

New Jersey

Start Strong Assessment–Science

(NJSSA–S)

TECHNICAL BRIEF
Grades 6, 9, and 12

2022

JULY 2023

PTM XXXX.XX



State of New Jersey
Department of Education

Copyright © 2023 by New Jersey Department of Education
All rights reserved

Contents

Contents	2
Part 1: Description of the NJSSA–S	5
1.1 Content Domains and Scientific Practices	5
1.2 Crosscutting Concepts.....	10
Part 2: Test Development and Test Administration	11
2.1 Test Specifications	11
2.1.1 Test Blueprints	11
2.1.2 Item Types	12
2.2 Item Development Processes.....	13
2.3 Test Construction Process.....	14
2.3.1 Test Construction — First Draft	14
2.3.2 Test Construction Content Review	16
2.3.3 Test Construction NJDOE Review.....	16
2.4 2020 NJSSA–S Test Construction	16
2.4.1 Grade 6 Test Construction	17
2.4.2 Grade 9 Test Construction	19
2.4.3 Grade 12 Test Construction	21
2.5 Test Administration	23
2.6 Test Registration.....	24
2.7 Test Accessibility Features and Accommodations	24
2.7.1 Accessibility Features	24
2.7.2 Accommodations	25
2.8 Scores and Score Reports.....	27
2.8.1 Machine-Scored Items	27
2.8.2 Adjudication.....	27
2.8.3 Test Scores.....	28
2.8.4 Support Level	28
2.8.5 Student-Level Reports	29
2.8.6 Classroom-, School-, and District-Level Reports.....	29
Part 3: Item and Test Statistics	31
3.1 Classical Test Theory Statistics	31

3.1.1 Item Difficulty and Discrimination Descriptive Statistics	31
3.1.2 Speededness	42
3.1.3 Operational DIF Analysis	42
3.2 Item Response Theory	47
3.2.1 Unidimensionality	48
3.2.2 Partial-Credit Model-Fit Statistics	51
3.2.3 Local Independence	62
3.2.4 Descriptive Statistics — Raw Score	63
Part 4: Scale Stability	64
4.1 Delta Plot Method	64
4.2 0.3 Logits Absolute Difference Method	67
Part 5: Reliability	70
5.1 Classical Test Theory Reliability Estimates.....	70
5.1.1 Reliability and Measurement Error.....	70
5.1.2 Raw Score Internal Consistency.....	71
5.2 Item Response Theory Reliability	75
5.2.1 Test Information Functions	75
5.2.2 Item Maps.....	78
5.3 Reliability of Performance Classifications	81
5.3.1 Conditional Standard Error of Measurement at Each Cut-Score	82
5.3.2 Classification Consistency Indices.....	83
Part 6: Validity	84
6.1 Evidence Based on Test Content	85
6.2 Evidence Based on Response Processes	86
6.3 Evidence Based on Internal Structure	87
6.3.1 Intercorrelations	87
6.3.2 Other Internal Structure Evidence	88
6.4 Evidence Based on Relationships to Other Variables.....	88
6.5 Evidence Based on the Consequences of Testing	88
6.6 Other Validity Evidence.....	89
6.7 Summary of Validity Evidence.....	90
Appendix A: Detailed Test Maps	92

Appendix B: Raw Score Cumulative Frequency Distributions.....	98
Appendix C: Item Parameter Estimates and Model Fit Tables.....	107
Appendix D: Scale Stability Results Tables	110
Appendix E: Raw-to-Theta Score Tables	116
Appendix F: Conditional Standard Error of Measurement and Test Characteristic Curve Graphs	119
References.....	124

Part 1: Description of the NJSSA–S

The New Jersey Start Strong Assessment–Science (NJSSA–S) assesses students at the beginning of grades 6, 9, and 12 on their understanding and explanations of scientific phenomena and scenarios. The tests cover a range of material based upon the National Research Council’s (NRC; 2012) *Framework for K–12 Science Education*, which identifies the science knowledge and skills that all K–12 students should know, and the Next Generation Science Standards (NGSS Lead States, 2013), developed collaboratively by stakeholders across 25 states. To accomplish the necessary scope, each test item requires students to address multiple underlying variables, with items representing an interaction of Disciplinary Core Ideas (DCIs — within the domains of Physical, Life, and Earth and Space Science), Scientific and Engineering Practices (SEPs — Investigating, Sensemaking or Critiquing), and Crosscutting Concepts (CCC). Every test item counts toward the students’ performance in exactly one reported domain and one reported practice. Additionally, each item is also aligned to a CCC and the CCC concepts. The knowledge, skills, and abilities associated with the CCC and CCC concepts contribute to the overall score; however, the NJSSA–S does not report a specific CCC performance indicator. All items are machine-scored and consist of a mixture of multiple-choice (MC) and technology-enhanced (TE) items. Section 1.1 details the content domains and scientific practices, as well as how they are grouped to form the NJSSA–S reporting categories; next, Section 1.2 briefly describes the Crosscutting Concepts.

1.1 Content Domains and Scientific Practices

Although the NJSSA–S is a unidimensional test, six distinct, foundational sub-categories represent the three science content domains (Earth and Space, Life, and Physical) and the three scientific and engineering practices (Sensemaking, Critiquing, and Investigating).

Science content domains. Disciplinary Core Ideas can be classified into three major science content domains: Earth and Space Science, Life Science, and Physical Science. The NJSSA–S is designed to measure student performance in each of the three science content domains. The test development processes focus on balancing each science content domain equally. Furthermore, within each content domain, each DCI is balanced.

1. Earth and Space Science. The *Framework* states that “Earth and space sciences (ESS) investigate processes that operate on Earth and also address its place in the solar system” (National Research Council, 2012, p. 169). Table 1.1.1 shows the three ESS DCIs as well as the topics that are delineated within each.

Table 1.1.1: Earth and Space Science DCIs

DCI Topic Description	
ESS1: Earth’s Place in the Universe	
ESS1.A:	The universe and its stars
ESS1.B:	Earth and the solar system
ESS1.C:	The history of planet Earth
ESS2: Earth’s Systems	
ESS2.A:	Earth materials and systems
ESS2.B:	Plate tectonics and large-scale system interactions
ESS2.C:	The roles of water in Earth’s surface processes
ESS2.D:	Weather and climate
ESS2.E:	Biogeology
ESS3: Earth and Human Activity	
ESS3.A:	Natural Resources
ESS3.B:	Natural Hazards
ESS3.C:	Human Impacts on Earth Systems

2. Life Science. The *Framework* for the Life Sciences (LS) states that the DCIs “focus on patterns, processes, and relationships of living organisms” (National Research Council, 2012, p. 139). Table 1.1.2 presents the four LS DCIs and their underlying topics.

Table 1.1.2: Life Science DCIs

DCI Topic Description	
LS1: From Molecules to Organisms: Structures and Processes	
LS1.A:	Structure and function
LS1.B:	Growth and development of organisms
LS1.C:	Organization for matter and energy flow in organisms
LS1.D:	Information processing
LS2: Ecosystems: Interactions, Energy, and Dynamics	
LS2.A:	Interdependent relationships in ecosystems
LS2.B:	Cycles of matter and energy transfer in ecosystems
LS2.C:	Ecosystem dynamics, functioning, and resilience
LS2.D:	Social interactions and group behavior
LS3: Heredity: Inheritance and Variation of Traits	
LS3.A:	Inheritance of traits
LS3.B:	Variation of traits
LS4: Biological Evolution: Unity and Diversity	
LS4.A:	Evidence of common ancestry and diversity
LS4.B:	Natural selection
LS4.C:	Adaptation
LS4.D:	Biodiversity and humans

3. Physical Science. According to the *Framework*, the goal of learning Physical Science (PS) “is to help students see that there are mechanisms of cause and effect in all systems and processes that can be understood through a common set of physical and chemical principles” (National Research Council, 2012, p. 103). Table 1.1.3 illustrates the four PS DCIs along with the associated detailed topics for each.

Table 1.1.3: Physical Science DCIs

DCI Topic Description
PS1: Matter and its Interactions
PS1.A: Structure and matter
PS1.B: Chemical reactions
PS2: Motion and Stability: Force and Interactions
PS2.A: Force and motion
PS2.B: Types of interactions
PS2.C: Stability and instability in physical systems
PS3: Energy
PS3.A: Definitions of energy
PS3.B: Conservation of energy and energy transfer
PS3.C: Relationship between energy and forces
PS3.D: Energy in chemical processes and everyday life
PS4: Waves and their Applications in Technologies for Information Transfer
PS4.A: Wave properties
PS4.B: Electromagnetic radiation
PS4.C: Information technologies and instrumentation

Scientific and engineering practices. The *Framework* contains eight different Scientific and Engineering Practices (SEPs). One of the goals of the SEPs is to help “students understand how scientific knowledge develops; such direct involvement gives them an appreciation of the wide range of approaches that are used to investigate, model, and explain the world” (National Research Council, 2012, p. 42). Within the context of the NJSSA–S, the SEPs are consolidated into three categories of scientific practices: Investigating, Sensemaking, and Critiquing. Table 1.1.4, adapted from the work of McNeill et al. (2015), shows how the eight *Framework* SEPs were consolidated for the purposes of the NJSSA–S.

Table 1.1.4: SEP Consolidation

SEP	Grouping
Asking questions and defining problems (AQDP)	Investigating
Planning and carrying out investigations (PACI)	Investigating
Using mathematics and computational thinking (UMCT)	Investigating
Analyzing and interpreting data (AID)	Sensemaking
Constructing explanations and designing solutions (CEDS)	Sensemaking
Developing and using models (DUM)	Sensemaking
Engaging in argument from evidence (EAE)	Critiquing
Obtaining evaluating and communicating information (OECI)	Critiquing

1. Investigating. Investigating Practices (McNeill et al., 2015) involve asking questions, conducting investigations, and using mathematical skills to probe naturally occurring phenomena. Table 1.1.5 delineates the *Framework* (National Research Council, 2012) definition of each of the Investigating Practices.

Table 1.1.5: Investigating Practices

SEP	National Research Council Framework
Asking questions and defining problems (AQDP)	Students at any grade level should be able to ask questions of each other about the texts they read, the features of the phenomena they observe, and the conclusions they draw from their models or scientific investigations. For engineering, they should ask questions to define the problem to be solved and to elicit ideas that lead to the constraints and specifications for its solution. (p.56)
Planning and carrying out investigations (PACI)	Students should have opportunities to plan and carry out several different kinds of investigations during their K–12 years. At all levels, they should engage in investigations that range from those structured by the teacher — in order to expose an issue or question that they would be unlikely to explore on their own (e.g., measuring specific properties of materials) — to those that emerge from students’ own questions. (p. 61)
Using mathematics and computational thinking (UMCT)	Although there are differences in how mathematics and computational thinking are applied in science and in engineering, mathematics often brings these two fields together by enabling engineers to apply the mathematical form of scientific theories and by enabling scientists to use powerful information technologies designed by engineers. Both kinds of professionals can thereby accomplish investigations and analyses and build complex models, which might otherwise be out of the question. (p. 65)

2. Sensemaking. Sensemaking Practices (McNeill et al., 2015) are conceptualized as analyzing the data that is produced from an investigation and developing models and explanations that can explain naturally occurring phenomena. Table 1.1.6 illustrates the *Framework* (National Research Council, 2012) definition of each of the Sensemaking Practices.

3. **Table 1.1.6 Sensemaking Practices**

SEP	National Research Council Framework
Developing and using models (DUM)	Modeling can begin in the earliest grades, with students’ models, progressing from concrete “pictures” and/or physical scale models (e.g., a toy car) to more abstract representations of relevant relationships in later grades, such as a diagram representing forces on a particular object in a system. (p. 58)
Analyzing and interpreting data (AID)	Once collected, data must be presented in a form that can reveal any patterns and relationships and that allows results to be communicated to others. Because raw data as such have little meaning, a major practice of scientists is to organize and interpret data through tabulating, graphing, or statistical analysis. Such analysis can bring out the meaning of data — and their relevance — so that they may be used as evidence. (p. 61)
Constructing explanations and designing solutions (CEDS)	Asking students to demonstrate their own understanding of the implications of a scientific idea by developing their own explanations of phenomena, whether based on observations they have made or models they have developed, engages them in an essential part of the process by which conceptual change can occur. (p. 68)

3. Critiquing. Critiquing Practices (McNeill et al., 2015) are conceptualized as the ability of students to evaluate information, engage in argument, and communicate whether the models, explanations, or interpretations are adequate representations of naturally occurring phenomena. Table 1.1.7 shows the *Framework* (National Research Council, 2012) definition of each of the Critiquing Practices.

Table 1.1.7 Critiquing Practices

SEP	National Research Council Framework
Engaging in argument from evidence (EAE)	The study of science and engineering should produce a sense of the process of argument necessary for advancing and defending a new idea or an explanation of a phenomenon and the norms for conducting such arguments. In that spirit, students should argue for the explanations they construct, defend their interpretations of the associated data, and advocate for the designs they propose. (p. 73)
Obtaining evaluating and communicating information (OEI)	Any education in science and engineering needs to develop students’ ability to read and produce domain-specific text. As such, every science or engineering lesson is in part a language lesson, particularly reading and producing the genres of texts that are intrinsic to science and engineering. (p. 76)

1.2 Crosscutting Concepts

The *Framework* contains seven different Crosscutting Concepts (CCCs). They were selected to help “students with an organizational framework for connecting knowledge from the various disciplines into a coherent and scientifically based view of the world” (National Research Council, 2012, p. 83). Due to reporting constraints, the CCCs are the lowest priority of the three dimensions described in the *Framework*. However, because each item is aligned to a CCC, the CCC concepts and the knowledge, skills, and abilities associated with them are still being assessed by the NJSSA–S and contribute to the overall NJSSA–S score. Table 1.2.1 shows the CCCs being measured by the NJSSA–S.

Table 1.2.1: Crosscutting Concepts

CCC	National Research Council Framework (p. 84)
Patterns	Observed patterns of forms and events guide organization and classification, and they prompt questions about relationships and the factors that influence them.
Cause and Effect	Events have causes, sometimes simple, sometimes multifaceted. A major activity of science is investigating and explaining causal relationships and the mechanisms by which they are mediated. Such mechanisms can then be tested across given contexts and used to predict and explain events in new contexts.
Scale, Proportion, and Quantity	In considering phenomena, it is critical to recognize what is relevant at different measures of size, time, and energy and to recognize how changes in scale, proportion, or quantity affect a system’s structure or performance.
Systems and System Models	Defining the system under study — specifying its boundaries and making explicit a model of that system — provides tools for understanding and testing ideas that are applicable throughout science and engineering.
Energy and Matter	Tracking fluxes of energy and matter into, out of, and within systems helps one understand the systems’ possibilities and limitations.
Structure and Function	The way in which an object or living thing is shaped and its substructure determine many of its properties and functions.
Stability and Change	For natural and built systems alike, conditions of stability and determinants of rates of change or evolution of a system are critical elements of study.

