

## Section 7. Field Test Analyses

Following the receipt of the final score file from Measurement Incorporated and Pearson for the January administration and from Pearson for the May administration, the field test analyses for SR and SPR items from January and May were completed. The analyses consisted of four components: classical item analyses, differential item functioning (DIF), calibration, and scaling. All the analyses were completed using *GENASYS*, an ETS proprietary software program. The analysis procedures for each component are described in detail below. All valid records available were used as samples for the analyses, including those for students learning English as a second language, students with IEP or 504 plans, and students receiving accommodations. Only records invalidated by the test administrator and records with no item responses to the first five items were excluded from the analysis sample.

### Classical Item Analyses

Classical item analyses involve computing a set of statistics based on classical test theory for every item in each form. The statistics provide key information about the quality of the items from an empirical perspective. The statistics estimated for the HSA field test items, and associated criteria used to flag items for the content specialists' review, are described below.

Classical item difficulty (“p-value”): This statistic indicates the mean item score expressed as a proportion of the maximum obtainable item score. For SR and SPR items, it is equivalent to the proportion of examinees in the sample that answered the item correctly. Desired p-values generally fall within the range of 0.25 to 0.90. Occasionally, items that fall outside this range can be justified for inclusion in an item bank based upon the quality and educational importance of the item content or the ability to measure students with very high or low achievement, especially if the students have not yet received instruction in the content.

Item-total correlation of the correct response option for SR and SPR items: This statistic describes the relationship between performance on the specific item and performance on the total test, including the item under study. It is sometimes referred to as a discrimination index. For SR and SPR items, the item-total correlation is the point-biserial correlation. Values less than 0.15 were flagged for a weaker than desired relationship and receive careful consideration by ETS staff and MSDE before including them on future forms. Items with negative correlations can indicate there are serious problems with the item content (e.g., multiple correct answers, unusually complex content), there is an incorrect key, or students have not been taught the content.

Proportion of students choosing each response option (SR items): This statistic indicates the percent of examinees selecting each answer choice, or option. Options not selected by any students or selected by a very low proportion of students indicate problems with plausibility of the option. Items that do not have all answer options functioning may be discarded or revised and field tested again.

Point-biserial correlation of incorrect response option (SR items) with the total raw score: These statistics describe the relationship between selecting an incorrect response option for a specific item and performance on the total test, including the item under study. Typically, the correlation between an incorrect answer and total test performance is weak or negative. Values are typically compared and contrasted with the discrimination index. When the magnitude of these point-biserial correlations for the incorrect answer is stronger relative to the correct answer, the item will be carefully reviewed for content-related problems. Alternatively, positive point-biserial correlations on incorrect options may indicate that students have not had sufficient opportunity to learn the material.

Percent of students omitting an item: This statistic is useful for identifying problems with test features, such as testing time and item/test layout. Typically, it is assumed that if students have an adequate amount of testing time, 95 percent should attempt to answer each question. When a pattern of omit percentages exceeds 5 percent for a series of items at the end of a timed section, this may indicate that there was insufficient time for students to complete all items. For individual items, if the omit percentage is greater than 5 percent for a single SR or SPR item, this could be an indication of an item/test layout problem. For example, students might accidentally skip an item that follows a lengthy stem.

In addition, a series of flags was created to identify items with extreme values. Flagged items were subject to additional scrutiny prior to the inclusion of the items in the final calibrations. The following flagging criteria were applied to all items tested in the 2009 assessments:

- *Difficulty flag:* P-values less than 0.25 or greater than 0.90.
- *Discrimination flag:* Item-total correlation less than 0.15.
- *Distractor flag:* SR point-biserial correlation positive for incorrect option.
- *Omit flag:* Percent omitted is greater than 5 for SR and SPR items.

Distributions of p-values and item-total correlations for the field test items administered in January 2009 are shown in Tables 7.1 and 7.2, respectively. Corresponding results for the field test items administered in May 2009 are shown in Tables 7.3 and 7.4, respectively.

