small and random deviations off the curves are good indicators of model fit.  Again, the fit seems excellent.

Assessment of model fit can not only be checked at item level, but also at the test level. Comparison between observed and expected total score distributions provide another view of model fit. Figure 6.3 compares the distribution of observed total scores with the distribution of averaged predicted scores from 100 simulations, given the estimated item and ability parameters. The corresponding comparison of cumulative distributions is provided in Figure 6.4. Both figures show close approximations of the predicted score distribution to the observed distribution.

**Figure 6.3 Comparison of Total Score and Predicted Test Score Distributions**



**Figure 6.4  Comparison of Cumulative Score Distributions**



In summary, these analyses are highly reflective of model fit to the data.

## 7.  Test Information and Conditional Standard Errors

Test information and conditional standard errors are shown in Figure 7.1 and Figure

7.2. The test information is high from the middle through the higher end along the

proficiency continuum.  As a rough rule of thumb, information above 10 is highly desirable

and this is the case for proficiency scores between -1.0 and 2.0. The average student on the test is located around a scaled score of 0.0. Here the test information is high suggesting an above level of precision than is usually observed with achievement tests. This is excellent news and may suggest in the future that the addition of a few additional easy questions, could strengthen precision for lower performing candidates without taking anything significant away from proficiency estimation in the middle range of scores.

**Figure 7.1 Test Information Function**



Test Information Function

**Figure 7.2 Conditional Standard Errors**

Standard Error of Measurement



## 8. Identification of Differentially Functioning Items

Differential item functioning (DIF) analyses becomes a routine part in the analyses of large-scale assessments as an effort to enhance test fairness. DIF exists if individuals with the same ability, but from different subgroups, have different probabilities in answering the item correctly. DIF is different from bias in that more evidence needs to be collected to spot the sources of DIF. The sources of DIF could be the construct intended to measure, or could be construct-irrelevant factors. The latter one is often considered as bias.

It is important to check the means and standard deviations of test scores across different subgroups before completing DIF analyses. Table 8.1 shows the summary statistics of the scores for the subgroups.

**Table 8.1   Descriptive Statistics of the Test Scores**

|  | Group | Number | Minimum | Maximum | Mean | Std. Deviation |
|---|---|---|---|---|---|---|
| Gender | Female | 7559 | 1 | 60 | 30.56 | 11.97 |
|  | Male | 7762 | 1 | 60 | 30.33 | 12.89 |
| Ethnicity | White | 10568 | 1 | 60 | 33.34 | 11.41 |
|  | Hispanic | 1780 | 1 | 56 | 20.74 | 10.02 |
|  | Black | 1939 | 1 | 56 | 21.34 | 10.65 |
|  | Asian | 981 | 4 | 60 | 34.90 | 12.56 |
|  | Native | 44 | 8 | 55 | 28.73 | 12.60 |

In addition to score distributions, item parameters calibrated from different subgroups were plotted. Figure 8.3 to Figure 8.5 show the difficulty parameter (b) scatter plots between the groups, male/female, white/Black, and white/Hispanic.

**Figure 8.3   b-plot from Male and Female Samples**
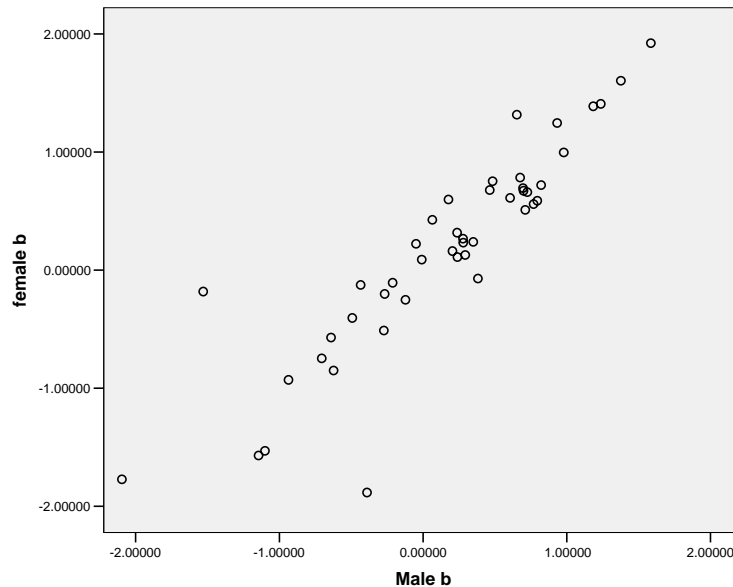
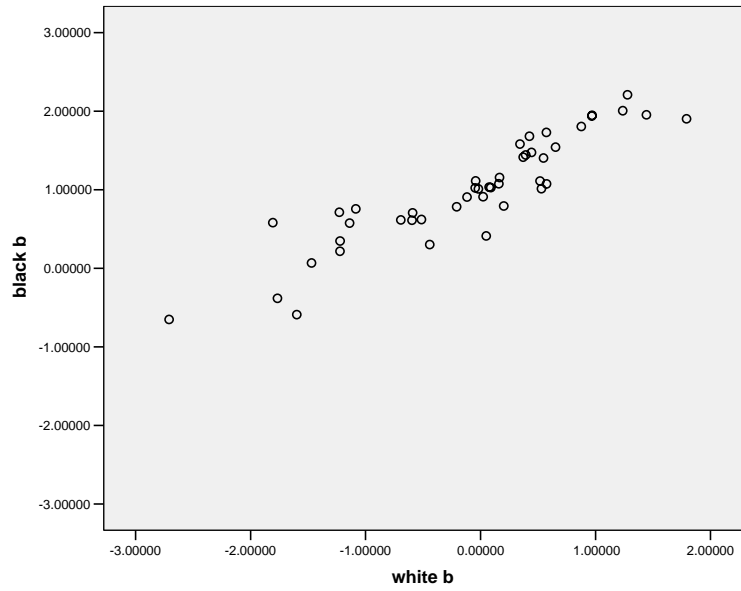**Figure 8.4   b-plot from White and Black Samples**



**Figure 8.5   b-plot from White and Hispanic Samples**



The statistics show that for gender groups, males and females achieve comparable scores, while for ethnicity groups, whites perform better than Blacks and Hispanics. In addition, item difficulty parameter plots show good linear relationship between the scores of

subgroups suggesting a minimum amount of DIF.   Clearly then we were not surprised by the results that follow using more traditional DIF procedures.

DIF analyses were carried out for three pairs of subgroups, male/female, white/Black, and white/Hispanic.  The sample sizes of Asian and Native groups were not large enough to do the DIF analyses. The weighted two-stage conditional p-value comparison procedure was used to calculate the DIF index by STDIF, a DOS-based program written by Frederic Robin (2001). The procedure consists of two stages. In Stage 1, for students in the focal and reference groups with the same total scores, the difference in proportion-correct at each score point were calculated and summed up. A statistic of unsigned DIF indices (UDIF) (all differences are treated as positive and contribute to DIF) was obtained for each item. The items with absolute values of UDIF larger than 0.075 were flagged as potential DIF items and were excluded from the total score in Stage2. In Stage 2, the comparisons were repeated based on the newly adjusted total scores. Items with absolute values of UDIF equal to or larger than 0.1 were flagged  and identified as DIF items .Table 8.2 shows the UDIF indices of the items for the three pair wise groups in Stage 1 and Stage 2.