Part 2: Test Development and Test Administration

The NJSSA–S is aligned to the New Jersey Student Learning Standards for Science (NJSLS–S), adopted in 2014, which in turn are based upon the National Research Council’s (2012) *Framework for K–12 Science Education* and the Next Generation Science Standards (NGSS; NGSS Lead States, 2013).

The Test Design and Development chapter within the *Standards for Educational and Psychological Testing* (American Educational Research Association, the American Psychological Association, and the National Council on Measurement in Education [AERA, APA, NCME], 2014) outlines a series of five primary phases of the test development process: (1) test specifications; (2) item development and review; (3) assembling and evaluating test forms; (4) development of procedures and materials for test administration and scoring; and (5) test revisions (p. 83). Sections 2.1, 2.2, 2.3, and 2.4 detail the NJSSA–S test specifications, item development processes and both the test construction processes and their results in 2022, respectively. Additionally, the development of procedures and materials for test administration scoring and score reports is covered in Sections 2.5 and 2.8, respectively. No test revisions were documented.

2.1 Test Specifications

According to the *Standards*, “[t]he term *test specifications* is sometimes limited to description of the content and format of the test. In the *Standards*, test specifications are defined more broadly to also include documentation of the purpose and intended uses of the test, as well as detailed decisions about content, format, test length, psychometric characteristics of the items and test, delivery mode, administration, scoring, and score reporting” (AERA, APA, NCME, 2014, p. 76).

As was described in Part 1 of this document, despite being administered to students entering grades 6, 9, and 12, the NJSSA–S was developed to measure the knowledge, skills, and abilities (KSAs) identified in the NJSLS–S in grades 5, 8, and 11. The test is designed to provide reporting information for the three student support levels (Strong Support, Some Support, and Less Support) at each of the three content domains (Earth and Space Science, Life Science, and Physical Science) and the three scientific practices (Investigating, Sensemaking, and Critiquing). The test specifications call for a balanced test design that prioritizes each science content domain and each DCI, each scientific practice, and each SEP, as well as all seven CCCs. The detailed information recommended in the NJSLS–S is presented in the sections that follow.

2.1.1 Test Blueprints

Table 2.1.1 depicts the NJSSA–S test blueprint for all grades. The table summarizes the ideal range of the numbers of items on the operational NJSSA–S for each of the six reporting categories.

Table 2.1.1: NJSSA–S Test Blueprints

Domain	Practice	Grade 6	Grade 9	Grade 12
PS	Investigating <i>AQDP, PACI, UMCT</i>	2–5	2–5	2–5
PS	Sensemaking <i>DUM, AID, CEDS</i>	2–5	2–5	2–5
PS	Critiquing <i>EAE, OECI</i>	2–5	2–5	2–5
PS	Total Items	7–10	7–10	7–10
LS	Investigating <i>AQDP, PACI, UMCT</i>	2–5	2–5	2–5
LS	Sensemaking <i>DUM, AID, CEDS</i>	2–5	2–5	2–5
LS	Critiquing <i>EAE, OECI</i>	2–5	2–5	2–5
LS	Total Items	7–10	7–10	7–10
ESS	Investigating <i>AQDP, PACI, UMCT</i>	2–5	2–5	2–5
ESS	Sensemaking <i>DUM, AID, CEDS</i>	2–5	2–5	2–5
ESS	Critiquing <i>EAE, OECI</i>	2–5	2–5	2–5
ESS	Total Items	7–10	7–10	7–10

2.1.2 Item Types

Two types of items comprise the NJSSA–S: multiple-choice (MC) and technology-enhanced (TE).

- MC items all have a key (A, B, C, or D) associated with them, and students are asked to select the best of the four options. MC items are scored dichotomously (0 or 1).
- TE items require students to interact with more complex methods of answering the items. Examples of TE item interactions include drop-down choice, hot spot, fill-in-the-blank, drag-and-drop, multiple selection, and ordering. Not all TE interaction types will be used on a single NJSSA–S form. TE items are either dichotomously scored or polytomously scored with at least three score points (e.g., 0, 1, or 2).

Table 2.1.2 describes each NJSSA–S item type.

Table 2.1.2: NJSSA–S Item Types

Item Type	Description
MC: Multiple Choice	Select one response from four possible options (A, B, C, D).
TE: Multiple Selection	Select two or more answer options.
TE: Short Answer	Type a brief constrained response to the question.
TE: Drop-Down Choice	Select from a drop-down menu embedded in the prompt.
TE: Ordering	Drag text or image-based options into a particular order.
TE: Drag and Drop	Place one or more text or graphic choices into blank spots within a sentence, table, or diagram.
TE: Matching in a Table	Check a box in the table to match the row to the column.
TE: Fill in the Blank	Type a response to fill in a blank within a text-based prompt.
TE: Scatter Plot	Plot one or more points on a graph.
TE: Bar Graph	Drag each bar to the correct length on the graph.
TE: Line Graph	Plot one or more lines on a graph.
TE: Slider	Slide an area within a graphic to change its length.
TE: Hot Spot	Select one or more regions on a graphic or image to identify an answer.
TE: Hot Text	Select one or more sentences within a paragraph of text.

2.2 Item Development Processes

All 2022 NJSSA–S items were originally used on the 2019 NJSLA–S operational assessment; thus, they were subjected to the NJSLA–S item development process, which is described in detail in the [2019 NJSLA-S Technical Report](#) (New Jersey Department of Education [NJDOE], 2020). The item development for the 2019 NJSLA–S was conducted by Measurement Incorporated (MI) and Pearson with oversight from NJDOE staff and the New Jersey Science Advisory Committee (NJSAC). The NJSLA–S item development process is extremely rigorous and involves item writers, content specialists, editors, graphic artists, programmers, scoring experts, and psychometricians. The resulting products are phenomenon-based scenarios (PBS) and items that are aligned to the NJSLA–S and the NJSSA–S reporting categories. The PBSs and their items are all housed in Pearson’s Assessment Banking for Building and Interoperability (ABBI) item banking system. ABBI is specifically designed to handle online, interactive, and accessible content.

The steps in the item development process and how they incorporate the principles of universal design (Thompson et al., 2002), are detailed in Part 2 of the 2019 NJSLA–S Technical Report details. Between the NJSAC and the New Jersey Bias and Sensitivity Committee (NJBSC), New Jersey educators and administrators were intimately and actively involved in the item development process and had to review and approve each NJSLA–S/NJSSA–S item multiple times.

2.3 Test Construction Process

The NJSSA–S test construction process ensures that the test forms balance the specifications set forth in the test blueprint, along with other psychometric constraints. Each form is built to measure students across the whole spectrum of ability levels and to foster valid interpretations of test scores in adherence to the *Standards for Educational and Psychological Testing* (AERA, APA & NCME, 2014). The steps and constraints associated with constructing the NJSSA–S operational tests are detailed in the following sections. An evaluation of the results of the test construction process is presented in Section 2.4.

2.3.1 Test Construction — First Draft

The first step in the NJSSA–S test construction process involves selecting items that had been previously used operationally on the NJSLA–S that best matched the NJSSA–S test blueprint and statistical constraints. The process of selecting items is contingent upon the availability of previously operationalized NJSLA–S items at each grade level. If specific content constraints are challenging to fulfill, then those content constraints are given priority in the initial selection of items. Next, items are selected iteratively based on which content constraints need to be fulfilled while simultaneously balancing the various statistical constraints. Detailed descriptions of the statistical constraints are presented in the sections below.

2.3.1.1 Test Construction Statistical Constraints

To ensure that the NJSSA–S operational test form is reliable and fosters valid interpretations, the following statistical requirements are used during the test construction process. Table 2.3.1 provides a summary of the NJSSA–S test construction requirements.

Item difficulty. Each test form is constructed to a specific difficulty level. The most important decisions made from the NJSSA–S are at the Some Support and Less Support cut scores. To maximize the reliability of decisions at those cut scores, the target average item difficulty parameter for each test form is at the point on the NJSSA–S scale that maximizes test information at both of those decision points. See Section 3.1.1 for more detailed information about item difficulty and Section 5.2.1 for more information about the test information function.

Item discrimination. In this context, item discrimination refers to the ability of an item to differentiate between individuals with different levels of knowledge. An item that discriminates poorly could indicate ineffective measurement of the NJSSA–S scale and reduce test form reliability. Item discrimination is measured via the item-total correlation, which can range from –1.0 to 1.0. Item-total correlations close to 1.0 represent items that discriminate between individuals well. See Section 3.1.1 for more detailed information about item discrimination and item-total correlations. For the NJSSA–S, items with item-total correlations below 0.2 are only selected for placement on the NJSSA–S if no other viable options are available.

IRT model fit. The NJSSA–S uses an Item Response Theory (IRT) model called the Partial Credit Model (PCM; Masters, 1982) to estimate student ability levels. The PCM makes certain assumptions that, if violated, could impact the validity of interpretations made from NJSLA–S test scores. Statistical constraints based on PCM model-fit statistics include infit, outfit, discrimination and lower asymptote. See Section 3.2 for more details about IRT and Section 3.2.2.1 for more information about infit and outfit.

During test construction, the mean item infit, outfit, and discrimination statistics are all constrained to be as close to 1.0 as possible. If an individual item has an infit or outfit statistic outside the acceptable range of 0.7 to 1.3 or a discrimination statistic outside the acceptable range of 0.5 to 1.5, it is only used if no other viable options are available. The lower asymptote statistic is constrained to be as close to zero as possible; any item whose lower asymptote is greater than 0.1 is flagged and only used if necessary.

Time on items. The NJSSA–S is not designed to be a speeded test; consequently, almost all students should be able to finish it within the allotted time. Items are selected to ensure the median time spent on the test is well below the time limit. If the median time spent on items is greater than the total test time minus 30 minutes, then items that are taking students too long are replaced by items that take less time, unless no other options are available.

Differential Item Functioning. Differential Item Functioning (DIF) exists when different groups of students have different probabilities of getting an item correct, after accounting for their ability levels. NJSSA–S comparison groups include Male/Female, White/Black, White/Hispanic, and White/Asian. If any item favors one group over another based on the ETS Mantel-Haenszel (Dorans & Holland, 1993; Zieky, 1993) method for dichotomously scored items, or Penfield’s (2007) DIF classification method for polytomously scored items, that item is classified as demonstrating either “B” or “C” level DIF. All items classified as exhibiting either “B” or “C” level DIF are reviewed by the New Jersey Bias and Sensitivity Committee (NJBSC) during the statistical review process. If the NJBSC deems an item biased, then it is ineligible for placement on the operational NJSSA–S regardless of DIF classification. It should be noted that a small number of “B” items can be used to maintain the test blueprint, whereas “C” items are not used on the operational NJSSA–S. See Section 3.1.3 for more information on differential item functioning.

Table 2.3.1: Summary of NJSSA–S Test Construction Statistical Constraints

Statistical Constraint	Description
Item Difficulty	Average item difficulty maximizes information at both the Some Support and Less Support cut scores.
Item Discrimination	Items have item-total correlations greater than 0.2.
IRT Model Fit	<ul style="list-style-type: none"> • Item Infit and Outfit statistics range from 0.7 to 1.3 and average 1.0. • Item Discrimination statistics range from 0.5 to 1.5 and average 1.0. • Item Lower Asymptote statistics < 0.1 and average as close to 0.0 as possible.
Time On Items	Total median time on items < (total test time — 30 minutes).
DIF	<ul style="list-style-type: none"> • “B” items are only used if necessary. • “C” items are not used.

2.3.2 Test Construction Content Review

After MI’s psychometric staff finishes the first draft of the NJSSA–S test forms, content specialists at each grade level check the forms to ensure that no items cue each other or have content that is too similar. The content review is an iterative process between content specialists and psychometricians. If during the review content specialists identify items that are too similar or that cue each other, they alert MI’s psychometric staff, and the items are replaced. The content review then resumes until the test matches NJSSA–S content and statistical requirements.

2.3.3 Test Construction NJDOE Review

All of the NJSSA–S test forms are reviewed and approved by NJDOE. Once the content and psychometrics teams have agreed upon the operational test forms, they are sent to NJDOE for approval. After NJDOE approves the test forms, they are released for final editorial review and publishing.

2.4 2020 NJSSA–S Test Construction

2020 was the first year of NJSSA–S test construction. Overall, the test construction process achieved forms that matched the balance required by the test blueprint. The science content domains were well balanced at each grade level. Moreover, all grade levels met the requirement that no more than 50% of the items be multiple choice. However, some constraints were more difficult to achieve. At all three grade levels, it was challenging to identify enough Critiquing items that were acceptable from a content and statistical perspective to balance out the three scientific practice reporting categories.

A final content constraint that was not met was the balance between the three content domains across the three scientific practice reporting categories, as shown in the test blueprint in Section 2.1.1. The items associated with each scientific practice were meant to be balanced across all content domains. Table 2.4.1 shows this lack of balance. At each grade level, at least

one scientific practice was over-represented for a given content domain. For instance, of the eight Earth and Space Science points available on the grade 6 test, six were aligned to the Sensemaking practice, whereas only one point each was aligned to the Investigating or Critiquing practices.

Table 2.4.1: NJSSA–S Points Available by Domain and Practice

Grade	Practice	Earth and Space Science	Life Science	Physical Science
6	Investigating	1	3	4
6	Sensemaking	6	1	4
6	Critiquing	1	3	2
9	Investigating	4	1	4
9	Sensemaking	2	5	3
9	Critiquing	2	1	3
12	Investigating	3	0	5
12	Sensemaking	5	4	2
12	Critiquing	2	3	1

2.4.1 Grade 6 Test Construction

At grade 6 the science content domains were balanced, as illustrated in Table 2.4.2. The least balanced content domain was Life Science, and it still made up 7 points of the 25 total score points. The scientific practices were slightly less balanced, with only 6 out of 25 points being allocated to the Critiquing reporting category and 11 out of 25 points allocated to the Sensemaking reporting category. Table 2.4.2 details the item and point totals for each of the six reporting categories. Tables 2.4.3 through 2.4.5 show the distributions of DCIs, SEPs, and CCCs.