Following the classical item analyses, items with poor item statistics and items that were not scored as per MSDE's instructions were removed from further analyses (see Table 7.5). These items have been identified for revision and possible additional field testing. Table 7.6 presents the number of items that were retained for further analyses and evaluation after being flagged for statistical reasons, including extreme p-values, low item-total correlations, and/or high omits rates. Calibration results indicated the items were estimated reasonably, and therefore they were not removed from scaling.

**Table 7.1** Distribution of P-Values for the MD HSA January 2009 Field Test Items

| P-Value                | Percentage and Number of Items |   |         |    |         |   |            |    |
|------------------------|--------------------------------|---|---------|----|---------|---|------------|----|
|                        | Algebra <sup>a</sup>           |   | Biology |    | English |   | Government |    |
|                        | %                              | N | %       | N  | %       | N | %          | N  |
| $P < 0.25$             | 20                             | 3 | 3       | 1  | 5       | 1 | 0          | 0  |
| $0.25 \leq P < 0.35$   | 27                             | 4 | 3       | 1  | 14      | 3 | 10         | 3  |
| $0.35 \leq P < 0.45$   | 27                             | 4 | 41      | 12 | 18      | 4 | 28         | 8  |
| $0.45 \leq P < 0.55$   | 13                             | 2 | 21      | 6  | 18      | 4 | 24         | 7  |
| $0.55 \leq P < 0.65$   | 7                              | 1 | 14      | 4  | 27      | 6 | 34         | 10 |
| $0.65 \leq P < 0.75$   | 0                              | 0 | 17      | 5  | 14      | 3 | 3          | 1  |
| $0.75 \leq P < 0.85$   | 7                              | 1 | 0       | 0  | 0       | 0 | 0          | 0  |
| $P \geq 0.85$          | 0                              | 0 | 0       | 0  | 5       | 1 | 0          | 0  |
| Descriptive Statistics |                                |   |         |    |         |   |            |    |
| N Items                | 15                             |   | 29      |    | 22      |   | 29         |    |
| Mean                   | 0.36                           |   | 0.48    |    | 0.51    |   | 0.49       |    |
| SD                     | 0.16                           |   | 0.14    |    | 0.17    |   | 0.11       |    |
| Min                    | 0.08                           |   | 0.22    |    | 0.22    |   | 0.26       |    |
| Max                    | 0.75                           |   | 0.72    |    | 0.88    |   | 0.67       |    |

<sup>a</sup>SPR items included**Table 7.2** Distribution of Item-Total Correlations for the MD HSA January 2009 Field Test Items

| Correlation            | Percentage and Number of Items |   |         |    |         |    |            |   |
|------------------------|--------------------------------|---|---------|----|---------|----|------------|---|
|                        | Algebra <sup>a</sup>           |   | Biology |    | English |    | Government |   |
|                        | %                              | N | %       | N  | %       | N  | %          | N |
| $R < 0.15$             | 20                             | 3 | 0       | 0  | 5       | 1  | 0          | 0 |
| $0.15 \leq R < 0.25$   | 7                              | 1 | 7       | 2  | 18      | 4  | 17         | 5 |
| $0.25 \leq R < 0.35$   | 40                             | 6 | 17      | 5  | 55      | 12 | 31         | 9 |
| $0.35 \leq R < 0.45$   | 33                             | 5 | 45      | 13 | 18      | 4  | 24         | 7 |
| $0.45 \leq R < 0.55$   | 0                              | 0 | 31      | 9  | 5       | 1  | 24         | 7 |
| $0.55 \leq R < 0.65$   | 0                              | 0 | 0       | 0  | 0       | 0  | 3          | 1 |
| $0.65 \leq R < 0.75$   | 0                              | 0 | 0       | 0  | 0       | 0  | 0.0        | 0 |
| $R \geq 0.75$          | 0                              | 0 | 0       | 0  | 0       | 0  | 0.0        | 0 |
| Descriptive Statistics |                                |   |         |    |         |    |            |   |
| N Items                | 15                             |   | 29      |    | 22      |    | 29         |   |
| Mean                   | 0.28                           |   | 0.40    |    | 0.30    |    | 0.37       |   |
| SD                     | 0.11                           |   | 0.09    |    | 0.08    |    | 0.11       |   |
| Min                    | 0.10                           |   | 0.18    |    | 0.13    |    | 0.17       |   |
| Max                    | 0.42                           |   | 0.52    |    | 0.45    |    | 0.56       |   |