**Table 8.2. UDIF Indices in Stage 1 and Stage 2***

| Item | UDIF(G) Stage 1 | UDIF(G) Stage 2 | UDIF(B) Stage 1 | UDIF(B) Stage 2 | UDIF(H) Stage 1 | UDIF(H) Stage 2 |
|---|---|---|---|---|---|---|
| 1 | 0.062 | 0.063 | 0.079** | 0.059 | 0.062 | 0.040 |
| 2 | 0.045 | 0.053 | 0.061 | 0.056 | 0.061 | 0.056 |
| 3 | 0.051 | 0.053 | 0.057 | 0.058 | 0.069 | -0.078 |
| 4 | 0.031 | 0.029 | 0.072 | 0.063 | 0.061 | -0.063 |
| 5 | 0.035 | 0.041 | 0.086** | 0.075 | 0.075 | 0.078 |
| 6 | -0.043 | -0.042 | -0.064 | -0.064 | -0.049 | -0.063 |
| 7 | 0.127 ** | 0.127 ** | 0.068 | 0.063 | 0.057 | -0.056 |
| 8 | -0.050 | -0.046 | -0.130** | -0.104** | -0.100** | -0.074 |
| 9 | 0.044 | 0.040 | -0.077** | 0.081 | 0.063 | -0.064 |
| 10 | -0.052 | -0.055 | -0.126** | -0.099 | -0.098** | -0.084 |
| 11 | -0.029 | -0.028 | 0.042 | 0.042 | 0.047 | 0.037 |
| 12 | 0.037 | 0.051 | -0.075 | -0.067 | -0.081** | -0.056 |
| 13 | -0.052 | -0.044 | -0.060 | -0.066 | -0.058 | -0.063 |
| 14 | 0.048 | 0.050 | 0.072 | 0.066 | 0.066 | -0.052 |
| 15 | 0.031 | 0.031 | 0.076** | 0.081 | 0.058 | 0.071 |
| 16 | 0.052 | 0.049 | 0.064 | 0.068 | 0.061 | 0.055 |
| 17 | -0.037 | -0.043 | -0.061 | -0.055 | 0.052 | 0.062 |
| 18 | -0.038 | -0.038 | -0.059 | -0.049 | -0.077** | 0.063 |
| 19 | 0.050 | 0.051 | -0.071 | -0.064 | -0.069 | -0.060 |
| 20 | -0.059 | -0.053 | -0.103** | -0.085 | -0.084** | -0.075 |
| 21 | -0.046 | -0.046 | -0.075 | -0.066 | -0.061 | -0.061 |
| 22 | 0.066 | 0.064 | -0.062 | -0.044 | -0.061 | -0.058 |
| 23 | 0.062 | 0.066 | 0.055 | -0.055 | 0.079** | 0.073 |
| 24 | 0.065 | 0.069 | -0.072 | -0.073 | -0.048 | -0.047 |
| 25 | -0.030 | -0.026 | -0.024 | -0.023 | 0.030 | -0.025 |
| 26 | -0.042 | -0.037 | -0.029 | -0.033 | 0.029 | -0.024 |
| 27 | 0.022 | 0.017 | 0.048 | -0.044 | -0.056 | -0.057 |
| 28 | 0.068 | 0.068 | 0.063 | -0.059 | 0.058 | 0.073 |
| 29 | 0.082 ** | 0.086 | 0.075 | 0.050 | -0.066 | -0.066 |
| 30 | 0.051 | 0.053 | 0.070 | 0.066 | 0.072 | 0.060 |
| 31 | 0.046 | 0.043 | 0.060 | -0.050 | -0.068 | -0.043 |
| 32 | -0.042 | -0.041 | 0.064 | 0.052 | 0.046 | 0.041 |
| 33 | 0.045 | 0.047 | 0.059 | -0.062 | -0.061 | -0.049 |
| 34 | -0.090 ** | -0.089 | -0.104** | -0.086 | -0.083** | -0.081 |
| 35 | -0.059 | -0.057 | -0.091** | -0.065 | -0.099** | -0.080 |
| 36 | 0.038 | 0.034 | 0.056 | 0.057 | 0.062 | 0.050 |
| 37 | 0.044 | 0.044 | 0.051 | 0.051 | 0.049 | 0.048 |
| 38 | -0.039 | -0.044 | -0.057 | -0.071 | 0.061 | 0.051 |
| 39 | -0.017 | -0.014 | 0.031 | 0.022 | 0.024 | 0.019 |
| 40 | -0.052 | -0.053 | 0.071 | 0.053 | 0.062 | -0.052 |
| 41 | 0.054 | 0.048 | 0.059 | 0.042 | -0.064 | -0.065 |
| 42 | 0.048 | 0.040 | -0.061 | -0.056 | 0.070 | 0.056 |
| 43 | -0.050 | -0.050 | -0.088** | -0.076 | -0.089** | -0.072 |
| 44 | 0.049 | 0.049 | 0.056 | 0.046 | -0.040 | -0.056 |
| 45 | 0.041 | 0.036 | 0.053 | -0.039 | -0.059 | -0.072 |

* G = Groups of Male/Female, B = Groups of White/Black, H = Groups of White/Hispanic.
** Flagged items where |UDIF| > 0.075 at Stage 1, and |UDIF| ≥ 0.1 at Stage 2.

Table 8.3 displays the distributions of UDIF indices at Stage 2. For gender groups, one item was flagged with UDIF larger than 0.1 favoring males over females. For the White/Black group comparison, one item was flagged with UDIF smaller than -0.1 favoring Blacks over whites. For the White/Hispanic group comparison, there was no item flagged showing high DIF.

**Table 8.3  Distribution of UDIF Indices in Stage 2**

|  | UDIF ≤ -.1 | -.1<UDIF <-.075 | -.075≤UDIF≤.075 | .075<UDIF<.1 | UDIF≥.1 |
|---|---|---|---|---|---|
|  | **High DIF** | **Low DIF** | **Non-DIF** | **Low DIF** | **High DIF** |
| **Gender** | 0 | 1(2.2%) | 42(93.3%) | 1(2.2%) | 1(2.2%) |
| **W/B** | 1(2.2%) | 4(8.9%) | 38(84.4%) | 2(4.4%) | 0 |
| **W/H** | 0 | 4(8.9%) | 40(88.9%) | 1(2.2%) | 0 |

Note: UDIF < 0 favoring focal groups (Female, Black or Hispanic). UDIF > 0 favoring reference groups (Male or White)

Figure 8.9 to Figure 8.11 illustrate the magnitudes and directions of UDIF indices across the 45 items.

**Figure 8.9   Plot of Male/Female UDIF Indices**



29

**Figure 8.10   Plot of White/Black UDIF Indices**



**Figure 8.11   Plot of White/Hispanic UDIF Indices**



The figures below are p-value plots conditional on adjusted total score and only two items were flagged from the analyses.  The plots for these potentially problematic items allow for a closer look at whether the DIF is uniform (consistent in direction), whether the magnitude is the same across test scores, or whether there is any interaction between the DIF and proficiency level (as evidenced by intersecting displays of data), etc.

**Figure 8.12  The Conditional p-value Plot of Flagged Item in the Gender Group**

**( UDIF = 0.127 )**



**Conditional p-value plot Item7**

**Figure 8.13  The Conditional p-value Plot of Flagged Item in White/Black Group**

**( UDIF= -0.104 )**



**Conditional p-value Item 8**

Both plots show more or less uniform conditional differences.  Also, the plot of ethnicity group comparisons showed more fluctuation than that of the gender group comparison. This is probably because of the smaller sample size at each score point in the ethnicity groups. In fact, the group sizes in the ethnicity groups are less than 20 at test scores

above 40. A further check of the item type shows that both items are multiple-choice items and the flagged item in the gender group comparison includes reference components. However, further studies are required in determining the causes of DIF, and whether the items are really biased or not. The detection of two problematic items in three DIF comparisons of 45 items each does not suggest anything approaching a problem with bias in the Introductory Physics Test.

## 9. Conclusions

In summary, the various analyses in this report suggest that the 2006 MCAS Grade 9/10 Introductory Physics Test is psychometrically sound—in fact, we would describe the statistical analyses as reflecting an excellent test. In classical item analyses, the items are of appropriate difficulty levels, properly discriminate between high and low performers, and the test scores are satisfactorily reliable for the multiple-choice items, performance items, and the total scores. From the item response theory (IRT) analyses, all of the indicators suggest a strong first factor. Also, the fits of the model to the data were excellent. The fit between model and data was checked in different ways using graphical procedures. The results show excellent model fit at both the item and test levels. In addition, the test information function shows that measurement precision is excellent across the scale. The differential item functioning (DIF) analyses show very little evidence of DIF. Only two items were identified, one from the male/female group comparison (favoring males), and the other from the white/Black group comparison (favoring Blacks). These two items suggest potential bias and need further checking. But the identification of only two items is not above the level that might be expected by chance. In summary, we believe our analyses have revealed that the Introductory Physics Test in 2006 is excellent in every statistical respect. We noticed that a slight shift in the test information function could be helpful, and our DIF analyses revealed that two items might be checked further, not because of their small to trivial impact on the

test scores in 2006 but because something might be learned to eliminate these small problems

from future tests.

**References**

Brooks, G. P., & Johanson, G. A. (2003). Test analysis program. *Applied Psychological Measurement, 27*, 305-306.

Hambleton, R. K., Zhao, Y., Smith, Z., Lam, W., & Deng, N. (2008). *Psychometric analyses of the 2006 MCAS high school science tests* (Center for Educational Assessment Research Report No. 649). Amherst, MA: University of Massachusetts, Center for Educational Assessment.

Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage.

Massachusetts Department of Education. (2006). *The Massachusetts Science and Technology/Engineering Curriculum Framework*. Malden, MA: Massachusetts Department of Education.

Robin, F. (2001). *STDIF: Standardization-DIF analysis program* [Computer program]. Amherst, MA: University of Massachusetts, Center for Educational Assessment.

Zenisky, A. L., & Hambleton, R. K. (2007). *Differential item functioning analyses with STDIF: User's guide* (Unpublished report). Amherst, MA: University of Massachusetts, Center for Educational Assessment.

Zenisky, A., Hambleton, R. K. & Robin, F. (2003). Detection of differential item functioning in large-scale state assessments: A study evaluating a two-stage approach. *Educational and Psychological Measurement, 63*(1), 51-64.

Zenisky, A., Hambleton, R. K., & Robin, F. (2004). DIF detection and interpretation in large-scale science assessments: Informing item writing practices. *Educational Assessment, 9*(1&2), 61-78.