Table 2.4.2: NJSSA–S Grade 6 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	Items	Points
Earth and Space Science	4	4	8	8
Life Science	1	6	7	7
Physical Science	2	8	10	10
Total — Domains	7	18	25	25
Investigating	2	6	8	8
Sensemaking	4	7	11	11
Critiquing	1	5	6	6
Total — Practices	7	18	25	25

Table 2.4.3: NJSSA–S Grade 6 DCIs

DCI	Items	Points
ESS1	0	0
ESS2	8	8
ESS3	0	0
LS1	4	4
LS2	0	0
LS3	0	0
LS4	3	3
PS1	3	3
PS2	3	3
PS3	4	4
PS4	0	0

Table 2.4.4: NJSSA–S Grade 6 SEPs

SEP	Items	Points
AQDP	2	2
PACI	4	4
UMCT	3	3
DUM	1	1
AID	6	6
CEDS	3	3
EAE	6	6
OECI	0	0

Table 2.4.5: NJSSA–S Grade 6 CCCs

CCC	Items	Points
C & E	5	5
E & M	4	4
Patterns	9	9
S & SM	0	0
S, P, & Q	2	2
SC	1	1
SF	4	4

The statistical constraints for the grade 6 NJSSA–S operational test form were generally met. The only exception was that the item difficulty maximized information close to the Some Support cut score instead of between the Some and Less Support cuts. Thus, the test was measuring the Some Support cut score at a higher level of reliability than the Less Support cut score. Otherwise, all items had item-total correlations above the 0.2 threshold, and very few items were flagged for divergent item-fit statistics. The median test time of 30.33 minutes was well below the 45-minute threshold, and out of 100 DIF classifications, there were zero “C” values and only two “B” values. All “B” DIF items were approved for operational test use by the

NJBSC as described in Section 2.3.1.1. Tables 2.4.6 and 2.4.7 summarize the test construction and DIF statistics.

Table 2.4.6: NJSSA–S Grade 6 Test Construction Statistics

Statistics	Average	Target	Flags
Item Difficulty	−0.29	0.31	N/A
Item Total Correlation	0.42	> 0.35	0
Infit	0.97	1.00	0
Outfit	0.99	1.00	2
Item Discrimination	1.05	1.00	0
Lower Asymptote	0.02	0.00	1
Median Time	30.33	< 45	N/A

Table 2.4.7: NJSSA–S Grade 6 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	23	2	0
White/Black	25	0	0
White/Hispanic	25	0	0
White/Asian	25	0	0

2.4.2 Grade 9 Test Construction

The grade 9 NJSSA–S content constraints were generally met. The science content domains were balanced. Each content domain was represented by 7 to 10 points worth of items. The scientific practices were slightly less balanced, with only 6 of 25 points allocated to the Critiquing reporting category. All eleven DCIs, seven of eight SEPs, and all seven CCCs were represented. Table 2.4.8 details the item and point totals for each of the six reporting categories; Tables 2.4.9 through 2.4.11 show the distributions of DCIs, SEPs, and CCCs for grade 8.

Table 2.4.8: NJSSA–S Grade 9 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	Items	Points
Earth and Space	4	4	8	8
Life	1	6	7	7
Physical	3	6	9	10
Total — Domains	8	16	24	25
Investigating	6	3	9	9
Sensemaking	1	9	10	10
Critiquing	1	4	5	6
Total — Practices	8	16	24	25

Table 2.4.9: NJSSA–S Grade 9 DCIs

DCI	Items	Points
ESS1	2	2
ESS2	3	3
ESS3	3	3
LS1	2	2
LS2	3	3
LS3	1	1
LS4	1	1
PS1	1	1
PS2	3	3
PS3	3	4
PS4	2	2

Table 2.4.10: NJSSA–S Grade 9 SEPs

SEP	Items	Points
AQDP	5	5
PACI	2	2
UMCT	2	2
DUM	4	4
AID	3	3
CEDS	3	3
EAE	5	6
OECI	0	0

Table 2.4.11: NJSSA–S Grade 9 CCCs

CCC	Items	Points
C & E	5	5
E & M	5	6
Patterns	4	4
S & SM	2	2
S, P, & Q	2	2
SC	1	1
SF	5	5

The statistical constraints for the grade 9 NJSSA–S operational test form were met. The average item difficulty parameter was only 0.05 logits from the target value. Only one grade 9 item was flagged for having item-total correlations below the 0.2 threshold. The infit, outfit, and PCM item discrimination model-fit statistics were all close to their ideal values of 1.00. The median test time of 28.67 minutes was well below the 45-minute threshold, and out of 96 DIF classifications, there were zero “C” values and only one “B” value. The “B” DIF item was approved for operational test use by the NJBSC as described in Section 2.3.1.1. Tables 2.4.12 and 2.4.13 summarize the test construction and DIF statistics.

Table 2.4.12: NJSSA–S Grade 9 Test Construction Statistics

Statistics	Average	Target	Flags
Item Difficulty	–0.20	–0.25	N/A
Item Total Correlation	0.39	> 0.35	1
Infit	0.97	1.00	0
Outfit	0.97	1.00	2
Item Discrimination	1.06	1.00	1
Lower Asymptote	0.02	0.00	0
Median Time	28.67	< 45	N/A

Table 2.4.13: NJSSA–S Grade 9 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	23	1	0
White/Black	24	0	0
White/Hispanic	24	0	0
White/Asian	24	0	0

2.4.3 Grade 12 Test Construction

The grade 12 NJSSA–S content constraints were generally met. The content domains were close to being balanced. Each content domain had between 7 and 10 points worth of items. The scientific practices were less balanced, with only 6 out of 25 points being allocated to the Critiquing reporting category. The Sensemaking practice was overrepresented with 11 out of 25 points. Table 2.4.14 details the item and point totals for each of the six reporting categories; Tables 2.4.15 through 2.4.17 show the distributions of DCIs, SEPs, and CCCs at grade 12.

Table 2.4.14: NJSSA–S Grade 12 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	Items	Points
Earth and Space	5	5	10	10
Life	1	6	7	7
Physical	4	4	8	8
Total – Domains	10	15	25	25
Investigating	4	4	8	8
Sensemaking	5	6	11	11
Critiquing	1	5	6	6
Total – Practices	10	15	25	25

Table 2.4.15: NJSSA–S Grade 12 DCIs

DCI	Items	Points
ESS1	3	3
ESS2	3	3
ESS3	4	4
LS1	0	0
LS2	3	3
LS3	0	0
LS4	4	4
PS1	3	3
PS2	2	2
PS3	3	3
PS4	0	0

Table 2.4.16: NJSSA–S Grade 12 SEPs

SEP	Items	Points
AQDP	2	2
PACI	2	2
UMCT	4	4
DUM	3	3
AID	7	7
CEDS	1	1
EAE	2	2
OECI	4	4

Table 2.4.17: NJSSA–S Grade 12 CCCs

CCC	Items	Points
C & E	4	4
E & M	1	1
Patterns	5	5
S & SM	8	8
S, P, & Q	4	4
SC	3	3
SF	0	0

The statistical constraints for the grade 12 NJSSA–S operational test form were more challenging to meet than for the other two grades. The average item difficulty parameter was 0.36 logits below the target, meaning the test was more reliable towards the Some Support cut score than at the Less Support cut score. One grade 12 item was flagged for having an item-total correlation below the 0.2 threshold. The infit, outfit, PCM item discrimination, and lower asymptote model-fit statistics did not meet their goals, and more items were flagged than in grades 6 and 9 combined. The increase in model-fit flags was due to the relatively large percentage of items that were flagged for the 2019 NJSLA–S operational test, as was documented in the 2019 NJSLA–S Technical Report (NJDOE, 2019). The median test time was

only 20.28 minutes, which was over 20 minutes below the 45-minute constraint. Of 100 DIF classifications, there were zero “C” values and seven “B” values. All “B” DIF items were approved for operational test use by the NJBSC as described in Section 2.3.1.1. Tables 2.4.18 and 2.4.19 summarize the test construction and DIF statistics for grade 12.

Table 2.4.18: NJSSA–S Grade 12 Test Construction Statistics

Statistics	Average	Target	Flags
Item Difficulty	-0.28	0.08	N/A
Item Total Correlation	0.43	> 0.35	1
Infit	1.02	1.00	1
Outfit	1.01	1.00	3
Item Discrimination	0.95	1.00	4
Lower Asymptote	0.05	0.00	6
Median Time	20.28	< 45	N/A

Table 2.4.19: NJSSA–S Grade 12 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	24	1	0
White/Black	23	2	0
White/Hispanic	23	2	0
White/Asian	23	2	0

2.5 Test Administration

The Start Strong Assessments were available for administration using Pearson’s TestNav online delivery system. Support for educators, students, and caregivers was available at the [New Jersey Assessments Resource Center — Start Strong](#). The Getting Started Online section included tutorials and practice tests to familiarize students with the online testing experience.

Administration procedures were standardized as much as possible. Two administration guides and a policy guide were made available — the Start Strong PBT User Guide, the Start Strong CBT User Guide, and the Start Strong Administration Policies Guide. The user guides provided detailed instructions for starting and ending the administration as well as allowable and unallowable supports that may be provided to students while taking the assessment. Like the NJSLA summative assessment, assistance that supported student responses was discouraged because any deviation from normal administration conditions threatens inferences made from the results.

Additionally, if recommended by the teacher, an app-based testing lockdown of the desktop may have occurred to provide a focused testing experience. However, a lockdown was not a requirement for the Start Strong assessments. Therefore, when interpreting the results, it is important to note that the Start Strong administrations are not considered secure.

2.6 Test Registration

Student registration consisted of a simplified version of the normal NJSLA summative assessment registration process with a streamlined Student Registration/Personal Needs Profile (SR/PNP) using PearsonAccess^{next}. Teachers created sessions at the classroom level, generated testing tickets, and provided login information for students to take the assessment at home or in the classroom. Students accessed the assessments online with teacher-provided usernames and passwords.

2.7 Test Accessibility Features and Accommodations

Standard 3.9 states that “[t]est developers and/or test users are responsible for developing and providing test accommodations, when appropriate and feasible, to remove construct-irrelevant barriers that otherwise would interfere with examinees’ ability to demonstrate their standing on the target constructs” (AERA, APA, NCME, 2014, p. 67). Federal and state regulations require that all students — including those classified as English learners (EL) and those with disabilities — be included in the statewide assessment program and assessed annually. The Every Student Succeeds Act of 2015 (ESSA) mandates that all states must test science one time each in three different grade bands: 3–5, 6–8, and 9–12. To ensure that the diverse population of students taking the NJSSA–S is tested under appropriate conditions and to adhere to the principles of universal design (Thompson et al., 2002), NJDOE has adopted test accommodations and accessibility features that may be used when testing special populations of students. The content of the test remains the same, but administration procedures, setting, and answer modes may be adapted. Students requiring accommodations may be tested in a separate location from general education students.

The [*NJSLA Accessibility Features and Accommodations Manual \(AF&A Manual\)*](#); NJDOE, 2022) is available online at the [New Jersey Assessments Resource Center](#) under Educator Resources > Accessibility Features and Accommodations (AF&A) Resources > *NJSLA & NJGPA Accessibility Features and Accommodations, 11th Edition* — Updated April 2023. It contains detailed information about each accessibility feature and accommodation. Schools must refer to the *AF&A Manual* for full information about identifying and administering accessibility features and accommodations.

2.7.1 Accessibility Features

The purpose of accessibility features is to ensure that a diverse population of students is being tested fairly and that construct-irrelevant factors are not unduly impacting their test scores. According to the *NJSLA–S AF&A Manual*, accessibility features are defined as “tools or preferences that are either built into the testing platform or provided externally by Test Administrators” (NJDOE, 2022, p. 2). All students have access to accessibility features. However, for some accessibility features to be available for students during testing, an administrator must have identified the student as needing the accessibility feature prior to testing. It is

essential that students using accessibility features get to practice with them prior to operational testing. Thus, NJSLA–S practice tests that contain the accessibility features are available throughout the year at the following link: measinc-nj-science.com.

2.7.1.1 Text-to-Speech

The most used NJSSA–S accessibility feature is Text-to-Speech (TTS). Prior to testing, an administrator activates the TTS accessibility feature for individual students. When the selected student gets placed into a testing session, their form automatically defaults to the designated TTS form. During testing the student can select the TTS player, and the test will be read aloud to them via the TTS software embedded within TestNav. Students using the TTS accessibility feature must be wearing headphones. The items on the TTS form all contain the same phenomenon-based scenarios, item stems, and response options as are presented to the students taking the traditional computer-based test (CBT) form. All final TTS forms are verified by NJDOE to verify that the TTS functionality is working correctly.

2.7.2 Accommodations

The role of accommodations is to minimize the impact of students’ disabilities or English language proficiency level on their assessment performance. The NJSLA–S *AF&A Manual* defines an accommodation as “ ...[an adjustment] to the testing conditions, test format, or test administration that provide[s] equitable access during assessments for students with disabilities and students who are ELs” (NJDOE, 2022, p. 14). Accommodations are only available to students who have an Individualized Education Program (IEP), a Section 504 plan, or an English learner (EL) plan.

Different accommodations are necessary depending on whether the test was administered using a CBT or paper-based test (PBT) format. Per NJDOE policy, each student who received a PBT version of the NJSSA–S had an appropriate accommodation. No physical test materials were automatically shipped for Start Strong. Test coordinators placed orders in PearsonAccess^{next} for braille and large-print test kits.

A comprehensive explanation of each NJSSA–S accommodation is presented in the NJSLA–S *AF&A Manual* (NJDOE, 2022). The NJSSA–S’ CBT accommodations include:

- Assistive Technology — Screen Reader
- Assistive Technology — Non-Screen Reader
- American Sign Language (ASL) Text-to-Speech (TTS)
- Human Reader
- Spanish
- Spanish Text-to-Speech
- Spanish Human Reader

PBT accommodations are received as kits, which include:

- Braille
- Large Print
- Alternate Representation
- Spanish
- Spanish Large Print

2.7.2.1 Accommodated Test Form Development

The *Standards* state that “an appropriate accommodation is one that responds to specific individual characteristics but does so in a way that does not change the construct the test is measuring or the meaning of the scores” (AERA, APA, NCME, 2014, p. 67). Each of the accommodated test forms requires specific processes to ensure they are addressing the needs of their intended users. After NJDOE approval, the accommodated test forms are sent to various subcontractors so that they can adapt the items to Spanish, Braille, and American Sign Language (ASL). The adaptation processes for those forms are presented in Sections 2.7.2.1.1 through 2.7.2.1.3. The PBT form adaptation process is presented in Section 2.7.2.1.4. Following adaptation, NJDOE verifies each accommodated test form.