<sup>a</sup>SPR items included

**Table 7.3** Distribution of P-Values for the MD HSA May 2009 Field Test Items

| P-Value                | Percentage and Number of Items |    |         |    |         |    |            |    |
|------------------------|--------------------------------|----|---------|----|---------|----|------------|----|
|                        | Algebra <sup>a</sup>           |    | Biology |    | English |    | Government |    |
|                        | %                              | N  | %       | N  | %       | N  | %          | N  |
| $P < 0.25$             | 4                              | 7  | 1       | 3  | 1       | 3  | 3          | 7  |
| $0.25 \leq P < 0.35$   | 9                              | 15 | 7       | 16 | 5       | 15 | 6          | 17 |
| $0.35 \leq P < 0.45$   | 22                             | 38 | 11      | 28 | 10      | 32 | 11         | 30 |
| $0.45 \leq P < 0.55$   | 19                             | 33 | 20      | 48 | 14      | 47 | 16         | 42 |
| $0.55 \leq P < 0.65$   | 22                             | 38 | 26      | 64 | 19      | 62 | 19         | 51 |
| $0.65 \leq P < 0.75$   | 14                             | 25 | 19      | 47 | 24      | 78 | 23         | 60 |
| $0.75 \leq P < 0.85$   | 8                              | 14 | 9       | 22 | 21      | 69 | 13         | 34 |
| $P \geq 0.85^b$        | 2                              | 4  | 7       | 18 | 7       | 22 | 9          | 25 |
| Descriptive Statistics |                                |    |         |    |         |    |            |    |
| N Items                | 174                            |    | 246     |    | 328     |    | 266        |    |
| Mean                   | 0.53                           |    | 0.59    |    | 0.63    |    | 0.60       |    |
| SD                     | 0.16                           |    | 0.16    |    | 0.16    |    | 0.18       |    |
| Min                    | 0.18                           |    | 0.21    |    | 0.16    |    | 0.11       |    |
| Max                    | 0.90                           |    | 0.95    |    | 0.95    |    | 0.94       |    |

<sup>a</sup> SPR items included; <sup>b</sup> P-value > 0.90: 7 Biology , 4 English, and 3 Government items

**Table 7.4** Distribution of Item-Total Correlations for the MD HSA May 2009 Field Test Items

| Correlation            | Percentage and Number of Items |    |         |    |         |     |            |     |
|------------------------|--------------------------------|----|---------|----|---------|-----|------------|-----|
|                        | Algebra <sup>a</sup>           |    | Biology |    | English |     | Government |     |
|                        | %                              | N  | %       | N  | %       | N   | %          | N   |
| $R < 0.15$             | 2                              | 3  | 2       | 5  | 4       | 12  | 5          | 12  |
| $0.15 \leq R < 0.25$   | 5                              | 9  | 10      | 25 | 11      | 37  | 9          | 25  |
| $0.25 \leq R < 0.35$   | 15                             | 26 | 26      | 63 | 29      | 94  | 18         | 47  |
| $0.35 \leq R < 0.45$   | 33                             | 57 | 38      | 93 | 42      | 138 | 43         | 114 |
| $0.45 \leq R < 0.55$   | 30                             | 53 | 24      | 58 | 14      | 47  | 25         | 66  |
| $0.55 \leq R < 0.65$   | 14                             | 24 | 1       | 2  | 0       | 0   | 1          | 2   |
| $0.65 \leq R < 0.75$   | 1                              | 2  | 0       | 0  | 0       | 0   | 0          | 0   |
| $R \geq 0.75$          | 0                              | 0  | 0       | 0  | 0       | 0   | 0          | 0   |
| Descriptive Statistics |                                |    |         |    |         |     |            |     |
| N Items                | 174                            |    | 246     |    | 328     |     | 266        |     |
| Mean                   | 0.43                           |    | 0.37    |    | 0.35    |     | 0.38       |     |
| SD                     | 0.12                           |    | 0.10    |    | 0.10    |     | 0.11       |     |
| Min                    | 0.10                           |    | 0.08    |    | 0.08    |     | 0.08       |     |
| Max                    | 0.67                           |    | 0.58    |    | 0.54    |     | 0.59       |     |