# Appendix A.  A Detailed Distracter Analysis with Frequencies and Percentages

* is keyed answer, # is option that discriminates better than keyed answer

| Item | Group | A | B | C | D |
|---|---|---|---|---|---|
| 1 | TOTAL | 4437*(0.844) | 89 (0.017) | 627 (0.119) | 83 (0.016) |
|  | High | 1648 (0.946) | 4 (0.002) | 81 (0.046) | 8 (0.005) |
|  | Low | 1054 (0.700) | 71 (0.047) | 306 (0.203) | 59 (0.039) |
|  | Diff | 594 (0.246) | -67(-0.045) | -225(-0.157) | -51(-0.035) |
|  |  |  |  |  |  |
| 2 | TOTAL | 891 (0.170) | 2285*(0.435) | 1072 (0.204) | 966 (0.184) |
|  | High | 202 (0.116) | 1221 (0.701) | 161 (0.092) | 155 (0.089) |
|  | Low | 313 (0.208) | 337 (0.224) | 493 (0.328) | 335 (0.223) |
|  | Diff | -111(-0.092) | 884 (0.477) | -332(-0.235) | -180(-0.134) |
|  |  |  |  |  |  |
| 3 | TOTAL | 2765*(0.526) | 1307 (0.249) | 656 (0.125) | 485 (0.092) |
|  | High | 1426 (0.819) | 132 (0.076) | 111 (0.064) | 68 (0.039) |
|  | Low | 374 (0.249) | 644 (0.428) | 269 (0.179) | 192 (0.128) |
|  | Diff | 1052 (0.570) | -512(-0.352) | -158(-0.115) | -124(-0.089) |
|  |  |  |  |  |  |
| 4 | TOTAL | 458 (0.087) | 3756*(0.715) | 670 (0.128) | 338 (0.064) |
|  | High | 16 (0.009) | 1621 (0.931) | 96 (0.055) | 9 (0.005) |
|  | Low | 285 (0.189) | 655 (0.435) | 286 (0.190) | 251 (0.167) |
|  | Diff | -269(-0.180) | 966 (0.495) | -190(-0.135) | -242(-0.162) |
|  |  |  |  |  |  |
| 5 | TOTAL | 537 (0.102) | 390 (0.074) | 3841*(0.731) | 444 (0.085) |
|  | High | 31 (0.018) | 18 (0.010) | 1658 (0.952) | 32 (0.018) |
|  | Low | 308 (0.205) | 252 (0.167) | 651 (0.433) | 259 (0.172) |
|  | Diff | -277(-0.187) | -234(-0.157) | 1007 (0.519) | -227(-0.154) |
|  |  |  |  |  |  |
| 6 | TOTAL | 385 (0.073) | 332 (0.063) | 447 (0.085) | 4057*(0.772) |
|  | High | 31 (0.018) | 20 (0.011) | 89 (0.051) | 1601 (0.919) |
|  | Low | 237 (0.157) | 196 (0.130) | 206 (0.137) | 837 (0.556) |
|  | Diff | -206(-0.140) | -176(-0.119) | -117(-0.086) | 764 (0.363) |
|  |  |  |  |  |  |
| 7 | TOTAL | 2626 (0.500) | 207 (0.039) | 303 (0.058) | 2078*(0.396) |
|  | High | 626 (0.359) | 14 (0.008) | 34 (0.020) | 1065 (0.611) |
|  | Low | 856 (0.569) | 135 (0.090) | 153 (0.102) | 328 (0.218) |
|  | Diff | -230(-0.209) | -121(-0.082) | -119(-0.082) | 737 (0.393) |
|  |  |  |  |  |  |
| 8 | TOTAL | 2607*(0.496) | 1074 (0.204) | 972 (0.185) | 538 (0.102) |
|  | High | 1275 (0.732) | 217 (0.125) | 144 (0.083) | 97 (0.056) |
|  | Low | 432 (0.287) | 396 (0.263) | 406 (0.270) | 227 (0.151) |
|  | Diff | 843 (0.445) | -179(-0.139) | -262(-0.187) | -130(-0.095) |
|  |  |  |  |  |  |
| 9 | TOTAL | 361 (0.069) | 652 (0.124) | 225 (0.043) | 3976*(0.757) |
|  | High | 22 (0.013) | 54 (0.031) | 15 (0.009) | 1650 (0.947) |
|  | Low | 225 (0.150) | 362 (0.241) | 159 (0.106) | 721 (0.479) |

|    |       |               |               |               |               |
|----|-------|---------------|---------------|---------------|---------------|
|    | Diff  | -203(-0.137)  | -308(-0.210)  | -144(-0.097)  | 929 (0.468)   |
| 10 | TOTAL | 2605*(0.496)  | 1333 (0.254)  | 865 (0.165)   | 393 (0.075)   |
|    | High  | 1317 (0.756)  | 291 (0.167)   | 99 (0.057)    | 31 (0.018)    |
|    | Low   | 445 (0.296)   | 399 (0.265)   | 407 (0.270)   | 213 (0.142)   |
|    | Diff  | 872 (0.460)   | -108(-0.098)  | -308(-0.214)  | -182(-0.124)  |
|    |       |               |               |               |               |
| 12 | TOTAL | 663 (0.126)   | 493 (0.094)   | 1168 (0.222)  | 2873*(0.547)  |
|    | High  | 83 (0.048)    | 47 (0.027)    | 255 (0.146)   | 1356 (0.778)  |
|    | Low   | 314 (0.209)   | 278 (0.185)   | 400 (0.266)   | 462 (0.307)   |
|    | Diff  | -231(-0.161)  | -231(-0.158)  | -145(-0.119)  | 894 (0.471)   |
|    |       |               |               |               |               |
| 13 | TOTAL | 3143*(0.598)  | 1288 (0.245)  | 412 (0.078)   | 356 (0.068)   |
|    | High  | 1277 (0.733)  | 358 (0.206)   | 55 (0.032)    | 50 (0.029)    |
|    | Low   | 622 (0.413)   | 439 (0.292)   | 216 (0.144)   | 178 (0.118)   |
|    | Diff  | 655 (0.320)   | -81(-0.086)   | -161(-0.112)  | -128(-0.090)  |
|    |       |               |               |               |               |
| 14 | TOTAL | 472 (0.090)   | 2451*(0.467)  | 1329 (0.253)  | 946 (0.180)   |
|    | High  | 63 (0.036)    | 1277 (0.733)  | 234 (0.134)   | 166 (0.095)   |
|    | Low   | 257 (0.171)   | 366 (0.243)   | 462 (0.307)   | 370 (0.246)   |
|    | Diff  | -194(-0.135)  | 911 (0.490)   | -228(-0.173)  | -204(-0.151)  |
|    |       |               |               |               |               |
| 15 | TOTAL | 523 (0.100)   | 3667*(0.698)  | 646 (0.123)   | 359 (0.068)   |
|    | High  | 48 (0.028)    | 1611 (0.925)  | 54 (0.031)    | 27 (0.015)    |
|    | Low   | 310 (0.206)   | 583 (0.387)   | 339 (0.225)   | 220 (0.146)   |
|    | Diff  | -262(-0.178)  | 1028 (0.537)  | -285(-0.194)  | -193(-0.131)  |
|    |       |               |               |               |               |
| 16 | TOTAL | 2979*(0.567)  | 274 (0.052)   | 384 (0.073)   | 1552 (0.295)  |
|    | High  | 1290 (0.741)  | 10 (0.006)    | 11 (0.006)    | 430 (0.247)   |
|    | Low   | 512 (0.340)   | 214 (0.142)   | 276 (0.183)   | 445 (0.296)   |
|    | Diff  | 778 (0.400)   | -204(-0.136)  | -265(-0.177)  | -15(-0.049)   |
|    |       |               |               |               |               |
| 17 | TOTAL | 667 (0.127)   | 823 (0.157)   | 1073 (0.204)  | 2616*(0.498)  |
|    | High  | 139 (0.080)   | 213 (0.122)   | 212 (0.122)   | 1175 (0.675)  |
|    | Low   | 280 (0.186)   | 279 (0.185)   | 454 (0.302)   | 426 (0.283)   |
|    | Diff  | -141(-0.106)  | -66(-0.063)   | -242(-0.180)  | 749 (0.391)   |
|    |       |               |               |               |               |
| 18 | TOTAL | 642 (0.122)   | 1604 (0.305)  | 2550*(0.485)  | 390 (0.074)   |
|    | High  | 106 (0.061)   | 379 (0.218)   | 1210 (0.695)  | 45 (0.026)    |
|    | Low   | 299 (0.199)   | 515 (0.342)   | 425 (0.282)   | 206 (0.137)   |
|    | Diff  | -193(-0.138)  | -136(-0.125)  | 785 (0.412)   | -161(-0.111)  |
|    |       |               |               |               |               |
| 19 | TOTAL | 995 (0.189)   | 565 (0.108)   | 3079*(0.586)  | 539 (0.103)   |
|    | High  | 210 (0.121)   | 37 (0.021)    | 1405 (0.807)  | 87 (0.050)    |
|    | Low   | 337 (0.224)   | 352 (0.234)   | 503 (0.334)   | 244 (0.162)   |
|    | Diff  | -127(-0.103)  | -315(-0.213)  | 902 (0.472)   | -157(-0.112)  |
|    |       |               |               |               |               |
| 20 | TOTAL | 1112 (0.212)  | 2438*(0.464)  | 658 (0.125)   | 966 (0.184)   |