2.7.2.1.1 Spanish. All Spanish accommodations were made by Teneo Linguistics Company (TLC). TLC received the NJDOE-approved tests and created the translations within ABBI. Once the items were translated, a committee of Spanish teachers from New Jersey reviewed the items online, with TLC representatives in attendance. Edits were made during the review, and then the final versions of the online forms were verified by NJDOE. The translation that was created for the online version was then used to create the paper version of the Spanish tests.

2.7.2.1.2 Braille. All braille accommodations were created by the National Braille Press (NBP). NBP received downloadable paper versions of the operational test forms. NBP then provided MI with feedback about any items that were unable to be brailled. Once the tests were brailled, external reviewers received the draft braille versions and reviewed them for any issues a student might have taking the braille tests. For the 2022 NJSSA–S, all items were able to be brailled.

2.7.2.1.3 American Sign Language. All ASL accommodations were created by the ADS Group in Plymouth, MN. The ADS Group provided ASL video production with two ASL content specialist translators and one ASL proofer. Their video production engineer provided studio editing. Additionally, they provided proofing/QC services as well as closed captioning. Once NJDOE approved the operational test forms, the ADS Group created videos of American Sign Language for each item. These items were verified by external expert reviewers under the guidance of MI.

2.7.2.1.4 Paper-Based Test. The conversion of the NJSSA–S CBT into PBT form was undertaken by MI’s Editorial Department. Most PBT items were exactly the same as their CBT counterparts. However, some aspects needed adaptation. The following bullets represent the major changes made to the stimuli and items during the adaptation processes:

- All artwork was converted from color to grayscale.
- Video items were converted to still images. This was accomplished by MI's Editorial staff working in conjunction with content specialists to select specific frames from the video that effectively conveyed its essence. In some cases, the captured images were redrawn to ensure that no essential information was lost in the adaptation process.
- TE items were converted to PBT format via multiple methods depending on the TE item type.

2.8 Scores and Score Reports

This section details the item scoring, adjudication process, test scores, and reports for the NJSSA–S.

2.8.1 Machine-Scored Items

All multiple-choice (MC) and technology-enhanced (TE) items are machine-scored. Each item has a key (correct answer) associated with it, which has been supplied and verified by content specialists and approved by the Department prior to test administration. All student responses are machine-scored based on these prior approved keys. The data from the student responses are then screened via Pearson's Customer Data Quality (CDQ) team. The CDQ team verifies the accuracy of the student responses and metadata within two file types: the Summative Record File (SRF) and the Item Response File (IRF). Verification steps include validating variable acceptable ranges, computing raw overall scores and subscores, validating ID numbers and unique item numbers (UINs), and flagging inconsistent student records for investigation. Once the data have been verified, the files are placed on a Secure File Transfer Protocol site from which they are retrieved by MI's IT group. The IT group then prepares the files for psychometric analysis.

2.8.2 Adjudication

Adjudication involves the careful review of all student responses to an item to ensure that its key was applied correctly and that no possible correct answer has been overlooked in the many prior key checks. All machine-scored items are adjudicated. The Content Development's Psychometric departments use the SRF and the IRF to analyze the student response patterns for each item. The response patterns are simple for items with limited possible options; for instance, an MC item only has five possible student responses (A, B, C, D, or blank). However, some TE items can have thousands of different student responses. The student response data is used to produce one file for each operational item that contains a Response ID, the point value associated with it (i.e., 0, 1, or 2), the total number and percentage of students selecting each response, the text of the response (retrieved from the item's XML coding), and the item-total correlation associated with each response option that was selected more than 100 times. Item means and item-total correlations are also calculated at the item level, and items are flagged for aberrant behavior.

The role of the content specialist during the adjudication process is to use the information housed in the adjudication files to identify any possible miskeys. They are instructed to first check items that were flagged for having low item means and item-total correlations because those statistics could indicate that the item is not performing as intended. Next, they look at combinations of student responses that are keyed as receiving “0” points but have item-total correlations above 0. That combination of response-level data could also be an indication of a possible student response that deserves credit for a correct response but has been keyed as incorrect. Finally, through a sorting process, the content specialists can relatively quickly review all other combinations of student responses. If there are any miskeys, key changes are submitted to the Department, and upon approval, subsequently corrected in the SRF and IRF. These steps are essential to ensuring both the reliability of student test scores and their valid interpretations.

2.8.3 Test Scores

Student performance is reported using an overall raw score (i.e., the number of points earned). While the raw score can be used to compare students who took the same assessment (e.g., grade 6 science), it cannot be used to compare students from a science assessment to students who took the mathematics assessment, nor can it be used to compare students in the 6th grade to students in the 9th grade. Because the Start Strong Assessment is a classroom assessment for gauging where students are in their learning of previous content standards, converting the raw score to a percent correct for the purpose of assigning a grade is not appropriate.

2.8.4 Support Level

Students are categorized into one of three support levels based on their individual total raw scores. Each support level is defined by a range of overall raw scores. There are three support levels for the Start Strong Assessment:

- Level 3 — Less Support May Be Needed
- Level 2 — Some Support May Be Needed
- Level 1 — Strong Support May Be Needed

Students performing at Level 3 may not require additional academic/instructional support in the tested content area, while students in Level 1 will likely benefit from additional academic/instructional support in the tested content area. The Start Strong performance levels are meant to indicate the amount of support a student might require. While these performance levels leverage the NJSLA–S summative cut scores, they are not intended to assign proficiency or mastery because the purpose and blueprint of the Start Strong assessments are different from those of the NJSLA–S.

The Start Strong Assessment’s primary purpose was to provide instructional information to classroom teachers about students’ needs for additional support upon returning to school in the fall of 2022. The information provided by this assessment is a snapshot of a student’s

understanding and should only be used with other supporting evidence (assignments, homework, etc.) when drawing conclusions about a student’s overall academic performance. Examples and further documentation of each report are available on the [New Jersey Assessments Resource Center — Practice Tests](#) and [Start Strong Testing Resources](#).

2.8.5 Student-Level Reports

Three student-level reports were produced for the Start Strong Assessment and are available via PearsonAccess^{next} (PAN).

- The **OnDemand Student Report** (ODR) is the first report for the Start Strong Assessments. It shows the student’s support level and the scores on each reporting concept. Only students who received a score will receive an OnDemand Student Report.
- The **Student Performance Item Level Report** allows users to compare the support level assigned to individual students within a group and then drill down to an individual student’s response to each item. This can be useful for understanding what misconceptions students may have.
- The **Individual Student Reports** (ISRs) are the last type of report to be released for the Start Strong Assessments. Users will be able to download PDFs of ISRs from PAN; school districts will also receive hard copies to distribute to students’ parents or guardians. ISRs will be shipped for both testing sites and accountable schools if different schools are involved.

2.8.6 Classroom-, School-, and District-Level Reports

In addition to the student-level reports described above, appropriate users will also have access to the Results by Question Reports and Support Level Reports via PearsonAccess^{next}.

- The **Results by Question Report** provides users with group-level information about student performance on specific items or standards. The Results by Question Report has two different ways to view information: the question list and the student list. The default view is the question list. You can switch between the two views by checking the “Show Students” checkbox at the bottom of the list. Drilling down to scores on individual test items enables the teacher to corroborate, verify, or otherwise build upon test information to identify instructional needs at the individual student or group level. This aids in the design and delivery of effective educational methods to meet these needs.
 - **Question List** — The question list shows items in numerical order, along with the standard(s) to which each item is aligned, the reporting concept(s) the item is associated with, and the number and percentage of students who answered the item correctly, incorrectly, and partially correctly (for those items that are worth more than 1 point).

- **Show Students** — Users can view individual student results by question. Selected students are sorted by last name, first name, middle name, and then the Statewide Student Identifier (SSI). Questions for only a single standard can be displayed at one time, and a standard is automatically selected by default. You may select a different standard in the drop-down box above the student list.
- The **Support Level Reports** display the overall distribution of support levels for a group of students on a particular test which can be filtered by school, grade, or demographic information (i.e., gender, ethnicity, students with disabilities, etc.). The groupings are completely flexible and can be defined to include any specific students of interest.

Part 3: Item and Test Statistics

Standard 5.0 states that “[t]est scores should be derived in a way that supports the interpretations of test scores for proposed uses of tests. Test developers and users should document evidence of fairness, reliability, and validity of test scores for their proposed uses” (AERA, APA, NCME, 2014, p.102). The NJSSA–S was designed to support inferences based on the classification of students into three support levels, as has been described throughout this technical brief. The interpretations of the support level classifications are dependent upon the test performing as intended. As was described in Section 2.3, the NJSSA–S was constructed using a combination of statistics and numerous content constraints. The following sections detail how well the 2022 NJSSA–S performed based on those statistics and other criteria. Detailed test maps containing item metadata and various statistics are presented in [Appendix A](#).

The data for these and all subsequent analyses were verified by Pearson’s Customer Data Quality (CDQ) team. Responses from students who did not attempt any items or who had their test scores voided were removed from the data set prior to analysis. NJDOE set the threshold for attemptedness as any student who made a legitimate student response to at least one item. Student responses were voided for cheating, security breaches, or other reasons.

3.1 Classical Test Theory Statistics

For each administration, a set of statistics based on CTT was generated prior to item calibration and scaling. The statistics can be grouped into measures of four concepts:

- Item Difficulty
- Item Discrimination
- Speededness
- Differential Item Functioning

These statistics were calculated for every operational item; each statistic provides some key information about the quality of each item from an empirical perspective. Descriptions of each type of statistic appear in the following sections.

3.1.1 Item Difficulty and Discrimination Descriptive Statistics

Monitoring item difficulty is essential for ensuring that the test is reliable and will foster valid test score interpretations. If items tend to be too challenging or too easy for a population of test-takers, then the reliability and validity of test score interpretations will suffer. In CTT, dichotomous item difficulty is assessed via the p-value, which is defined as the proportion of students who answered an item correctly. P-values can range from 0 to 1.00; an item with a high p-value is easier to answer correctly, whereas one with a low p-value is more challenging. Dichotomous items with p-values either below .25 or above .90 were flagged for review. For 0–2-point TE items, item difficulty is expressed as an item mean. The flagging criteria for polytomously scored items involve converting the item mean to a proportion by dividing it by the maximum points possible on the item (i.e., making it a p-value), then flagging the item if its

converted p-value falls outside of the .25 to .90 range. It should be noted that the flagging criteria are intended as a recommendation.

Item discrimination is also important to monitor because if items are unable to discriminate between students with different ability levels, then both the reliability of the items and the validity of test score interpretations might suffer. CTT item discrimination is expressed as the correlation between item scores and the total score of the remaining items on the test, the latter being a proxy for overall student ability. This item statistic is calculated as a point-biserial correlation, which is commonly denoted by r_{pb} or RPB. The values of RPB can range from -1.00 to 1.00 . Dichotomously scored items with values below .2 are flagged for review during the adjudication process. Polytomously scored items are expected to have higher item-total correlations; as such, the 0–2-point TE items are flagged if they show RPB values below .25.

Tables 3.1.1 through 3.1.6 summarize the average item difficulty and discrimination of the 2022 NJSSA–S items by item type. The averages within each of these tables are disaggregated by content domain and scientific practice. Tables 3.1.7 through 3.1.18 summarize frequency distributions for MC and TE item difficulty and discrimination; they are also disaggregated by content domain and scientific practice.

The average item difficulties and discriminations indicate that the items are productive for measuring students in New Jersey. At grade 6, the TE items tended to be more challenging but equally discriminating when compared to the MC items. At grade 9, the TE items were more challenging than the MC items; however, the MC items were slightly more discriminating. At grade 12 the TE items were, on average, more challenging and more discriminating than the MC items.

The frequency distributions of the item-total correlations also indicate that the items are productive for discriminating between high- and low-achieving students. Two items (one MC and one TE) at grade 9 and three items (two MC and one TE) at grade 12 had correlations below .20. Grade 6 had zero items with item-total correlations below .20. The p-value distributions, however, were less positive. At grade 12 no items had a p-value above .75 or below .25, indicating that there were almost no items at the easier and harder ends of the scale. At grade 9, five of the 16 TE items had p-values below .25, meaning that almost 33% of the grade 9 TE items were extremely challenging for New Jersey students. Most of the grade 6 items fell between .25 and .75; only two of the TE items had p-values below .25.

Table 3.1.1: Grade 6 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, MC

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	7	.66	.10	.44
Earth and Space Science	4	.63	.12	.41
Life Science	1	.71	N/A	.50
Physical Science	2	.70	.06	.48
Investigating	2	.61	.19	.40
Sensemaking	4	.68	.07	.45
Critiquing	1	.71	N/A	.50

Table 3.1.2: Grade 6 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, TE

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	18	.49	.16	.45
Earth and Space Science	4	.44	.17	.39
Life Science	6	.56	.12	.43
Physical Science	8	.47	.19	.50
Investigating	7	.49	.07	.49
Sensemaking	6	.56	.21	.41
Critiquing	5	.41	.19	.45

Table 3.1.3: Grade 9 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, MC

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	8	.46	.15	.37
Earth and Space Science	4	.51	.19	.40
Life Science	1	.37	N/A	.26
Physical Science	3	.42	.10	.37
Investigating	6	.40	.11	.34
Sensemaking	1	.73	N/A	.46
Critiquing	1	.52	N/A	.50

Table 3.1.4: Grade 9 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, TE

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	16	.38	.20	.34
Earth and Space Science	4	.39	.18	.31
Life	6	.47	.26	.42
Physical	6	.28	.10	.29
Investigating	3	.19	.05	.28
Sensemaking	9	.48	.19	.37
Critiquing	4	.29	.14	.33

Table 3.1.5: Grade 12 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, MC

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	10	.52	.11	.31
Earth and Space Science	5	.58	.10	.40
Life Science	1	.63	N/A	.37
Physical Science	4	.42	.04	.19
Investigating	4	.55	.13	.28
Sensemaking	5	.52	.12	.35
Critiquing	1	.38	N/A	.29

Table 3.1.6: Grade 12 Item Difficulty and Discrimination Summary Statistics by Domain/Practice, TE

Domain/Practice	# Items	Item Difficulty Mean	Item Difficulty S.D.	Item Discrimination Mean
NJSSA–S	15	.47	.12	.42
Earth and Space Science	5	.40	.12	.46
Life Science	6	.56	.11	.40
Physical Science	4	.42	.06	.42
Investigating	4	.45	.07	.49
Sensemaking	6	.48	.09	.39
Critiquing	5	.47	.19	.41