<sup>a</sup> SPR items included

**Table 7.5 MD HSA 2009 Field Test Items Excluded from Calibration**

| Administration | Content    | ItemID | Form | Sequence | Response Type | Reason   |
|----------------|------------|--------|------|----------|---------------|--|
| January        | Biology    | 79501  | A, B | 18       | SR            | R_ITT = -0.02  |
|                | English    | 108772 | A    | 55       | SR            | Faulty item; MSDE approved item be suppressed in IRT |
|                | English    | 251242 | A, B | 22       | SR            | Faulty item; MSDE approved item be suppressed in IRT |
|                | English    | 251243 | A, B | 23       | SR            | Faulty item; MSDE approved item be suppressed in IRT |
|                | English    | 251244 | A, B | 24       | SR            | Faulty item; MSDE approved item be suppressed in IRT |
|                | Government | 302865 | A, B | 63       | SR            | R_ITT = 0.04   |
| May            | Algebra    | 282463 | H    | 63       | SR            | R_ITT = 0.01   |
|                | Algebra    | 268716 | J    | 15       | SR            | R_ITT = 0.05   |
|                | Biology    | 297528 | D    | 89       | SR            | R_ITT = -0.19  |
|                | Biology    | 271125 | E    | 44       | SR            | R_ITT = -0.14  |
|                | Biology    | 263127 | K    | 42       | SR            | R_ITT = -0.03  |
|                | Biology    | 256519 | M    | 74       | SR            | R_ITT = 0.04   |
|                | Biology    | 264041 | N    | 73       | SR            | R_ITT = -0.05  |
|                | English    | 288639 | C    | 21       | SR            | R_ITT = 0.05   |
|                | English    | 281409 | E    | 40       | SR            | R_ITT = 0.06   |
|                | English    | 261667 | F    | 88       | SR            | R_ITT = -0.04  |
|                | English    | 281757 | H    | 6        | SR            | R_ITT = 0.04   |
|                | English    | 264668 | H    | 18       | SR            | R_ITT = -0.15  |
|                | English    | 264669 | H    | 19       | SR            | R_ITT = 0.06   |
|                | English    | 281386 | H    | 84       | SR            | R_ITT = 0.04   |
|                | English    | 285440 | J    | 88       | SR            | R_ITT = -0.00  |
|                | English    | 288647 | K    | 40       | SR            | R_ITT = -0.05  |
|                | English    | 288672 | L    | 22       | SR            | R_ITT = 0.01   |
|                | English    | 285495 | M    | 89       | SR            | R_ITT = 0.07   |
|                | English    | 281404 | N    | 22       | SR            | R_ITT = 0.05   |
|                | English    | 285617 | N    | 39       | SR            | R_ITT = -0.10  |
|                | Government | 297121 | C    | 104      | SR            | R_ITT = -0.08  |
|                | Government | 296522 | D    | 69       | SR            | R_ITT = 0.06   |
|                | Government | 79700  | D    | 80       | SR            | R_ITT = 0.06   |
| Government     | 263975     | G      | 60   | SR       | R_ITT = -0.01 |  |
| Government     | 283278     | H      | 104  | SR       | R_ITT = 0.04  |  |
| Government     | 297436     | J      | 52   | SR       | R_ITT = 0.03  |  |
| Government     | 296486     | M      | 52   | SR       | R_ITT = 0.02  |  |
| Government     | 263974     | M      | 68   | SR       | R_ITT = -0.07 |  |
| Government     | 279834     | N      | 6    | SR       | R_ITT = 0.05  |  |