|  |  |  |  |  |  |
|---|---|---|---|---|---|
|  | High | 369 (0.212) | 1200 (0.689) | 53 (0.030) | 118 (0.068) |
|  | Low | 323 (0.215) | 379 (0.252) | 341 (0.227) | 391 (0.260) |
|  | Diff | 46(-0.003) | 821 (0.437) | -288(-0.196) | -273(-0.192) |
|  |  |  |  |  |  |
| 21 | TOTAL | 2694*(0.513) | 978 (0.186) | 1113 (0.212) | 370 (0.070) |
|  | High | 1187 (0.681) | 207 (0.119) | 299 (0.172) | 43 (0.025) |
|  | Low | 463 (0.308) | 384 (0.255) | 355 (0.236) | 219 (0.146) |
|  | Diff | 724 (0.374) | -177(-0.136) | -56(-0.064) | -176(-0.121) |
|  |  |  |  |  |  |
| 22 | TOTAL | 1266 (0.241) | 1043 (0.199) | 1912*(0.364) | 941 (0.179) |
|  | High | 383 (0.220) | 188 (0.108) | 1010 (0.580) | 160 (0.092) |
|  | Low | 382 (0.254) | 365 (0.243) | 325 (0.216) | 347 (0.231) |
|  | Diff | 1(-0.034) | -177(-0.135) | 685 (0.364) | -187(-0.139) |
|  |  |  |  |  |  |
| 23 | TOTAL | 433 (0.082) | 572 (0.109) | 1051 (0.200) | 3108*(0.592) |
|  | High | 17 (0.010) | 50 (0.029) | 268 (0.154) | 1404 (0.806) |
|  | Low | 260 (0.173) | 326 (0.217) | 371 (0.247) | 469 (0.312) |
|  | Diff | -243(-0.163) | -276(-0.188) | -103(-0.093) | 935 (0.494) |
| 24 | TOTAL | 436 (0.083) | 1169 (0.222) | 688 (0.131) | 2861*(0.545) |
|  | High | 73 (0.042) | 166 (0.095) | 89 (0.051) | 1407 (0.808) |
|  | Low | 218 (0.145) | 450 (0.299) | 331 (0.220) | 421 (0.280) |
|  | Diff | -145(-0.103) | -284(-0.204) | -242(-0.169) | 986 (0.528) |
|  |  |  |  |  |  |
| 27 | TOTAL | 4585*(0.873) | 210 (0.040) | 156 (0.030) | 164 (0.031) |
|  | High | 1727 (0.991) | 6 (0.003) | 1 (0.001) | 8 (0.005) |
|  | Low | 953 (0.633) | 164 (0.109) | 139 (0.092) | 116 (0.077) |
|  | Diff | 774 (0.358) | -158(-0.106) | -138(-0.092) | -108(-0.072) |
|  |  |  |  |  |  |
| 28 | TOTAL | 1175 (0.224) | 581 (0.111) | 2533*(0.482) | 809 (0.154) |
|  | High | 303 (0.174) | 74 (0.042) | 1193 (0.685) | 168 (0.096) |
|  | Low | 348 (0.231) | 325 (0.216) | 388 (0.258) | 299 (0.199) |
|  | Diff | -45(-0.057) | -251(-0.173) | 805 (0.427) | -131(-0.102) |
|  |  |  |  |  |  |
| 29 | TOTAL | 3387*(0.645) | 730 (0.139) | 682 (0.130) | 299 (0.057) |
|  | High | 1577 (0.905) | 92 (0.053) | 60 (0.034) | 13 (0.007) |
|  | Low | 466 (0.310) | 348 (0.231) | 340 (0.226) | 204 (0.136) |
|  | Diff | 1111 (0.596) | -256(-0.178) | -280(-0.191) | -191(-0.128) |
|  |  |  |  |  |  |
| 30 | TOTAL | 433 (0.082) | 644 (0.123) | 685 (0.130) | 3320*(0.632) |
|  | High | 25 (0.014) | 38 (0.022) | 39 (0.022) | 1639 (0.941) |
|  | Low | 245 (0.163) | 347 (0.231) | 385 (0.256) | 370 (0.246) |
|  | Diff | -220(-0.148) | -309(-0.209) | -346(-0.233) | 1269 (0.695) |
|  |  |  |  |  |  |
| 31 | TOTAL | 751 (0.143) | 619 (0.118) | 3177*(0.605) | 545 (0.104) |
|  | High | 91 (0.052) | 42 (0.024) | 1571 (0.902) | 37 (0.021) |
|  | Low | 328 (0.218) | 327 (0.217) | 405 (0.269) | 294 (0.195) |
|  | Diff | -237(-0.166) | -285(-0.193) | 1166 (0.633) | -257(-0.174) |

| | | | | | |
|---|---|---|---|---|---|
| 33 | TOTAL | 838 (0.159) | 1987*(0.378) | 1567 (0.298) | 642 (0.122) |
| | High | 175 (0.100) | 964 (0.553) | 513 (0.294) | 90 (0.052) |
| | Low | 329 (0.219) | 351 (0.233) | 397 (0.264) | 235 (0.156) |
| | Diff | -154(-0.118) | 613 (0.320) | 116 (0.031) | -145(-0.104) |
| | | | | | |
| 34 | TOTAL | 468 (0.089) | 1176 (0.224) | 473 (0.090) | 2907*(0.553) |
| | High | 95 (0.055) | 174 (0.100) | 26 (0.015) | 1447 (0.831) |
| | Low | 206 (0.137) | 414 (0.275) | 302 (0.201) | 381 (0.253) |
| | Diff | -111(-0.082) | -240(-0.175) | -276(-0.186) | 1066 (0.577) |
| | | | | | |
| 35 | TOTAL | 332 (0.063) | 2593*(0.494) | 1589 (0.302) | 505 (0.096) |
| | High | 24 (0.014) | 1326 (0.761) | 337 (0.193) | 54 (0.031) |
| | Low | 212 (0.141) | 404 (0.268) | 448 (0.298) | 238 (0.158) |
| | Diff | -188(-0.127) | 922 (0.493) | -111(-0.104) | -184(-0.127) |
| | | | | | |
| 36 | TOTAL | 1988 (0.378) | 764 (0.145) | 904 (0.172) | 1334*(0.254) |
| | High | 542 (0.311) | 128 (0.073) | 272 (0.156) | 789 (0.453) |
| | Low | 513 (0.341) | 294 (0.195) | 266 (0.177) | 223 (0.148) |
| | Diff | 29(-0.030) | -166(-0.122) | 6(-0.021) | 566 (0.305) |
| | | | | | |
| 37 | TOTAL | 974 (0.185) | 1812 (0.345) | 816 (0.155) | 1406*(0.268) |
| | High | 296 (0.170) | 483 (0.277) | 238 (0.137) | 721 (0.414) |
| | Low | 274 (0.182) | 535 (0.355) | 276 (0.183) | 217 (0.144) |
| | Diff | 22(-0.012) | -52(-0.078) | -38(-0.047) | 504 (0.270) |
| | | | | | |
| 38 | TOTAL | 847 (0.161) | 2519*(0.479) | 835 (0.159) | 800 (0.152) |
| | High | 135 (0.077) | 1361 (0.781) | 142 (0.082) | 98 (0.056) |
| | Low | 318 (0.211) | 300 (0.199) | 326 (0.217) | 354 (0.235) |
| | Diff | -183(-0.134) | 1061 (0.582) | -184(-0.135) | -256(-0.179) |
| | | | | | |
| 40 | TOTAL | 965 (0.184) | 737 (0.140) | 2227*(0.424) | 751 (0.143) |
| | High | 325 (0.187) | 139 (0.080) | 1037 (0.595) | 195 (0.112) |
| | Low | 269 (0.179) | 284 (0.189) | 336 (0.223) | 243 (0.161) |
| | Diff | 56 (0.008) | -145(-0.109) | 701 (0.372) | -48(-0.050) |
| | | | | | |
| 41 | TOTAL | 665 (0.127) | 2243*(0.427) | 817 (0.156) | 988 (0.188) |
| | High | 132 (0.076) | 1159 (0.665) | 206 (0.118) | 202 (0.116) |
| | Low | 252 (0.167) | 299 (0.199) | 283 (0.188) | 316 (0.210) |
| | Diff | -120(-0.092) | 860 (0.467) | -77(-0.070) | -114(-0.094) |
| | | | | | |
| 42 | TOTAL | 2776*(0.528) | 761 (0.145) | 822 (0.156) | 331 (0.063) |
| | High | 1414 (0.812) | 68 (0.039) | 191 (0.110) | 21 (0.012) |
| | Low | 312 (0.207) | 361 (0.240) | 278 (0.185) | 193 (0.128) |
| | Diff | 1102 (0.604) | -293(-0.201) | -87(-0.075) | -172(-0.116) |
| | | | | | |
| 43 | TOTAL | 419 (0.080) | 596 (0.113) | 3271*(0.623) | 412 (0.078) |