Table 3.1.7: Grade 6 Difficulty Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	p<.25	.25<=p<.50	.50<=p<.75	.75<=p<.90	p>=.90
NJSSA-S	7	.71	0	1	6	0	0
Earth and Space Science	4	.66	0	1	3	0	0
Life Science	1	.71	0	0	1	0	0
Physical Science	2	.70	0	0	2	0	0
Investigating	2	.61	0	1	1	0	0
Sensemaking	4	.69	0	0	4	0	0
Critiquing	1	.71	0	0	1	0	0

Table 3.1.8: Grade 6 Difficulty Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	p<.25	.25<=p<.50	.50<=p<.75	.75<=p<.90	p>=.90
NJSSA-S	18	.53	2	6	9	1	0
Earth and Space Science	4	.44	0	2	2	0	0
Life Science	6	.55	0	2	3	1	0
Physical Science	8	.50	2	2	4	0	0
Investigating	7	.51	0	3	4	0	0
Sensemaking	6	.65	0	2	3	1	0
Critiquing	5	.48	2	1	2	0	0

Table 3.1.9: Grade 6 Discrimination Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<= RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA-S	7	.41	0	1	2	2	2
Earth and Space Science	4	.40	0	1	1	1	1
Life Science	1	.50	0	0	0	1	0
Physical Science	2	.48	0	0	1	0	1
Investigating	2	.40	0	0	1	1	0
Sensemaking	4	.46	0	1	1	0	2
Critiquing	1	.50	0	0	0	1	0

Table 3.1.10: Grade 6 Discrimination Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<= RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA-S	18	.46	0	3	2	7	6
Earth and Space Science	4	.35	0	2	0	1	1
Life Science	6	.42	0	1	1	3	1
Physical Science	8	.50	0	0	1	3	4
Investigating	7	.53	0	1	0	2	4
Sensemaking	6	.42	0	2	1	2	1
Critiquing	5	.47	0	0	1	3	1

Table 3.1.11: Grade 9 Difficulty Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	p<.25	.25<=p<.50	.50<=p<.75	.75<=p<.90	p>=.90
NJSSA–S	8	.39	0	5	3	0	0
Earth and Space Science	4	.49	0	2	2	0	0
Life Science	1	.37	0	1	0	0	0
Physical Science	3	.41	0	2	1	0	0
Investigating	6	.37	0	5	1	0	0
Sensemaking	1	.73	0	0	1	0	0
Critiquing	1	.52	0	0	1	0	0

Table 3.1.12: Grade 9 Difficulty Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	p<.25	.25<=p<.50	.50<=p<.75	.75<=p<.90	p>=.90
NJSSA–S	16	.37	5	8	2	1	0
Earth and Space Science	4	.39	1	2	1	0	0
Life Science	6	.39	1	3	1	1	0
Physical Science	6	.25	3	3	0	0	0
Investigating	3	.19	3	0	0	0	0
Sensemaking	9	.41	0	6	2	1	0
Critiquing	4	.26	2	2	0	0	0

Table 3.1.13: Grade 9 Discrimination Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<= RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA-S	8	.39	1	1	2	3	1
Earth and Space Science	4	.40	0	0	2	2	0
Life Science	1	.26	0	1	0	0	0
Physical Science	3	.42	1	0	0	1	1
Investigating	6	.36	1	1	2	2	0
Sensemaking	1	.46	0	0	0	1	0
Critiquing	1	.50	0	0	0	0	1

Table 3.1.14: Grade 9 Discrimination Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<=RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA-S	16	.36	1	4	7	3	1
Earth and Space Science	4	.31	0	2	2	0	0
Life Science	6	.42	0	0	3	3	0
Physical Science	6	.30	1	2	2	0	1
Investigating	3	.23	1	1	0	0	1
Sensemaking	9	.37	0	1	5	3	0
Critiquing	4	.33	0	2	2	0	0

Table 3.1.15: Grade 12 Difficulty Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	p<.25	.25≤p<.50	.50≤p<.75	.75≤p<.90	p>=.90
NJSSA–S	10	.50	0	5	5	0	0
Earth and Space Science	5	.63	0	1	4	0	0
Life Science	1	.63	0	0	1	0	0
Physical Science	4	.41	0	4	0	0	0
Investigating	4	.56	0	2	2	0	0
Sensemaking	5	.53	0	2	3	0	0
Critiquing	1	.38	0	1	0	0	0

Table 3.1.16: Grade 12 Difficulty Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	p<.25	.25≤p<.50	.50≤p<.75	.75≤p<.90	p>=.90
NJSSA–S	15	.50	0	7	8	0	0
Earth and Space Science	5	.36	0	3	2	0	0
Life Science	6	.56	0	1	5	0	0
Physical Science	4	.40	0	3	1	0	0
Investigating	4	.46	0	2	2	0	0
Sensemaking	6	.52	0	2	4	0	0
Critiquing	5	.41	0	3	2	0	0

Table 3.1.17: Grade 12 Discrimination Indices by Domain/Practice, MC

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<= RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA–S	10	.37	2	2	4	2	0
Earth and Space Science	5	.38	0	0	3	2	0
Life Science	1	.37	0	0	1	0	0
Physical Science	4	.20	2	2	0	0	0
Investigating	4	.31	1	1	1	1	0
Sensemaking	5	.37	1	0	3	1	0
Critiquing	1	.29	0	1	0	0	0

Table 3.1.18: Grade 12 Discrimination Indices by Domain/Practice, TE

Domain/Practice	# Items	Median	RPB <.20	.20<= RPB <.30	.30<= RPB <.40	.40<= RPB <.50	RPB >=.50
NJSSA–S	15	.44	1	1	4	4	5
Earth and Space Science	5	.48	0	0	1	3	1
Life Science	6	.38	0	1	2	1	2
Physical Science	4	0.46	1	0	1	0	2
Investigating	4	.51	0	0	1	1	2
Sensemaking	6	.41	1	0	2	1	2
Critiquing	5	.41	0	1	1	2	1

3.1.2 Speededness

The consequence of time limits on examinees' scores is called speededness. A test is "speeded" to the degree that those taking the test score lower than they would have had the test not been timed. A measure of the speededness of a test is the number of items that were not attempted by students. In each separately timed subsection of a test, if a student does not attempt the last item, it can be assumed that the student may have run out of time. The percentage of students omitting an item provides information about speededness, although it must be kept in mind that students can omit an item for reasons other than speededness (for example, choosing to not put effort into answering an item). Thus, if the percentage of omits is low, that implies that there is little speededness; if the percentage of omits is high, speededness, as well as other factors, may be the cause.

NJSSA–S was not designed to be a speeded test, but rather a power test. That is, all students are expected to have ample time to finish all items and prompts. NJSSA–S assessments were administered during a testing window with 60 minutes of testing time at each grade level. Students were assumed to have enough time to complete the test.

That assumption was tested by calculating the percentage of students omitting the last item on the test. As shown in Table 3.1.19, each grade level had less than 0.50% of students omitting the last item. This is clear evidence of the NJSSA–S being a non-speeded power test at all grade levels.

Table 3.1.19: Percentage of Students Omitting the Last Item

Grade	Location	%
6	25	0.45
9	24	0.45
12	25	0.22

3.1.3 Operational DIF Analysis

The *Standards* define Differential Item Functioning (DIF) as "when different groups of test-takers with similar overall ability, or similar status on an appropriate criterion, have, on average, systematically different responses to a particular item" (AERA, APA, NCME, 2014, p. 16). If items are performing differently for sub-groups of students, the test might disadvantage some groups of students over others.

Different methods are used for DIF detection depending on whether the item is dichotomous or polytomous. For dichotomous items, DIF was identified using the Mantel-Haenszel (Mantel & Haenszel, 1959) procedure in conjunction with the ETS classification system (Dorans & Holland, 1993; Zieky, 1993). The Mantel-Haenszel (MH) method is a non-parametric approach to DIF. The ETS categorization is applied to flag the significance of DIF effects (Dorans & Holland, 1993). The letters A, B, and C are used to denote the ETS categorizations. A-level indicates negligible

DIF, B-level indicates moderate DIF, and C-level indicates severe DIF and requires a careful review of the item for possible biases. For polytomous 0–2-point TE items, DIF was identified using the Liu-Agresti cumulative common log-odds ratio (Penfield, 2007). The Liu-Agresti method allows for the ETS categorization to be applied to polytomous items.

DIF detection for the NJSSA–S operational test focused on seven comparisons of students.

- Male/Female
- White/Black
- White/Hispanic
- White/Asian
- Non-English learner (EL-No)/English learner (EL-Yes)
- Students with disabilities (SWD-Yes)/students without disabilities (SWD-No)
- Not economically disadvantaged (EconDis-No)/economically disadvantaged (EconDis-Yes)

The results of the DIF analyses were positive except for three items classified as “C.” There were no C-DIF classifications for grade 6. However, B-DIF was identified for one male-to-female comparison at grade 6. At grade 9, B-DIF was identified for one students-with-disabilities comparison and two male-to-female comparisons. Additionally, two grade 9 TE items showed C-DIF for the non-English-learner-to-English-learner comparison. At grade 12, the non-English learner-to-English-learner comparison yielded five B-DIF items, and the male-to-female comparison yielded two B-DIF items. The only grade 12 item showing C-DIF was a TE item for the White-to-Asian comparison. Tables 3.1.20 through 3.1.22 show the DIF classifications for all seven comparison groups by item type for grades 6, 9, and 12, respectively.

Table 3.1.20: Grade 6 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
6	Male/Female	MC	7	0	0
6	Male/Female	TE	17	1	0
6	Male/Female	Total	24	1	0
6	White/Black	MC	7	0	0
6	White/Black	TE	18	0	0
6	White/Black	Total	25	0	0
6	White/Hispanic	MC	7	0	0
6	White/Hispanic	TE	18	0	0
6	White/Hispanic	Total	25	0	0
6	White/Asian	MC	7	0	0
6	White/Asian	TE	18	0	0
6	White/Asian	Total	25	0	0
6	EL-No/EL-Yes	MC	7	0	0
6	EL-No/EL-Yes	TE	18	0	0
6	EL-No/EL-Yes	Total	25	0	0
6	SWD-No/SWD-Yes	MC	7	0	0
6	SWD-No/SWD-Yes	TE	18	0	0
6	SWD-No/SWD-Yes	Total	25	0	0
6	EconDis-No/EconDis-Yes	MC	7	0	0
6	EconDis-No/EconDis-Yes	TE	18	0	0
6	EconDis-No/EconDis-Yes	Total	25	0	0

Table 3.1.21: Grade 9 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
9	Male/Female	MC	8	0	0
9	Male/Female	TE	14	2	0
9	Male/Female	Total	22	2	0
9	White/Black	MC	8	0	0
9	White/Black	TE	16	0	0
9	White/Black	Total	24	0	0
9	White/Hispanic	MC	8	0	0
9	White/Hispanic	TE	16	0	0
9	White/Hispanic	Total	24	0	0
9	White/Asian	MC	8	0	0
9	White/Asian	TE	16	0	0
9	White/Asian	Total	24	0	0
9	EL-No/EL-Yes	MC	8	0	0
9	EL-No/EL-Yes	TE	14	0	2
9	EL-No/EL-Yes	Total	24	0	0
9	SWD-No/SWD-Yes	MC	8	0	0
9	SWD-No/SWD-Yes	TE	15	1	0
9	SWD-No/SWD-Yes	Total	23	1	0
9	EconDis-No/EconDis-Yes	MC	8	0	0
9	EconDis-No/EconDis-Yes	TE	16	0	0
9	EconDis-No/EconDis-Yes	Total	24	0	0

Table 3.1.22: Grade 12 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
12	Male/Female	MC	9	1	0
12	Male/Female	TE	14	1	0
12	Male/Female	Total	23	2	0
12	White/Black	MC	10	0	0
12	White/Black	TE	15	0	0
12	White/Black	Total	25	0	0
12	White/Hispanic	MC	10	0	0
12	White/Hispanic	TE	15	0	0
12	White/Hispanic	Total	25	0	0
12	White/Asian	MC	10	0	0
12	White/Asian	TE	14	0	1
12	White/Asian	Total	24	0	1
12	EL-No/EL-Yes	MC	10	0	0
12	EL-No/EL-Yes	TE	10	5	0
12	EL-No/EL-Yes	Total	20	5	0
12	SWD-No/SWD-Yes	MC	10	0	0
12	SWD-No/SWD-Yes	TE	15	0	0
12	SWD-No/SWD-Yes	Total	25	0	0
12	EconDis-No/EconDis-Yes	MC	10	0	0
12	EconDis-No/EconDis-Yes	TE	15	0	0
12	EconDis-No/EconDis-Yes	Total	25	0	0

3.2 Item Response Theory

The grade-specific NJSSA–S student ability estimates are calibrated via Item Response Theory (IRT) statistical processes. This section explains how IRT is used in the context of the NJSSA–S. Additionally, the concept of IRT is explained along with the reasoning as to why it improves upon classical test theory. Then, the specific IRT model used for the NJSSA–S is described in conjunction with the assumptions that the model must meet in order to be applicable. The remainder of this section evaluates how well the assumptions of IRT were met.

IRT is conceptualized as a family of mathematical models that explain the relationship of student performance on test items to student latent ability level on the construct of interest (Hambleton & Swaminathan, 1985). While latent abilities (e.g., anxiety, intelligence, or mastery of the NJSLS–S) are not directly observable, student responses to items measuring these abilities are directly observable. IRT models presume that the directly observable item responses of examinees can be explained by an unobservable latent trait. Within the context of the NJSSA–S, the directly observable behaviors are the responses of students to the test items, and the latent trait that we are assuming those items estimate is student understanding of the New Jersey science curriculum: the NJSLS–S.

The logic behind making and meticulously checking these assumptions is that IRT addresses many of the limitations of classical test theory (CTT) and can improve both the construction and usage of tests (Hambleton & van der Linden, 1982); hence, IRT can improve the validity of the inferences made from tests. The CTT item statistics that were presented in Section 3.1 are sample-dependent, which means that they are susceptible to substantial changes depending on the students who are answering the items. The CTT test reliability statistics presented later in Section 5.1 are susceptible to sample dependency and can increase or decrease depending on the sample’s heterogeneity. Moreover, CTT reliability is also the same for all examinees, which means that the consistency of students’ test performance is assumed to be the same regardless of their ability level.