**Table 7.6 MD HSA 2009 Field Test Items with Statistical Flags Retained in Calibration**

|                | P-<br>Value<br>< 0.25 | P-<br>Value<br>> 0.90 | R_ITT<br>< 0.15 | Distractor<br>Pt-Bis<br>> 0 | Omit<br>Rate<br>> 5% | C-<br>Level<br>DIF | Missing<br>Response <sup>a</sup> | Total<br>Flags | N Items <sup>b</sup> |
|----------------|-----------------------|-----------------------|-----------------|-----------------------------|----------------------|--------------------|----------------------------------|----------------|----------------------|
| <b>January</b> |                       |                       |                 |                             |                      |                    |                                  |                |                      |
| Algebra        | 3                     | 0                     | 3               | 2                           | 6                    | 1                  | 0                                | 15             | 10                   |
| Biology        | 1                     | 0                     | 0               | 1                           | 0                    | 1                  | 0                                | 3              | 3                    |
| English        | 1                     | 0                     | 1               | 1                           | 0                    | 2                  | 0                                | 5              | 4                    |
| Government     | 0                     | 0                     | 0               | 3                           | 0                    | 0                  | 0                                | 3              | 3                    |
| <b>May</b>     |                       |                       |                 |                             |                      |                    |                                  |                |                      |
| Algebra        | 7                     | 0                     | 3               | 14                          | 14                   | 5                  | 0                                | 43             | 33                   |
| Biology        | 3                     | 7                     | 5               | 22                          | 0                    | 3                  | 0                                | 40             | 33                   |
| English        | 3                     | 4                     | 12              | 42                          | 0                    | 13                 | 0                                | 74             | 60                   |
| Government     | 7                     | 3                     | 12              | 32                          | 0                    | 12                 | 0                                | 66             | 47                   |

<sup>a</sup> SR option with 0 students; <sup>b</sup> Represents total number of unique items.

### Differential Item Functioning

Following the classical item analyses, differential item functioning analyses were completed. One goal of test development is to assemble a set of items that provides an estimate of student ability that is as fair and accurate as possible for all groups within the population. DIF statistics are used to identify items whereby identifiable groups of students with the same underlying level of ability (e.g., females, African Americans, Hispanics) have different probabilities of answering correctly. If the item is more difficult for an identifiable subgroup, the item may be measuring something different from the intended construct. However, it is important to recognize that DIF-flagged items might be related to actual differences in relevant knowledge or skill (item impact) or statistical Type I error. A subsequent review by MSDE and ETS content experts is conducted to investigate the source and meaning of evident differences.

ETS used the Mantel-Haenszel DIF detection method. As part of the Mantel-Haenszel procedure, the statistic described by Holland & Thayer (1988), known as MH D-DIF, was used<sup>8</sup>. This statistic is expressed as the difference between the focal and reference group performance on an

<sup>8</sup> The formula for the estimate of constant odds ratio is

$$\hat{\alpha}_{MH} = \frac{\left( \frac{\sum_m R_{fm} W_{rm}}{N_m} \right)}{\left( \frac{\sum_m R_{rm} W_{fm}}{N_m} \right)},$$

where

- $R_{rm}$  = number in reference group at ability level m answering the item right,
- $W_{fm}$  = number in focal group at ability level m answering the item wrong,
- $R_{fm}$  = number in focal group at ability level m answering the item right,
- $W_{rm}$  = number in reference group at ability level m answering the item wrong,
- $N_m$  = total group at ability level m.

This can then be used in the following formula (Holland & Thayer, 1985):

$$MH\ D - DIF = -2.35 \ln[\alpha_{MH}].$$

item after conditioning on total test score. Negative MH D-DIF statistics favor the reference group, and positive values favor the focal group. The classification logic used for flagging items is based on a combination of absolute differences and significance testing. Items that are not significantly different based on the MH D-DIF ( $p > 0.05$ ) are considered to have similar performance between the two studied groups and to be functioning appropriately. For items where the statistical test indicates significant differences ( $p < 0.05$ ), the effect size is used to determine the direction and severity of the DIF. The male and white groups were treated as the reference groups for gender and ethnicity, respectively; the female and other race and ethnic groups were considered the focal groups.