|     |       |              |              |               |              |
| --- | ----- | ------------ | ------------ | ------------- | ------------ |
|     | High  | 28 (0.016)   | 39 (0.022)   | 1604 (0.921)  | 26 (0.015)   |
|     | Low   | 222 (0.148)  | 286 (0.190)  | 412 (0.274)   | 219 (0.146)  |
|     | Diff  | -194(-0.131) | -247(-0.168) | 1192 (0.647)  | -193(-0.131) |
|     |       |              |              |               |              |
| 44  | TOTAL | 843 (0.160)  | 954 (0.182)  | 1515 (0.288)  | 1388*(0.264) |
|     | High  | 188 (0.108)  | 167 (0.096)  | 569 (0.327)   | 773 (0.444)  |
|     | Low   | 271 (0.180)  | 331 (0.220)  | 327 (0.217)   | 216 (0.144)  |
|     | Diff  | -83(-0.072)  | -164(-0.124) | 242 (0.109)   | 557 (0.300)  |
|     |       |              |              |               |              |
| 45  | TOTAL | 990 (0.188)  | 565 (0.108)  | 2678*(0.510)  | 474 (0.090)  |
|     | High  | 172 (0.099)  | 47 (0.027)   | 1421 (0.816)  | 55 (0.032)   |
|     | Low   | 362 (0.241)  | 262 (0.174)  | 328 (0.218)   | 201 (0.134)  |
|     | Diff  | -190(-0.142) | -215(-0.147) | 1093 (0.598)  | -146(-0.102) |

# Appendix B   Item Residual Plots

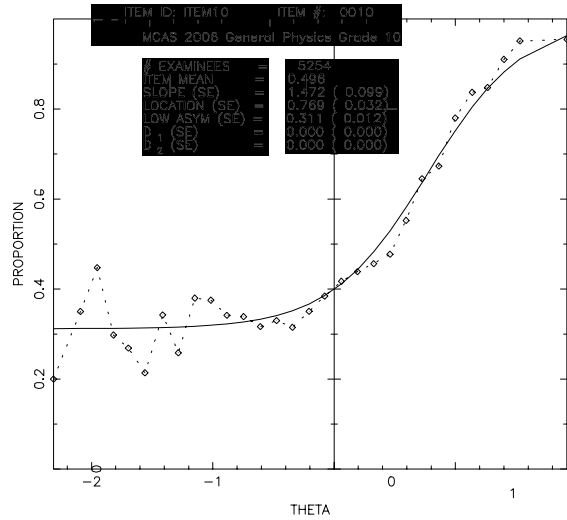LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM7        ITEM #: 0007

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.396 |
| SLOPE (SE) | = | 0.814 ( 0.071) |
| LOCATION (SE) | = | 1.149 ( 0.054) |
| LOW ASYM (SE) | = | 0.212 ( 0.017) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM8        ITEM #: 0008

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.496 |
| SLOPE (SE) | = | 1.098 ( 0.081) |
| LOCATION (SE) | = | 0.791 ( 0.043) |
| LOW ASYM (SE) | = | 0.288 ( 0.015) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM9        ITEM #: 0009

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.757 |
| SLOPE (SE) | = | 0.812 ( 0.042) |
| LOCATION (SE) | = | -0.892 ( 0.088) |
| LOW ASYM (SE) | = | 0.149 ( 0.039) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM10        ITEM #: 0010

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.496 |
| SLOPE (SE) | = | 1.472 ( 0.099) |
| LOCATION (SE) | = | 0.769 ( 0.032) |
| LOW ASYM (SE) | = | 0.311 ( 0.012) |
| D 1 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM11        ITEM #: 0011

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 2.616 |
| SLOPE (SE) | = | 0.985 ( 0.017) |
| LOCATION (SE) | = | -0.588 ( 0.017) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.698 ( 0.025) |
| D 2 (SE) | = | 0.417 ( 0.023) |
| D 3 (SE) | = | -0.208 ( 0.020) |
| D 4 (SE) | = | -0.909 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM11        ITEM #: 0011

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 2.616 |
| SLOPE (SE) | = | 0.985 ( 0.017) |
| LOCATION (SE) | = | -0.588 ( 0.017) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.698 ( 0.025) |
| D 2 (SE) | = | 0.417 ( 0.023) |
| D 3 (SE) | = | -0.208 ( 0.020) |
| D 4 (SE) | = | -0.909 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM11    ITEM #: 0011
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 2.616 |
| SLOPE (SE) = | 0.985 ( 0.017) |
| LOCATION (SE) = | -0.588 ( 0.017) |
| LOW ASYM (SE) = | 0.000 ( 0.000) |
| D 1 (SE) = | 0.698 ( 0.025) |
| D 2 (SE) = | 0.417 ( 0.023) |
| D 3 (SE) = | -0.206 ( 0.020) |
| D 4 (SE) = | -0.909 ( 0.020) |
| D 5 (SE) = | 0.000 ( 0.000) |

THETA



LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM11    ITEM #: 0011
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 2.616 |
| SLOPE (SE) = | 0.985 ( 0.017) |
| LOCATION (SE) = | -0.588 ( 0.017) |
| LOW ASYM (SE) = | 0.000 ( 0.000) |
| D 1 (SE) = | 0.698 ( 0.025) |
| D 2 (SE) = | 0.417 ( 0.023) |
| D 3 (SE) = | -0.206 ( 0.020) |
| D 4 (SE) = | -0.909 ( 0.020) |
| D 5 (SE) = | 0.000 ( 0.000) |

PROPORTION

THETA



LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM11    ITEM #: 0011
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 2.616 |
| SLOPE (SE) = | 0.985 ( 0.017) |
| LOCATION (SE) = | -0.588 ( 0.017) |
| LOW ASYM (SE) = | 0.000 ( 0.000) |
| D 1 (SE) = | 0.698 ( 0.025) |
| D 2 (SE) = | 0.417 ( 0.023) |
| D 3 (SE) = | -0.206 ( 0.020) |
| D 4 (SE) = | -0.909 ( 0.020) |
| D 5 (SE) = | 0.000 ( 0.000) |

THETA



LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM12    ITEM #: 0012
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 0.547 |
| SLOPE (SE) = | 0.672 ( 0.048) |
| LOCATION (SE) = | 0.224 ( 0.084) |
| LOW ASYM (SE) = | 0.174 ( 0.030) |
| D 1 (SE) = | 0.000 ( 0.000) |
| D 2 (SE) = | 0.000 ( 0.000) |

PROPORTION

THETA



LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM13    ITEM #: 0013
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 0.598 |
| SLOPE (SE) = | 0.352 ( 0.030) |
| LOCATION (SE) = | 0.287 ( 0.231) |
| LOW ASYM (SE) = | 0.134 ( 0.052) |
| D 1 (SE) = | 0.000 ( 0.000) |
| D 2 (SE) = | 0.000 ( 0.000) |

THETA



LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM14    ITEM #: 0014
MCAS 2006 General Physics Grade 10

| | |
|---|---|
| # EXAMINEES = | 5254 |
| ITEM MEAN | 0.487 |
| SLOPE (SE) = | 1.128 ( 0.073) |
| LOCATION (SE) = | 0.697 ( 0.037) |
| LOW ASYM (SE) = | 0.232 ( 0.014) |
| D 1 (SE) = | 0.000 ( 0.000) |
| D 2 (SE) = | 0.000 ( 0.000) |

PROPORTION

THETA

42

ITEM ID: ITEM15    ITEM #:  0015
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.698 | |
| SLOPE (SE) = | 0.929 | ( 0.050) |
| LOCATION (SE) = | -0.439 | ( 0.069) |
| LOW ASYM (SE) = | 0.201 | ( 0.030) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

THETA

ITEM ID: ITEM16    ITEM #:  0016
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.567 | |
| SLOPE (SE) = | 0.431 | ( 0.029) |
| LOCATION (SE) = | -0.143 | ( 0.130) |
| LOW ASYM (SE) = | 0.093 | ( 0.036) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

PROPORTION
THETA

ITEM ID: ITEM17    ITEM #:  0017
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.498 | |
| SLOPE (SE) = | 0.449 | ( 0.034) |
| LOCATION (SE) = | 0.314 | ( 0.115) |
| LOW ASYM (SE) = | 0.093 | ( 0.033) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

THETA

ITEM ID: ITEM18    ITEM #:  0018
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.485 | |
| SLOPE (SE) = | 0.720 | ( 0.060) |
| LOCATION (SE) = | 0.731 | ( 0.069) |
| LOW ASYM (SE) = | 0.227 | ( 0.023) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