IRT overcomes these shortcomings because its item difficulty parameters are independent of the students who took the test and its student ability estimates are independent of the test items (Hambleton & Swaminathan, 1985). If the assumptions of IRT are met, this allows students taking the NJSSA–S in the future to be placed on the same scale as the students who are taking it today. The ability to place students on the same scale allows for more meaningful year-to-year and form-to-form comparisons than CTT can offer. Moreover, unlike CTT, the reliability of IRT student ability estimates is different across the student ability spectrum as conceptualized by the test information function (TIF; see Section 5.2.1 for a more detailed explanation). This allows for test construction to be targeted to specific places on the student ability spectrum where decisions are most important in order to maximize the test’s ability to reliably classify examinees.

The increased power of IRT in comparison to CTT comes at a cost. IRT requires that certain assumptions be met. When the assumptions of IRT are not met, the data and the resulting test scores will be questionable, harming any interpretations of the test scores. Thus, it is imperative that assumptions be checked.

The NJSSA–S was constructed to meet the assumptions of a specific IRT model: the Rasch-based (Rasch, 1960) Partial Credit Model (PCM; Masters, 1982). The Rasch family of IRT models is a special case of other IRT models; Rasch models all assume that items discriminate equally and that guessing on items is minimal (Hambleton & Swaminathan, 1985). The PCM is a flexible, Rasch-based model that can be used with both dichotomous and polytomous item response data (Ostini & Nering, 2010). As described earlier, the NJSSA–S item types are designed to minimize guessing, and the test contains polytomous items (e.g., 0–2-point TE items). If the assumptions of the PCM are met, it is an appropriate and useful IRT model to use for the NJSSA–S.

The main assumptions of the PCM, as they apply to the NJSSA–S, are that the test is unidimensional, the items discriminate relatively equally, guessing on items is minimal, the response to each individual item is independent of the others, and the resulting item parameter estimates are invariant regardless of who answered the items. Each of these five major IRT assumptions will be explained in greater detail in the sections below as they relate to the PCM. The final component within this section shows disaggregated descriptive statistics of the raw scores. Overall, the results of the 2022 NJSSA–S indicate that the assumptions of the PCM were adequately met.

3.2.1 Unidimensionality

Unidimensionality was checked via multiple methods. First, the intercorrelations among the subscores were evaluated. High correlations would indicate a strong linear relationship among the subscore variables, providing evidence of unidimensionality. Second, the eigenvalues of the principal components analysis (PCA) were evaluated. A dominant first eigenvalue, in comparison to the other eigenvalues, is evidence of unidimensionality. Overall, there is ample evidence that the NJSSA–S is a unidimensional test and that the PCM assumption of unidimensionality has been met.

3.2.1.1 Intercorrelations. The Pearson product-moment correlations among the domains and practices are presented in Tables 3.2.1 through 3.2.2. High correlations would be evidence of a unidimensional test. Generally, more items in a cluster will lead to a higher correlation between that cluster and the total score. Furthermore, because each item is aligned to both a domain and a practice, the domain-to-domain and practice-to-practice intercorrelations will often be lower than the domain-to-practice and practice-to-domain intercorrelations.

At grade 6, all domains and practices correlated with the total test score at or above 0.85. Relatively high correlations between the domains or practices and the total test score were also

present at both grades 9 and 12. The lowest correlation between any subscore and the total test score was with Critiquing at grade 9. The intercorrelations among subscores provide strong evidence that the NJSSA–S is a unidimensional test.

Table 3.2.1: Correlation Matrix for Domains by Grade

Grade	Content	NJSSA–S	Earth and Space	Life	Physical
6	- Earth and Space	.87	1	-	-
	- Life	.86	.63	1	-
	- Physical	.92	.70	.70	1
9	- Earth and Space	.86	1	-	-
	- Life	.85	.61	1	-
	- Physical	.86	.58	.59	1
12	- Earth and Space	.91	1	-	-
	- Life	.84	.65	1	-
	- Physical	.81	.59	.54	1

Table 3.2.2: Correlation Matrix for Practices by Grade

Grade	Practice	NJSSA–S	Investigating	Sensemaking	Critiquing
6	- Investigating	.92	1	-	-
	- Sensemaking	.91	.74	1	-
	- Critiquing	.85	.70	.68	1
9	- Investigating	.84	1	-	-
	- Sensemaking	.89	.59	1	-
	- Critiquing	.78	.54	.58	1
12	- Investigating	.87	1	-	-
	- Sensemaking	.91	.66	1	-
	- Critiquing	.82	.61	.63	1

3.2.1.2 Principal Component Analysis. Principal Components Analysis (PCA) is a data reduction technique that attempts to account for the variance in measures by converting them into uncorrelated principal components (Brown, 2006). The first principal component accounts for as much measured variance as possible, and each succeeding factor does the same until there are as many principal components as original variables (Gorsuch, 1983). The resulting principal components can then be plotted and interpreted in a scree plot.

The results of each grade’s PCA provided further evidence of the unidimensionality of the NJSSA–S. The scree plots were interpreted by finding the place on the plot where the slope leveled off. Gorsuch (1983) noted that this method of interpretation works well when sample sizes are large and the factors are well-defined. The principal components to the left of the point on the plot where the slope leveled off were deemed practically significant. Each grade’s scree plot shows that only one major dimension is practically contributing to the variability in

student responses to items. The second most prominent eigenvalue for each grade level is close to 1, whereas the most prominent eigenvalues range from approximately 5–7.

Grade 6 Scree Plot

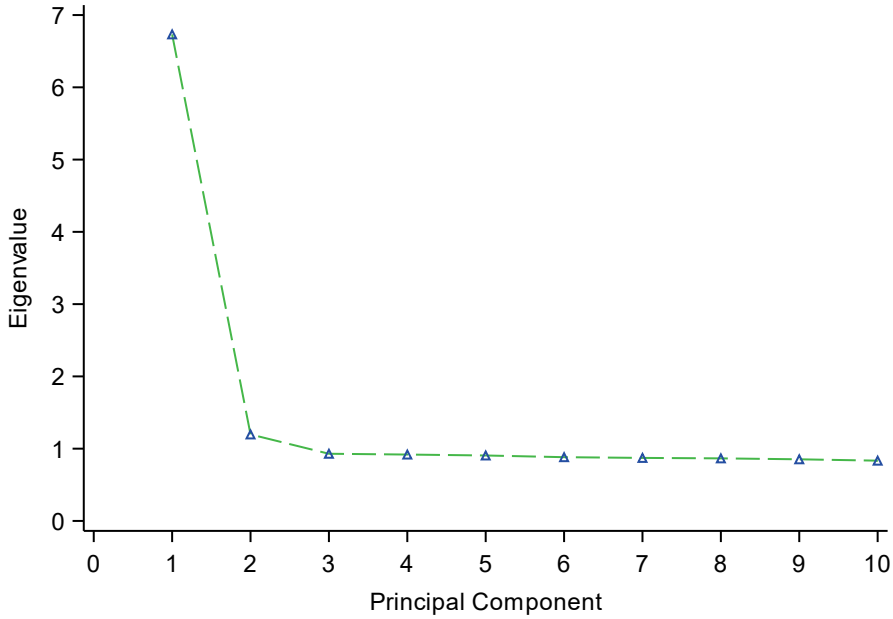


Figure 3.2.1. Grade 6 Scree Plot

Grade 9 Scree Plot

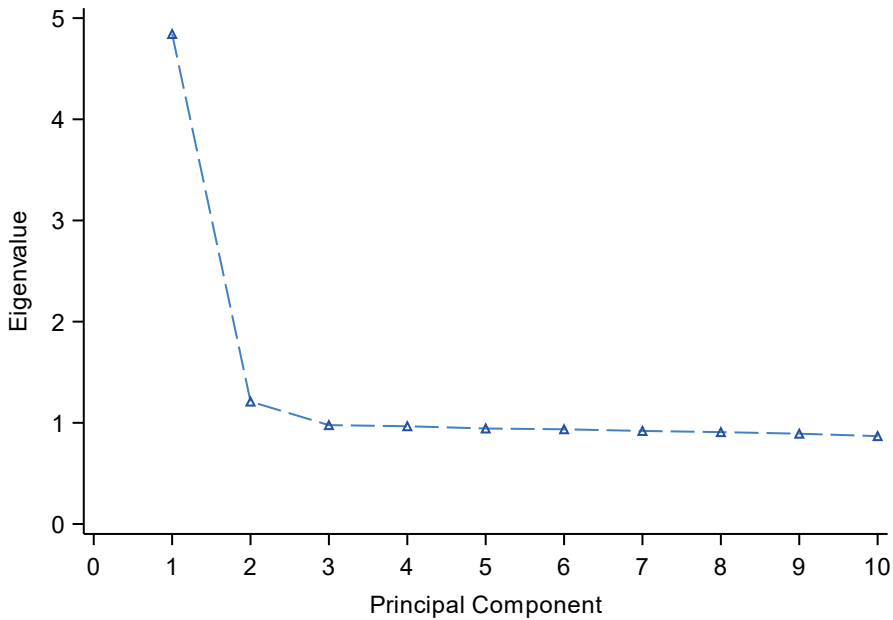


Figure 3.2.2. Grade 9 Scree Plot

Grade 12 Scree Plot

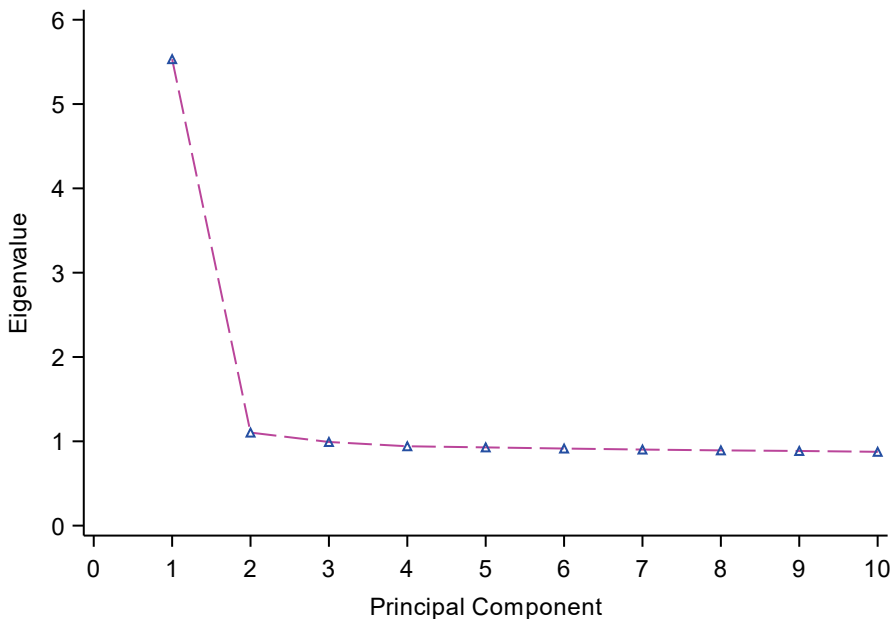


Figure 3.2.3. Grade 12 Scree Plot

3.2.2 Partial-Credit Model-Fit Statistics

Hambleton et al. (1991) noted that “[a] poorly fitting IRT model will not yield invariant item and ability parameters” (p. 53), which diminishes the beneficial properties inherent to IRT. PCM model fit was assessed at the item level via Rasch-based item infit and outfit, discrimination, and guessing statistics. At the person level, model fit was evaluated using Rasch-based person-infit and -outfit statistics. These statistics were calculated during the 2022 NJSSA–S IRT calibration processes via Winsteps 3.74 (Linacre, 2012). Detailed item parameter estimates and model-fit statistics are presented in [Appendix C](#). Overall, there is ample evidence that the items at all grades fit the assumptions of the PCM, as is described in the following sections.

3.2.2.1 Item infit and outfit. Rasch infit and outfit statistics range from zero to infinity, with one representing ideal model fit. For the NJSSA–S, items were flagged for having infit or outfit statistics outside of the 0.7 to 1.3 range (Wright and Linacre, 1994). Infit statistics are influenced by unexpected responses from students on items that are measuring near their ability level (Wright and Masters, 1982). Only one item across all grades was flagged for problematic infit statistics.

Outfit statistics are heavily influenced by unexpected student responses to items that are either relatively easy or relatively hard. The NJSSA–S outfit statistics resulted in more flagged items than the infit statistics: six grade 9 items were flagged, four grade 6 items were flagged, and three grade 12 items were flagged. Flagged outfit statistics are less of a threat to the validity of test score interpretations than problematic infit statistics, especially when the flagged items’ outfit statistics are only slightly outside the flagging thresholds, as was the case for all flagged

items. Thus, while there is clearly room for improvement regarding the item outfit, the infit and outfit statistics provide reasonable evidence that the assumptions of the PCM have been met. Table 3.2.3 provides a summary of item infit and outfit statistics at each grade level.

Table 3.2.3: Summary Item Infit and Outfit Statistics

Grade	Fit Statistics	Mean	Min	Max	Outside 0.7 to 1.3	% Flagged
6	Infit	0.97	0.77	1.23	0 out of 25	0.0%
6	Outfit	1.00	0.72	1.41	4 out of 25	16.0%
9	Infit	0.98	0.72	1.18	0 out of 24	0.0%
9	Outfit	1.00	0.59	1.52	6 out of 24	25.0%
12	Infit	1.01	0.82	1.36	1 out of 25	4.0%
12	Outfit	1.00	0.72	1.49	3 out of 25	12.0%

3.2.2.2 Rasch discrimination. The PCM assumes that all items discriminate equally. Practically, items never discriminate equally, but if they are within reasonable thresholds then the assumption holds. The assumption of equal discrimination can be tested with the Rasch discrimination statistic as well as the correlations presented earlier in the CTT section. Rasch discrimination statistics are centered at 1.0, which indicates that the item is discriminating exactly as expected by the PCM. Items are flagged when their discrimination statistics fall outside of the range of 0.5 to 1.5.

At each grade, the Rasch discrimination statistics looked outstanding. Only one grade 6 and two grade 12 items were flagged for having a value outside the 0.5 to 1.5 threshold. Table 3.2.4 provides a summary of the Rasch discrimination statistics at each grade level.

Table 3.2.4: Summary Rasch Discrimination Statistics

Grade	Fit Statistics	Mean	Min	Max	Outside 0.5 to 1.5	% Flagged
6	Discrimination	1.04	0.49	1.44	1 out of 25	4.0%
9	Discrimination	1.04	0.51	1.49	0 out of 24	0.0%
12	Discrimination	0.98	0.07	1.38	2 out of 25	8.0%

3.2.2.3 Rasch lower asymptote. The PCM assumes that there is minimal guessing on the test items. Practically, however, students guess, and sometimes they guess correctly. Thus, as with the assumption of equal discrimination, the guessing assumption is met if items remain within a reasonable threshold. The assumption of guessing can be tested with the Rasch lower asymptote statistics. The Rasch lower asymptote statistics are ideally 0.0, which indicates that an item is displaying little to no guessing. Items are flagged when their lower asymptote statistics fall outside of the range of 0.0 to 0.1.