Based on their DIF statistics, items are classified into one of three categories and assigned values of A, B, or C. Category A items contain negligible DIF, Category B items exhibit slight or moderate DIF, and Category C items have moderate to large DIF. Negative values imply that, conditional on the matching variable, the focal group has a lower mean item score than the reference group. In contrast, a positive value implies that, conditional on the matching variable, the reference group has a lower mean item score than the focal group.

Among the items field-tested in January, one Algebra item, one Biology item, and two English items were flagged for C-level DIF. Among the items field tested in May, five Algebra items, three Biology items, thirteen English items, and twelve Government items were flagged for C-level DIF. These flags were recorded in the item bank. The flagged items will be reviewed by ETS and MSDE content specialists as well as by ETS senior staff to determine their availability for future use.

### **IRT Calibration and Scaling**

One purpose of item calibration and scaling is to create a common scale for expressing the difficulty estimates of all the items across all versions of a test. The resulting scale has a mean score of 0 and a standard deviation of 1. This scale is often referred to as the “theta” metric and is not used for reporting purposes because the values typically range from  $-3$  to  $+3$ . Therefore, the scale is usually transformed to a reporting scale (also known as a scale score), which can be more meaningfully interpreted by students, teachers, and other stakeholders.

As noted previously, the IRT model used to calibrate the MD HSA test items was the 3-parameter logistic (3PL) model. Item response theory expresses the probability that a student will achieve a certain score on an item (such as correct or incorrect) as a function of the item’s statistical properties and the ability level (or proficiency level) of the student.

The 3PL model relates the probability that a person with ability  $\theta$  will respond correctly to item  $i$  as follows:

$$P_i(\theta) = c_i + \frac{1 - c_i}{1 + e^{-1.7a_i(\theta - b_i)}},$$

where

- $a_i$  is the slope parameter of item  $i$ , characterizing its discrimination;
- $b_i$  is the location parameter of item  $i$ , characterizing its difficulty; and
- $c_i$  is the lower asymptote parameter of item  $i$ , reflecting the chance that students with very low proficiency will select the correct answer, sometimes called the “pseudo-guessing” level.

A proprietary version of the *PARSCALE* computer program (Muraki & Bock, 1995) was used for all item calibration work. The resulting calibrations were then scaled to the bank estimates using Stocking and Lord’s (1983) test characteristic curve (TCC) method and the operational items as the anchor set.

The calibration and equating process is outlined in the steps below.

1. For each test, calibrate all items using a sparse matrix design that places all items on a common scale. Essentially, this means that the data were set up using the following format. In the diagram below, X’s represent items and spaces indicate missing data. For example, items included on version 2 but not on version 1, 3, 4, or 5 were treated as “not reached” for the purposes of the analyses and are denoted as “missing” in the diagram below.

| Common   | Unique 1 | Unique 2 | Unique 3 | Unique 4 | Unique 5 |
|----------|----------|----------|----------|----------|----------|
| XXXXXXXX | XXXXXXXX |          |          |          |          |
| XXXXXXXX |          | XXXXXXXX |          |          |          |
| XXXXXXXX |          |          | XXXXXXXX |          |          |
| XXXXXXXX |          |          |          | XXXXXXXX |          |
| XXXXXXXX |          |          |          |          | XXXXXXXX |

2. Once the items have been calibrated, results are reviewed to determine if any items failed to calibrate.
3. After the final calibration parameters were obtained, the items were then linked to the bank scale using the TCC method. Specifically, the banked parameters of the primary form operational items were used to place the field test items onto the operational reporting scale.

Once the items were calibrated and placed onto the operational scale, they were loaded into the item bank. Items that were not calibrated were listed as unavailable (see Table 7.5).