PROPORTION
THETA

ITEM ID: ITEM19    ITEM #:  0019
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.588 | |
| SLOPE (SE) = | 0.669 | ( 0.048) |
| LOCATION (SE) = | 0.065 | ( 0.094) |
| LOW ASYM (SE) = | 0.192 | ( 0.033) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

THETA

ITEM ID: ITEM20    ITEM #:  0020
MCAS 2006 General Physics Grade 10

| | | |
|---|---|---|
| # EXAMINEES = | 5254 | |
| ITEM MEAN = | 0.464 | |
| SLOPE (SE) = | 0.761 | ( 0.061) |
| LOCATION (SE) = | 0.768 | ( 0.060) |
| LOW ASYM (SE) = | 0.211 | ( 0.021) |
| D 1 (SE) = | 0.000 | ( 0.000) |
| D 2 (SE) = | 0.000 | ( 0.000) |

PROPORTION
THETA

ITEM ID: ITEM21      ITEM #:  0021
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.513 |
| SLOPE (SE) | = | 0.451 ( 0.040) |
| LOCATION (SE) | = | 0.368 ( 0.142) |
| LOW ASYM (SE) | = | 0.134 ( 0.040) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM22      ITEM #:  0022
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.384 |
| SLOPE (SE) | = | 1.397 ( 0.106) |
| LOCATION (SE) | = | 1.163 ( 0.034) |
| LOW ASYM (SE) | = | 0.237 ( 0.010) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM23      ITEM #:  0023
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.592 |
| SLOPE (SE) | = | 0.612 ( 0.030) |
| LOCATION (SE) | = | -0.294 ( 0.069) |
| LOW ASYM (SE) | = | 0.066 ( 0.024) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM24      ITEM #:  0024
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.545 |
| SLOPE (SE) | = | 1.005 ( 0.065) |
| LOCATION (SE) | = | 0.428 ( 0.048) |
| LOW ASYM (SE) | = | 0.257 ( 0.019) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM25      ITEM #:  0025
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 1.732 |
| SLOPE (SE) | = | 1.406 ( 0.022) |
| LOCATION (SE) | = | 0.268 ( 0.012) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.803 ( 0.017) |
| D 2 (SE) | = | 0.336 ( 0.015) |
| D 3 (SE) | = | -0.139 ( 0.015) |
| D 4 (SE) | = | -1.000 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM25      ITEM #:  0025
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 1.732 |
| SLOPE (SE) | = | 1.406 ( 0.022) |
| LOCATION (SE) | = | 0.268 ( 0.012) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.803 ( 0.017) |
| D 2 (SE) | = | 0.336 ( 0.015) |
| D 3 (SE) | = | -0.139 ( 0.015) |
| D 4 (SE) | = | -1.000 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

44

# EXAMINEES = 5254
ITEM MEAN = 1.732
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.268 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.803 ( 0.017)
D 2 (SE) = 0.336 ( 0.015)
D 3 (SE) = -0.139 ( 0.015)
D 4 (SE) = -1.000 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED
ITEM ID: ITEM25    ITEM #: 0025
MCAS 2006 General Physics Grade 10

# EXAMINEES = 5254
ITEM MEAN = 1.732
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.268 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.803 ( 0.017)
D 2 (SE) = 0.336 ( 0.015)
D 3 (SE) = -0.139 ( 0.015)
D 4 (SE) = -1.000 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED
ITEM ID: ITEM25    ITEM #: 0025
MCAS 2006 General Physics Grade 10

# EXAMINEES = 5254
ITEM MEAN = 1.732
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.268 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.803 ( 0.017)
D 2 (SE) = 0.336 ( 0.015)
D 3 (SE) = -0.139 ( 0.015)
D 4 (SE) = -1.000 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED
ITEM ID: ITEM25    ITEM #: 0025
MCAS 2006 General Physics Grade 10

# EXAMINEES = 5254
ITEM MEAN = 1.707
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.302 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.781 ( 0.016)
D 2 (SE) = 0.414 ( 0.015)
D 3 (SE) = -0.167 ( 0.015)
D 4 (SE) = -1.027 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED
ITEM ID: ITEM25    ITEM #: 0025
MCAS 2006 General Physics Grade 10

# EXAMINEES = 5254
ITEM MEAN = 1.707
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.302 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.781 ( 0.016)
D 2 (SE) = 0.414 ( 0.015)
D 3 (SE) = -0.167 ( 0.015)
D 4 (SE) = -1.027 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED
ITEM ID: ITEM25    ITEM #: 0025
MCAS 2006 General Physics Grade 10

# EXAMINEES = 5254
ITEM MEAN = 1.707
SLOPE (SE) = 1.406 ( 0.022)
LOCATION (SE) = 0.302 ( 0.012)
LOW ASYM (SE) = 0.000 ( 0.000)
D 1 (SE) = 0.781 ( 0.016)
D 2 (SE) = 0.414 ( 0.015)
D 3 (SE) = -0.167 ( 0.015)
D 4 (SE) = -1.027 ( 0.020)
D 5 (SE) = 0.000 ( 0.000)

PROPORTION

THETA

ITEM ID: ITEM26    ITEM #:  0026
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 1.707 |
| SLOPE (SE) | = | 1.406 ( 0.022) |
| LOCATION (SE) | = | 0.302 ( 0.012) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.781 ( 0.016) |
| D 2 (SE) | = | 0.414 ( 0.015) |
| D 3 (SE) | = | −0.167 ( 0.015) |
| D 4 (SE) | = | −1.027 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM26    ITEM #:  0026
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 1.707 |
| SLOPE (SE) | = | 1.406 ( 0.022) |
| LOCATION (SE) | = | 0.302 ( 0.012) |
| LOW ASYM (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.781 ( 0.016) |
| D 2 (SE) | = | 0.414 ( 0.015) |
| D 3 (SE) | = | −0.167 ( 0.015) |
| D 4 (SE) | = | −1.027 ( 0.020) |
| D 5 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM27    ITEM #:  0027
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.873 |
| SLOPE (SE) | = | 1.394 ( 0.078) |
| LOCATION (SE) | = | −1.124 ( 0.066) |
| LOW ASYM (SE) | = | 0.313 ( 0.037) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 1 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM28    ITEM #:  0028
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.482 |
| SLOPE (SE) | = | 0.652 ( 0.054) |
| LOCATION (SE) | = | 0.652 ( 0.078) |
| LOW ASYM (SE) | = | 0.191 ( 0.026) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM29    ITEM #:  0029
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.645 |
| SLOPE (SE) | = | 1.010 ( 0.053) |
| LOCATION (SE) | = | 0.195 ( 0.053) |
| LOW ASYM (SE) | = | 0.189 ( 0.024) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM30    ITEM #:  0030
MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
|---|---|---|
| ITEM MEAN | = | 0.832 |
| SLOPE (SE) | = | 1.473 ( 0.067) |
| LOCATION (SE) | = | −0.122 ( 0.031) |
| LOW ASYM (SE) | = | 0.179 ( 0.018) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM31    ITEM #:  0031

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      0.605
SLOPE (SE)      =      1.193 ( 0.060)
LOCATION (SE)   =      0.007 ( 0.040)
LOW ASYM (SE)   =      0.198 ( 0.018)
D 1 (SE)        =      0.000 ( 0.000)
D 2 (SE)        =      0.000 ( 0.000)

THETA

ITEM ID: ITEM32    ITEM #:  0032

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      1.388
SLOPE (SE)      =      0.984 ( 0.016)
LOCATION (SE)   =      0.702 ( 0.016)
LOW ASYM (SE)   =      0.000 ( 0.000)
D 1 (SE)        =      1.073 ( 0.020)
D 2 (SE)        =      0.341 ( 0.020)
D 3 (SE)        =     -0.393 ( 0.022)
D 4 (SE)        =     -1.020 ( 0.029)
D 5 (SE)        =      0.000 ( 0.000)

PROPORTION

THETA

ITEM ID: ITEM32    ITEM #:  0032

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      1.388
SLOPE (SE)      =      0.984 ( 0.016)
LOCATION (SE)   =      0.702 ( 0.016)
LOW ASYM (SE)   =      0.000 ( 0.000)
D 1 (SE)        =      1.073 ( 0.020)
D 2 (SE)        =      0.341 ( 0.020)
D 3 (SE)        =     -0.393 ( 0.022)
D 4 (SE)        =     -1.020 ( 0.029)
D 5 (SE)        =      0.000 ( 0.000)

THETA

ITEM ID: ITEM32    ITEM #:  0032

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      1.388
SLOPE (SE)      =      0.984 ( 0.016)
LOCATION (SE)   =      0.702 ( 0.016)
LOW ASYM (SE)   =      0.000 ( 0.000)
D 1 (SE)        =      1.073 ( 0.020)
D 2 (SE)        =      0.341 ( 0.020)
D 3 (SE)        =     -0.393 ( 0.022)
D 4 (SE)        =     -1.020 ( 0.029)
D 5 (SE)        =      0.000 ( 0.000)