At all grades, the lower asymptote statistics met the model assumptions. Only five items out of a total of 74 ($\approx 7\%$) were flagged for having values above the 0.1 threshold. The average lower asymptote statistics at each grade were close to the ideal value of 0.0. Grade 12 had the most items flagged with 3 out of 25. Table 3.2.5 provides a summary of the lower asymptote statistics at each grade level.

Table 3.2.5: Summary Rasch Lower Asymptote Statistics

Grade	Fit Statistics	Mean	Min	Max	Greater Than 0.1	% Flagged
6	Lower Asymptote	0.02	0.00	0.18	1 out of 25	4.0%
9	Lower Asymptote	0.03	0.00	0.11	1 out of 24	4.2%
12	Lower Asymptote	0.04	0.00	0.23	3 out of 25	12.0%

3.2.2.4 Rasch person infit and outfit. PCM person-fit statistics are useful for evaluating whether student response patterns are expected. Anomalous response patterns include response patterns that are either improbable or too probable given the underlying model. Multiple factors can cause distortions in the expected patterns of test scores, including:

- Carelessness — examinees miss items that they should have answered correctly.
- Cheating — examinees receive information to correctly answer items that they would have normally missed.
- Guessing — examinees correctly answer items without knowing the correct answer.
- Creative responses — examinees misinterpret the item.
- Test administration errors.

Two measures of PCM person-fit statistics were used: infit and outfit. Person infit is more influenced by responses to items that are targeted at the person’s ability level; outfit is more influenced by responses to items that are relatively easy or hard for a student given their ability level (Wright & Masters, 1982).

Person-fit statistics were evaluated based on the following demographics: gender, ethnicity, English learner (EL) status, economically disadvantaged (EconDis) status, students with disabilities (SWD) status, and all major forms. Figures 3.2.4 through 3.2.6 exhibit grade-level distributions of both the person-infit and -outfit statistics for all students. Tables 3.2.6 and 3.2.7 show person-infit and -outfit descriptive statistics by demographic variables. Tables 3.2.8 and 3.2.9 break down the person-infit and -outfit descriptive statistics by Computer-Based Testing (CBT), Paper-Based Testing (PBT), and Spanish forms.

Overall, few students were flagged for aberrant person-infit statistics at grades 6 and 12. Grade 9 showed higher numbers of students flagged at certain demographic variables. Less than 6% of students were flagged for person-infit statistics at all demographic variables at grade 6. At grade 12, less than 1% of 88,236 students were flagged. At grade 9, 11.61% of students were flagged across all demographic variables. Within demographic variables, 15% or more EL,

economically disadvantaged, black, Hispanic, and students with a disability status were flagged for aberrant infit statistics. As shown in Table 3.2.9, no form type at grades 6 or 12 had more than 5% of students flagged. Grade 9 showed more infit flags by form type, with 11.30% of CBT forms being flagged, 32.88% of PBT forms being flagged, and 27.30% of the SP forms being flagged.

The person-outfit statistics resulted in more students being flagged. At both grades 6 and 9, large percentages of students were flagged for outfit across several demographic variables and form types. Across all demographic variables, grade 6 showed 28.72% of students flagged, grade 9 showed 25.52% of students flagged, and grade 12 showed 9.20% of students flagged for aberrant outfit statistics. As stated earlier, aberrant person-infit statistics are more of a threat to the validity of the inferences made from test scores than are aberrant person-outfit statistics.

Table 3.2.6: Person-Infit Statistics by Demographic Group

Grade	Group	N	Mean Raw Score	Person-Infit Mean	Person-Infit Min	Person-Infit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
6	NJSSA–S	95,881	13.51	0.98	0.57	1.65	4,273	4.46	16.19
6	Male	48,951	13.56	0.99	0.57	1.65	2,215	4.52	16.15
6	Female	46,908	13.46	0.97	0.57	1.63	2,056	4.38	16.23
6	Am. Indian	164	13.62	0.98	0.63	1.41	6	3.66	16.83
6	Asian	10,324	17.65	0.98	0.57	1.63	565	5.47	18.21
6	Black	13,428	10.71	0.98	0.57	1.59	457	3.40	14.44
6	Hispanic	31,135	11.30	0.98	0.57	1.65	1,213	3.90	14.44
6	Pacific Islander	190	14.66	0.99	0.64	1.36	9	4.74	13.67
6	White	37,792	15.10	0.97	0.57	1.64	1,892	5.01	17.07
6	EL – Yes	7,465	7.47	1.01	0.61	1.65	194	2.60	10.22
6	EL – No	88,410	14.02	0.98	0.57	1.64	4,078	4.61	16.47
6	EconDis – Yes	31,367	10.89	0.98	0.57	1.65	1,160	3.70	14.20
6	EconDis – No	64,504	14.78	0.98	0.57	1.64	3,112	4.82	16.93
6	SWD – Yes	19,541	10.10	0.99	0.57	1.60	689	3.53	14.18
6	SWD – No	76,334	14.38	0.97	0.57	1.65	3,583	4.69	16.58
9	NJSSA–S	99,303	9.94	1.00	0.58	2.70	11,530	11.61	6.88
9	Male	50,525	10.19	1.01	0.58	2.70	6,273	12.42	6.89
9	Female	48,653	9.67	0.99	0.58	2.70	5,242	10.77	6.87
9	Am. Indian	169	10.02	1.01	0.64	2.00	18	10.65	8.33
9	Asian	10,609	13.63	0.96	0.58	2.70	790	7.45	12.62
9	Black	14,265	7.84	1.03	0.58	2.47	2,279	15.98	5.47
9	Hispanic	32,176	8.11	1.03	0.58	2.70	4,919	15.29	5.50
9	Pacific Islander	183	10.50	0.99	0.58	1.54	19	10.38	5.58
9	White	39,419	11.11	0.97	0.58	2.70	3,298	8.37	8.42
9	EL – Yes	6,062	5.48	1.14	0.58	2.70	1,683	27.76	4.26
9	EL – No	93,240	10.23	0.99	0.58	2.70	9,847	10.56	7.33

Grade	Group	N	Mean Raw Score	Person-Infit Mean	Person-Infit Min	Person-Infit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
9	EconDis – Yes	29,728	7.90	1.04	0.58	2.70	4,807	16.17	5.42
9	EconDis – No	69,574	10.81	0.98	0.58	2.70	6,723	9.66	7.93
9	SWD – Yes	18,671	7.75	1.06	0.58	2.70	3,570	19.12	5.20
9	SWD – No	80,630	10.44	0.98	0.58	2.70	7,960	9.87	7.64
12	NJSSA–S	88,236	12.26	1.01	0.64	1.53	702	0.80	11.46
12	Male	44,648	12.25	1.00	0.65	1.50	287	0.64	11.33
12	Female	43,498	12.26	1.02	0.64	1.53	415	0.95	11.55
12	Am. Indian	110	12.08	1.02	0.78	1.30	0	0.00	N/A
12	Asian	9,732	15.96	0.99	0.65	1.50	48	0.49	12.63
12	Black	11,482	9.83	1.02	0.65	1.49	106	0.92	10.95
12	Hispanic	26,030	10.33	1.02	0.64	1.53	224	0.86	10.75
12	Pacific Islander	196	13.06	1.01	0.66	1.43	3	1.53	11.33
12	White	38,973	13.29	1.01	0.65	1.50	309	0.79	11.97
12	EL – Yes	4,187	7.39	1.02	0.69	1.43	22	0.53	9.23
12	EL – No	84,049	12.50	1.01	0.64	1.53	680	0.81	11.53
12	EconDis – Yes	23,477	10.26	1.02	0.65	1.52	220	0.94	10.54
12	EconDis – No	64,758	12.98	1.01	0.64	1.53	482	0.74	11.88
12	SWD – Yes	16,986	10.00	1.02	0.65	1.52	138	0.81	10.49
12	SWD – No	71,247	12.80	1.01	0.64	1.53	564	0.79	11.70

Table 3.2.7: Person-Outfit Statistics by Demographic Group

Grade	Group	N	Mean Raw Score	Person-Outfit Mean	Person-Outfit Min	Person-Outfit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
6	NJSSA-S	95,881	13.51	1.00	0.23	9.73	27,541	28.72	13.81
6	Male	48,951	13.56	1.01	0.23	6.21	14,538	29.70	13.81
6	Female	46,908	13.46	0.98	0.23	9.73	12,992	27.70	13.80
6	Am. Indian	164	13.62	0.98	0.36	2.64	47	28.66	14.32
6	Asian	10,324	17.65	0.98	0.23	5.69	3,380	32.74	19.26
6	Black	13,428	10.71	1.02	0.23	9.73	3,818	28.43	9.65
6	Hispanic	31,135	11.30	1.02	0.23	9.73	8,782	28.21	10.45
6	Pacific Islander	190	14.66	1.03	0.35	2.64	58	30.53	15.31
6	White	37,792	15.10	0.98	0.23	6.21	10,589	28.02	16.18
6	EL – Yes	7,465	7.47	1.09	0.37	6.21	2,551	34.17	5.79
6	EL – No	88,410	14.02	0.99	0.23	9.73	24,986	28.26	14.62
6	EconDis – Yes	31,367	10.89	1.02	0.23	9.73	8,816	28.11	9.85
6	EconDis – No	64,504	14.78	0.99	0.23	9.73	18,719	29.02	15.67
6	SWD – Yes	19,541	10.10	1.05	0.23	9.73	6,143	31.44	9.06
6	SWD – No	76,334	14.38	0.99	0.23	6.21	21,393	28.03	15.17
9	NJSSA-S	99,303	9.94	1.00	0.12	8.41	25,343	25.52	8.16
9	Male	50,525	10.19	1.01	0.12	8.41	13,011	25.75	8.32
9	Female	48,653	9.67	0.98	0.12	8.19	12,301	25.28	7.98
9	Am. Indian	169	10.02	1.00	0.27	3.23	37	21.89	8.35
9	Asian	10,609	13.63	0.91	0.12	8.41	2,241	21.12	13.69
9	Black	14,265	7.84	1.07	0.12	8.19	4,379	30.70	6.12
9	Hispanic	32,176	8.11	1.06	0.12	8.19	9,546	29.67	6.29
9	Pacific Islander	183	10.50	0.97	0.20	2.91	50	27.32	8.84
9	White	39,419	11.11	0.95	0.12	8.19	8,523	21.62	9.70
9	EL – Yes	6,062	5.48	1.23	0.12	8.19	2,516	41.50	4.55
9	EL – No	93,240	10.23	0.98	0.12	8.41	22,826	24.48	8.56

Grade	Group	N	Mean Raw Score	Person-Outfit Mean	Person-Outfit Min	Person-Outfit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
9	EconDis – Yes	29,728	7.90	1.07	0.12	8.19	9,167	30.84	6.13
9	EconDis – No	69,574	10.81	0.97	0.12	8.41	16,175	23.25	9.31
9	SWD – Yes	18,671	7.75	1.11	0.12	8.19	6,218	33.30	5.96
9	SWD – No	80,630	10.44	0.97	0.12	8.41	19,123	23.72	8.87
12	NJSSA–S	88,236	12.26	1.00	0.40	3.63	8,116	9.20	12.91
12	Male	44,648	12.25	0.99	0.40	3.63	4,268	9.56	13.20
12	Female	43,498	12.26	1.01	0.40	3.63	3,843	8.83	12.58
12	Am. Indian	110	12.08	1.01	0.50	1.52	8	7.27	11.50
12	Asian	9,732	15.96	0.96	0.41	3.63	1,274	13.09	19.27
12	Black	11,482	9.83	1.02	0.40	3.62	1,023	8.91	7.93
12	Hispanic	26,030	10.33	1.01	0.40	3.63	2,112	8.11	8.75
12	Pacific Islander	196	13.06	1.00	0.49	3.31	20	10.20	14.50
12	White	38,973	13.29	1.00	0.40	3.31	3,512	9.01	14.48
12	EL – Yes	4,187	7.39	1.02	0.43	2.64	384	9.17	4.99
12	EL – No	84,049	12.50	1.00	0.40	3.63	7,732	9.20	13.30
12	EconDis – Yes	23,477	10.26	1.01	0.40	3.31	1,964	8.37	8.62
12	EconDis – No	64,758	12.98	1.00	0.40	3.63	6,152	9.50	14.28
12	SWD – Yes	16,986	10.00	1.02	0.40	3.62	1,643	9.67	8.72
12	SWD – No	71,247	12.80	1.00	0.40	3.63	6,472	9.08	13.97

Table 3.2.8: Person-Infit Statistics by Form

Grade	Group	N	Mean Raw Score	Person-Infit Mean	Person-Infit Min	Person-Infit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
6	CBT	93,918	13.64	0.98	0.57	1.65	4,233	4.51	16.24
6	PBT	111	10.12	1.05	0.77	1.45	4	3.60	8.75
6	SP	1,852	7.36	1.00	0.63	1.44	36	1.94	10.61
9	CBT	97,417	10.02	1.00	0.58	2.70	11,011	11.30	7.00
9	PBT	73	6.30	1.17	0.75	2.70	24	32.88	3.88
9	SP	1,813	5.58	1.14	0.58	1.99	495	27.30	4.35
12	CBT	87,081	12.32	1.01	0.64	1.53	689	0.79	11.49
12	PBT	71	9.18	1.02	0.81	1.25	0	0.00	N/A
12	SP	1,084	7.66	1.05	0.73	1.43	13	1.20	9.69

Table 3.2.9: Person-Outfit Statistics by Form

Grade	Group	N	Mean Raw Score	Person-Outfit Mean	Person-Outfit Min	Person-Outfit Max	Flagged N	Flagged Percent	Flagged Mean Raw Score
6	CBT	93,918	13.64	1.00	0.23	9.73	26,913	28.66	13.99
6	PBT	111	10.12	1.15	0.23	3.39	35	31.53	10.80
6	SP	1,852	7.36	1.06	0.37	6.21	593	32.02	5.48
9	CBT	97,417	10.02	0.99	0.12	8.41	24,592	25.24	8.27
9	PBT	73	6.30	1.34	0.66	3.92	29	39.73	3.97
9	SP	1,813	5.58	1.22	0.12	8.19	722	39.82	4.57
12	CBT	87,081	12.32	1.00	0.40	3.63	7,982	9.17	13.03
12	PBT	71	9.18	1.04	0.49	2.04	7	9.86	9.14
12	SP	1,084	7.66	1.07	0.43	1.89	127	11.72	5.54