PROPORTION

THETA

ITEM ID: ITEM32    ITEM #:  0032

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      1.388
SLOPE (SE)      =      0.984 ( 0.016)
LOCATION (SE)   =      0.702 ( 0.016)
LOW ASYM (SE)   =      0.000 ( 0.000)
D 1 (SE)        =      1.073 ( 0.020)
D 2 (SE)        =      0.341 ( 0.020)
D 3 (SE)        =     -0.393 ( 0.022)
D 4 (SE)        =     -1.020 ( 0.029)
D 5 (SE)        =      0.000 ( 0.000)

THETA

ITEM ID: ITEM32    ITEM #:  0032

MCAS 2006 General Physics Grade 10

# EXAMINEES     =      5254
ITEM MEAN       =      1.388
SLOPE (SE)      =      0.984 ( 0.016)
LOCATION (SE)   =      0.702 ( 0.016)
LOW ASYM (SE)   =      0.000 ( 0.000)
D 1 (SE)        =      1.073 ( 0.020)
D 2 (SE)        =      0.341 ( 0.020)
D 3 (SE)        =     -0.393 ( 0.022)
D 4 (SE)        =     -1.020 ( 0.029)
D 5 (SE)        =      0.000 ( 0.000)

PROPORTION

THETA

47

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM33        ITEM #:  0033

MCAS 2006 General Physics Grade 10

| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.378 |      |
| SLOPE (SE)    | = | 0.852 ( 0.082) |
| LOCATION (SE) | = | 1.340 ( 0.057) |
| LOW ASYM (SE) | = | 0.232 ( 0.015) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM34        ITEM #:  0034

MCAS 2006 General Physics Grade 10

| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.553 |      |
| SLOPE (SE)    | = | 0.940 ( 0.050) |
| LOCATION (SE) | = | 0.109 ( 0.046) |
| LOW ASYM (SE) | = | 0.148 ( 0.020) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM35        ITEM #:  0035

MCAS 2006 General Physics Grade 10

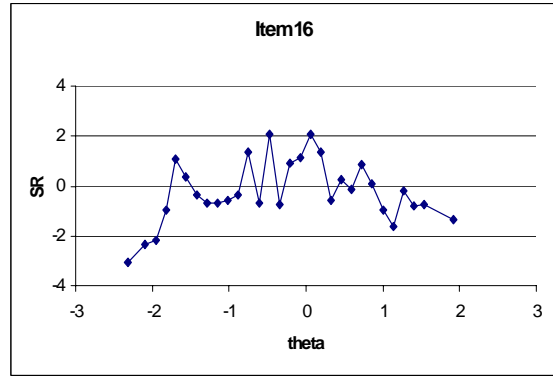| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.494 |      |
| SLOPE (SE)    | = | 1.240 ( 0.082) |
| LOCATION (SE) | = | 0.692 ( 0.036) |
| LOW ASYM (SE) | = | 0.276 ( 0.014) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM36        ITEM #:  0036

MCAS 2006 General Physics Grade 10

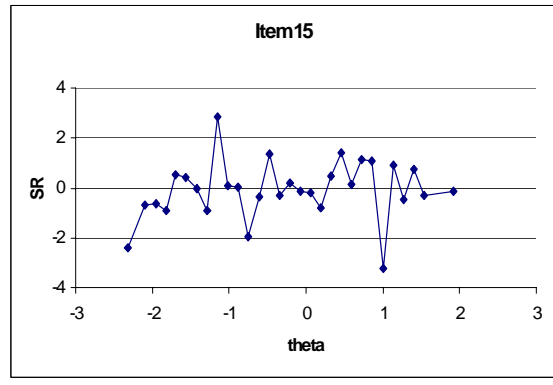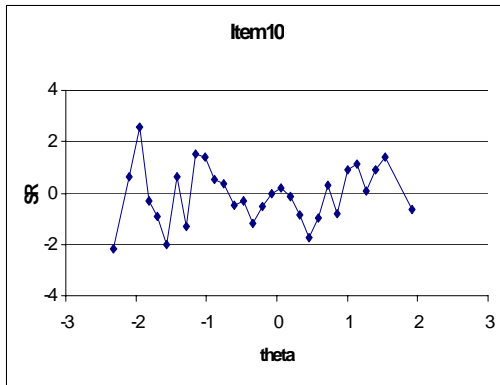| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.254 |      |
| SLOPE (SE)    | = | 1.745 ( 0.122) |
| LOCATION (SE) | = | 1.295 ( 0.026) |
| LOW ASYM (SE) | = | 0.149 ( 0.007) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

PROPORTION

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

ITEM ID: ITEM37        ITEM #:  0037

MCAS 2006 General Physics Grade 10

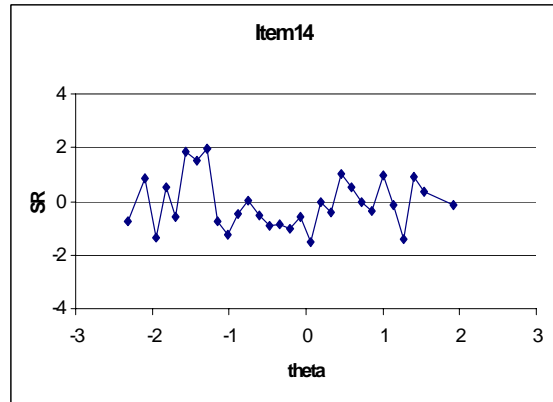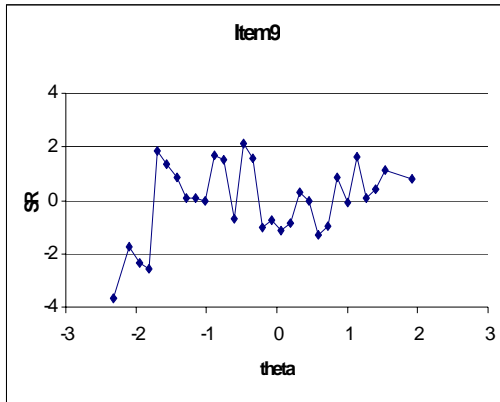| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.268 |      |
| SLOPE (SE)    | = | 0.803 ( 0.067) |
| LOCATION (SE) | = | 1.815 ( 0.090) |
| LOW ASYM (SE) | = | 0.117 ( 0.017) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

THETA

LOGISTIC G.R. MODEL WITH ITEM AND CATEGORY PARAMETERS SEPARATED

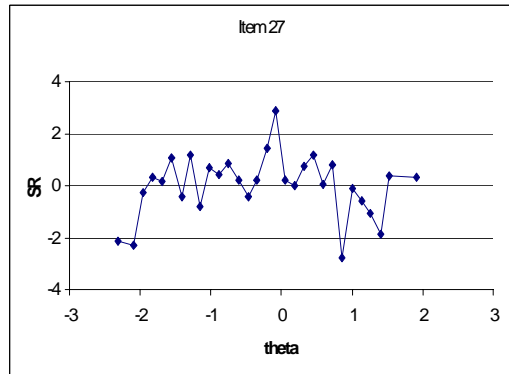ITEM ID: ITEM38        ITEM #:  0038

MCAS 2006 General Physics Grade 10

| # EXAMINEES   | = | 5254 |       |
| ITEM MEAN     | = | 0.479 |      |
| SLOPE (SE)    | = | 1.103 ( 0.059) |
| LOCATION (SE) | = | 0.419 ( 0.035) |
| LOW ASYM (SE) | = | 0.154 ( 0.015) |
| D 1 (SE)      | = | 0.000 ( 0.000) |
| D 2 (SE)      | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM39     ITEM #: 0039

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               1.064
SLOPE (SE)      =      1.141  ( 0.018)
LOCATION (SE)   =      1.062  ( 0.015)
LOW ASYM (SE)   =      0.000  ( 0.000)
$D_1$ (SE)              1.676  ( 0.019)
$D_2$ (SE)              0.089  ( 0.020)
$D_3$ (SE)             -0.523  ( 0.025)
$D_4$ (SE)             -1.242  ( 0.039)
$D_5$ (SE)              0.000  ( 0.000)

THETA

---

ITEM ID: ITEM39     ITEM #: 0039

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               1.064
SLOPE (SE)      =      1.141  ( 0.018)
LOCATION (SE)   =      1.062  ( 0.015)
LOW ASYM (SE)   =      0.000  ( 0.000)
$D_1$ (SE)              1.676  ( 0.019)
$D_2$ (SE)              0.089  ( 0.020)
$D_3$ (SE)             -0.523  ( 0.025)
$D_4$ (SE)             -1.242  ( 0.039)
$D_5$ (SE)              0.000  ( 0.000)

PROPORTION

THETA

---

ITEM ID: ITEM39     ITEM #: 0039

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               1.064
SLOPE (SE)             1.141  ( 0.018)
LOCATION (SE)         1.062  ( 0.015)
LOW ASYM (SE)         0.000  ( 0.000)
$D_1$ (SE)              1.676  ( 0.019)
$D_2$ (SE)              0.089  ( 0.020)
$D_3$ (SE)             -0.523  ( 0.025)
$D_4$ (SE)             -1.242  ( 0.039)
$D_5$ (SE)              0.000  ( 0.000)

THETA

---

ITEM ID: ITEM39     ITEM #: 0039

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               1.064
SLOPE (SE)             1.141  ( 0.018)
LOCATION (SE)         1.062  ( 0.015)
LOW ASYM (SE)         0.000  ( 0.000)
$D_1$ (SE)              1.676  ( 0.019)
$D_2$ (SE)              0.089  ( 0.020)
$D_3$ (SE)             -0.523  ( 0.025)
$D_4$ (SE)             -1.242  ( 0.039)
$D_5$ (SE)              0.000  ( 0.000)

PROPORTION

THETA

---

ITEM ID: ITEM39     ITEM #: 0039

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               1.064
SLOPE (SE)             1.141  ( 0.018)
LOCATION (SE)         1.062  ( 0.015)
LOW ASYM (SE)         0.000  ( 0.000)
$D_1$ (SE)              1.676  ( 0.019)
$D_2$ (SE)              0.089  ( 0.020)
$D_3$ (SE)             -0.523  ( 0.025)
$D_4$ (SE)             -1.242  ( 0.039)
$D_5$ (SE)              0.000  ( 0.000)

THETA

---

ITEM ID: ITEM40     ITEM #: 0040

MCAS 2006 General Physics Grade 10

# EXAMINEES      =      5254
ITEM MEAN               0.424
SLOPE (SE)             0.464  ( 0.043)
LOCATION (SE)         0.857  ( 0.098)
LOW ASYM (SE)         0.109  ( 0.050)
$D_1$ (SE)              0.000  ( 0.000)
$D_2$ (SE)              0.000  ( 0.000)

PROPORTION

THETA

ITEM ID: ITEM41      ITEM #:  0041

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.427 |
| SLOPE (SE) | = | 0.690 ( 0.050) |
| LOCATION (SE) | = | 0.698 ( 0.059) |
| LOW ASYM (SE) | = | 0.125 ( 0.021) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM42      ITEM #:  0042

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.528 |
| SLOPE (SE) | = | 0.890 ( 0.047) |
| LOCATION (SE) | = | 0.154 ( 0.047) |
| LOW ASYM (SE) | = | 0.121 ( 0.020) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

ITEM ID: ITEM43      ITEM #:  0043

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.623 |
| SLOPE (SE) | = | 1.291 ( 0.062) |
| LOCATION (SE) | = | −0.077 ( 0.036) |
| LOW ASYM (SE) | = | 0.189 ( 0.018) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

ITEM ID: ITEM44      ITEM #:  0044

MCAS 2006 General Physics Grade 10

| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.264 |
| SLOPE (SE) | = | 1.090 ( 0.092) |
| LOCATION (SE) | = | 1.508 ( 0.048) |
| LOW ASYM (SE) | = | 0.154 ( 0.010) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

PROPORTION

THETA

MCAS 2006 General Physics Grade 10

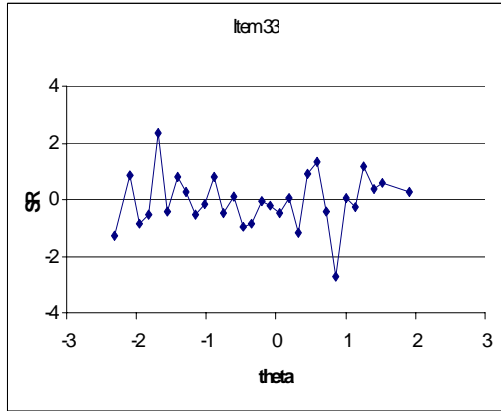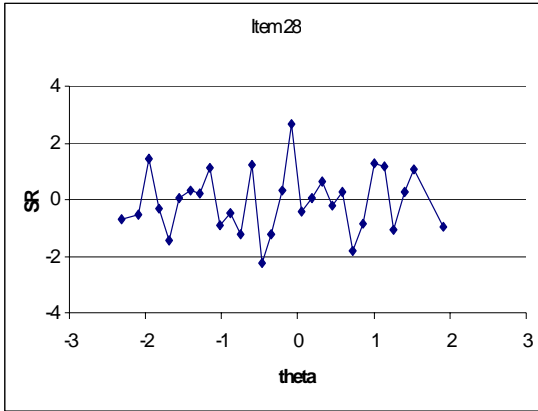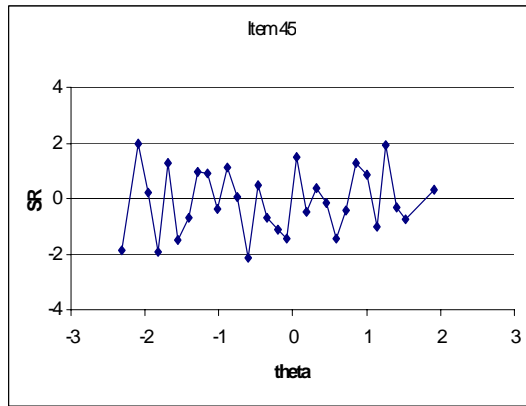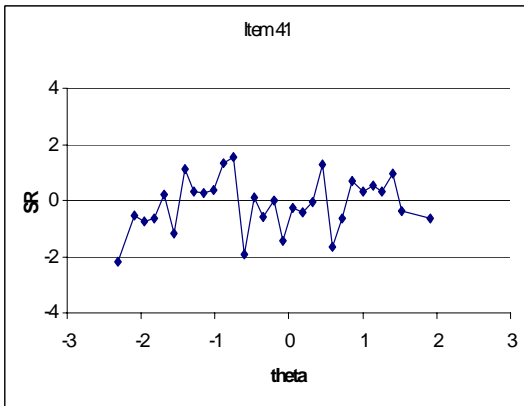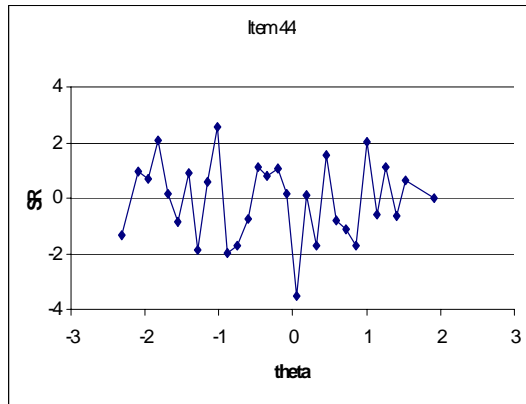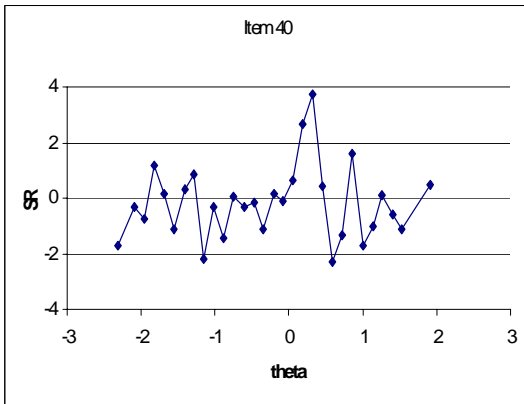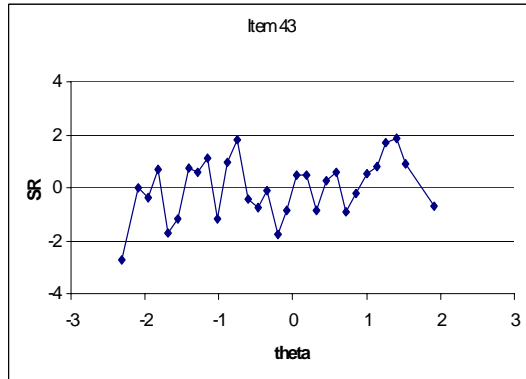| # EXAMINEES | = | 5254 |
| ITEM MEAN | = | 0.510 |
| SLOPE (SE) | = | 1.144 ( 0.063) |
| LOCATION (SE) | = | 0.364 ( 0.036) |
| LOW ASYM (SE) | = | 0.192 ( 0.016) |
| D 1 (SE) | = | 0.000 ( 0.000) |
| D 2 (SE) | = | 0.000 ( 0.000) |

THETA

50

# Appendix C.   Standardized Residual (SR) Plots

## Part I:  Standardized Residual (SR) Plots for Dichotomous Items



Item1



Item5



Item2



Item6



Item3



Item7



Item4



Item8

**Item9**

**Item14**

**Item10**

**Item15**

**Item12**

**Item16**

**Item13**

Item 17

Item 18

Item 22

Item 19

Item 23

Item 20

Item 24

Item 21

Item 27

Item 28

Item 33

Item 29

Item 34

Item 30

Item 35

Item 31

Item 36

Item 37

Item 42

Item 38

Item 43

Item 40

Item 44

Item 41

Item 45

# Part II: Standardized Residual (SR) Plots for Polytomous Items

Item 39(1)

Item 39 (2)

Item 39(3)

Item 39 (4)

Item 39(5)