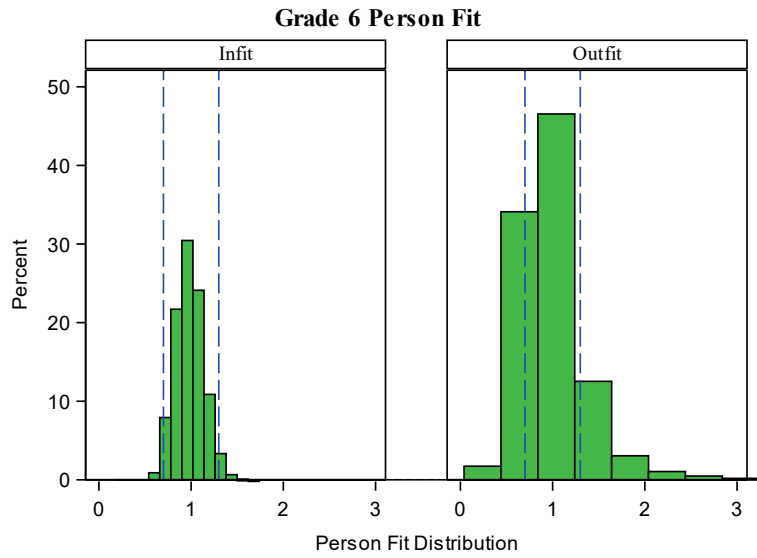


Figure 3.2.4. Grade 6 Person-Infit and -Outfit Distributions

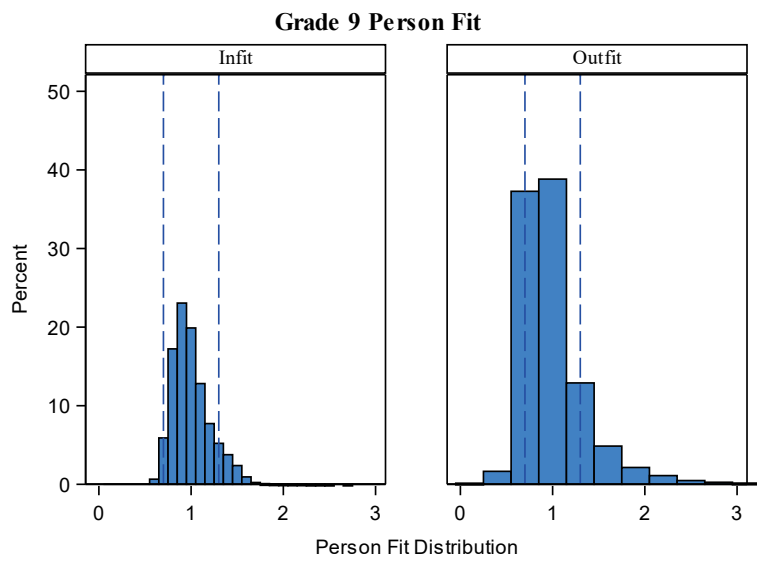


Figure 3.2.5. Grade 9 Person-Infit and -Outfit Distributions

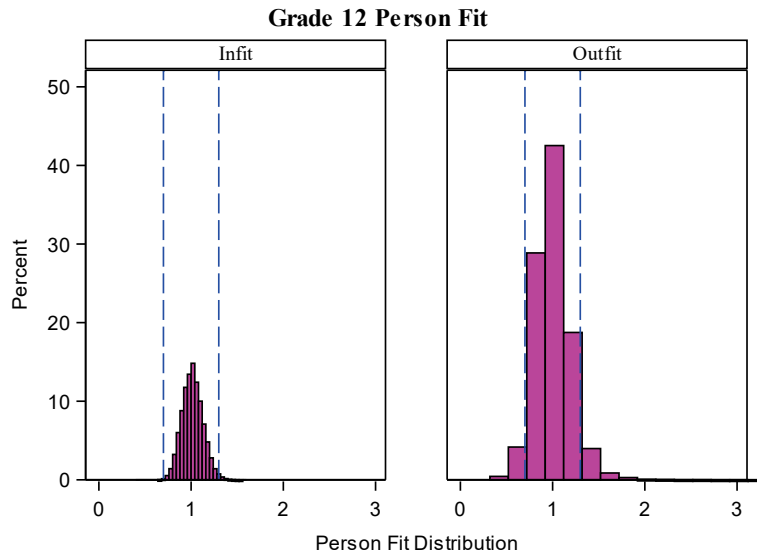


Figure 3.2.6. Grade 12 Person-Infit and -Outfit Distributions

3.2.3 Local Independence

The PCM (Masters, 1982) assumes that student responses to items are independent of responses to other items. In other words, the model assumes that student performance on one item does not affect performance on the other test items. This assumption is called the assumption of local independence. A violation of the assumption of local independence could increase the likelihood that the reliability of the assessment is overestimated and/or the item-total correlations are inflated; either outcome poses a threat to the validity of inferences made from test scores.

The assumption of local independence was tested via calculations of Yen's (1984) Q3, which is a residual correlation. All combinations of items were checked, and they were flagged if their Q3 value was greater than 0.2 or less than -0.2 (Chen & Thissen, 1997). The results at all grades indicate that the assumption of local independence was met because none of 876 residual correlations between items displayed a Q3 value outside the acceptable threshold. Table 3.2.10 summarizes Yen's Q3 statistics at each grade level.

Table 3.2.10: Summary of Yen's Q3 Statistics

Grade	Mean	Min	Max	Outside-0.2 to 0.2	% Flagged
6	-.04	-.12	.20	0 out of 300	0.0%
9	-.04	-.13	.09	0 out of 276	0.0%
12	-.04	-.15	.18	0 out of 300	0.0%

3.2.4 Descriptive Statistics — Raw Score

This section contains descriptive statistics for raw scores and support-level distributions by form and links to NJDOE documentation detailing student support-level percentages by demographic group.

3.2.4.1 Raw score distributions by form. Descriptive statistics for raw scores and percentage distributions of students’ support levels are summarized by form in Table 3.2.11. For all test forms, raw scores have a range of 0 to 25. The cut scores for each support level can be found in Section 5.3.1 of this report.

Table 3.2.11: Descriptive Statistics of Raw Scores and Students’ Support Levels by Form

Grade	Form*	N	Mean	SD	Min	Max	%Strong Support	%Some Support	%Less Support
6	CBT	93,949	13.63	5.97	0	25	41.15	33.50	25.35
6	PBT	111	10.12	5.71	1	24	67.57	21.62	10.81
6	SP	1,866	7.39	4.40	0	22	85.85	12.11	2.04
9	CBT	97,419	10.02	4.85	0	25	43.56	40.94	15.50
9	PBT	73	6.30	3.59	1	17	73.97	23.29	2.74
9	SP	1,813	5.58	2.71	0	19	86.49	13.46	0.06
12	CBT	87,085	12.32	5.50	0	25	52.56	21.91	25.53
12	PBT	71	9.18	5.01	2	24	73.24	15.49	11.27
12	SP	1,085	7.66	3.26	1	19	91.52	7.65	0.83

* CBT: Computer-Based Test; PBT: Paper-Based Test; SP: Spanish

3.2.4.2 Raw score distributions by demographic group. Percentage distributions of students’ support levels by demographic groups can be found on the [New Jersey Statewide Assessment Reports webpage](#). Raw score cumulative frequency distributions are attached as [Appendix B](#) in this technical brief.

Part 4: Scale Stability

Standard 5.6 states that “[t]esting programs that attempt to maintain a common scale over time should conduct periodic checks of the stability of the scale on which scores are reported” (AERA, APA, NCME, 2014, p.103). Moreover, as described in Part 2 of this technical brief, the items had previously been administered on the 2019 NJSLA–S operational test forms, as well as the 2018 NJSLA–S field test. That means different cohorts of students in New Jersey have been tested via this set of items for five consecutive years. Thus, in order to ensure the comparability of NJSSA–S test scores, the stability of the underlying item difficulty parameters warranted examination.

Kolen and Brennan (2004) recommend that scale stability be inspected both statistically and visually. The methods used for testing the scale stability of the NJSSA–S follow this recommendation. The first method is the Delta Plot method as described by Angoff (1982). The second method is the 0.3 Logit Absolute Difference criterion as described by Miller et al. (2004). The former is a CTT-based method for detecting discrepancies between group item means and p-values. The latter is a Rasch, IRT-based method for detecting item parameter drift.

The following sections describe both methods and present the results. Both the results of the Delta Plot and 0.3 Logit Absolute Difference methods were generally positive. Moreover, the person- and item-fit statistics presented in Sections 3.2.2.1 and 3.2.2.4 provide similarly positive evidence of scale stability. Overall, the NJSSA–S scale appears to have remained stable in 2022.

4.1 Delta Plot Method

The Delta Plot method — also referred to as the Transformed Item Difficulty Index — was introduced by Angoff (1972). It was originally conceptualized as a method for identifying biased test items within the framework of CTT. Per the recommendations of the New Jersey Technical Advisory Committee, the Delta Plot method was added to the scale stability analysis for the purpose of complementing the IRT methods, including the fit statistics described in Section 3.2 and the 0.3 Logits Absolute Difference Method described in Section 4.2.

The Delta Plot method uses multiple transformations to place p-values onto the “delta” scale, which is a common scale used by Educational Testing Service (ETS) that has a mean of 13 and a standard deviation of 4. Once the item p-values have been transformed, it is typical to plot the values onto a scatter plot and then create a trendline (Camilli & Shepard, 1994). The perpendicular distance (PD) of each item from the trendline is calculated, and items are flagged if their perpendicular distance is two standard deviation units away from the trendline.

The results of the Delta Plot method were generally positive. Table 4.1.1 shows a summary of the results, including the percentage of items that were flagged at each grade level. The “Flagging Distance” column shows the PD at which items were flagged by grade. Figures 4.1.1

to 4.1.3 show the delta plots for grades 6, 9, and 12, respectively. Only one item was flagged for grade 6, while grades 9 and 12 both had two items that were flagged. The flagged grade 6 item had comparable P-values across years and a PD slightly larger than the threshold of 0.43. Both items in grades 9 and 12 were more challenging than anticipated, given that students generally did worse on the items on the 2022 NJSSA–S in comparison to the students who took them on the 2019 NJSLA–S. [Appendix D](#) contains the detailed results of the Delta Plot method for each grade level. It should be noted that item means for 0–2-point TE items were converted to adjusted p-values (i.e., item mean divided by max score point) for running the Delta Plot method analyses.

Table 4.1.1: Summary of Delta Plot Method

Grade	Total Items	Flagged Items	Percent Flagged	Flagging Distance
6	25	1	4.00%	0.430
9	24	2	8.33%	0.528
12	25	2	8.00%	0.455

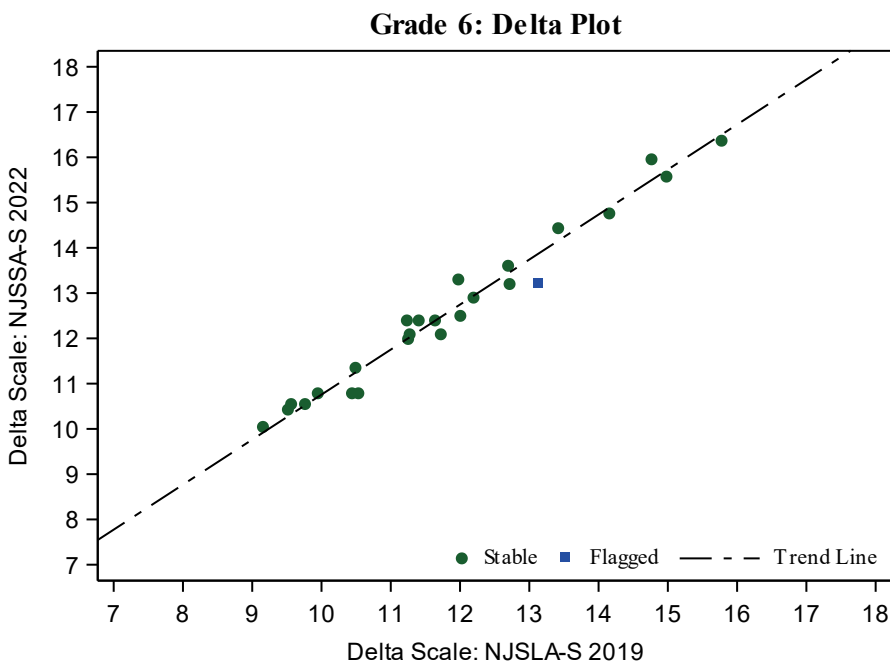


Figure 4.1.1. Grade 6 Delta Plot

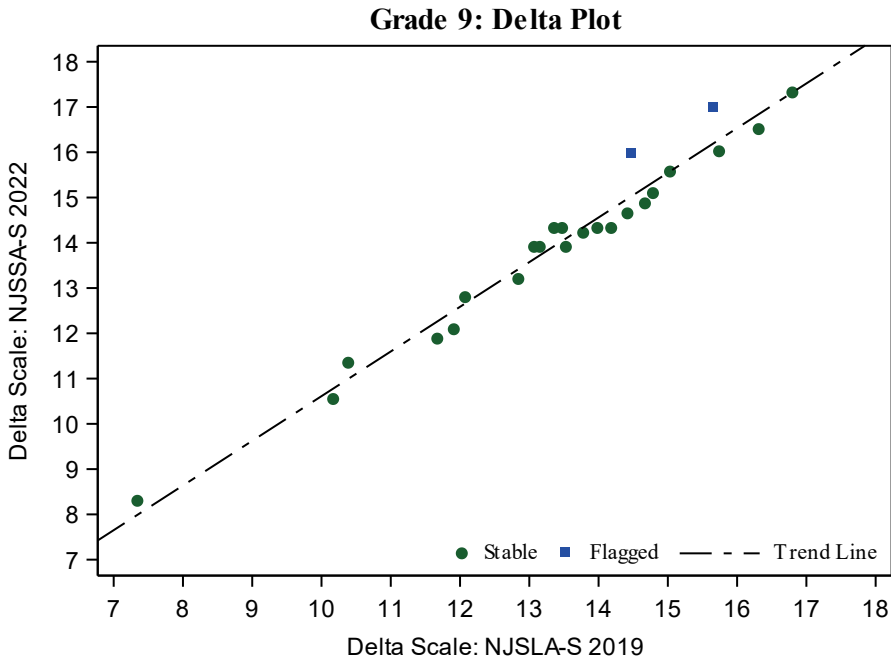


Figure 4.1.2. Grade 9 Delta Plot

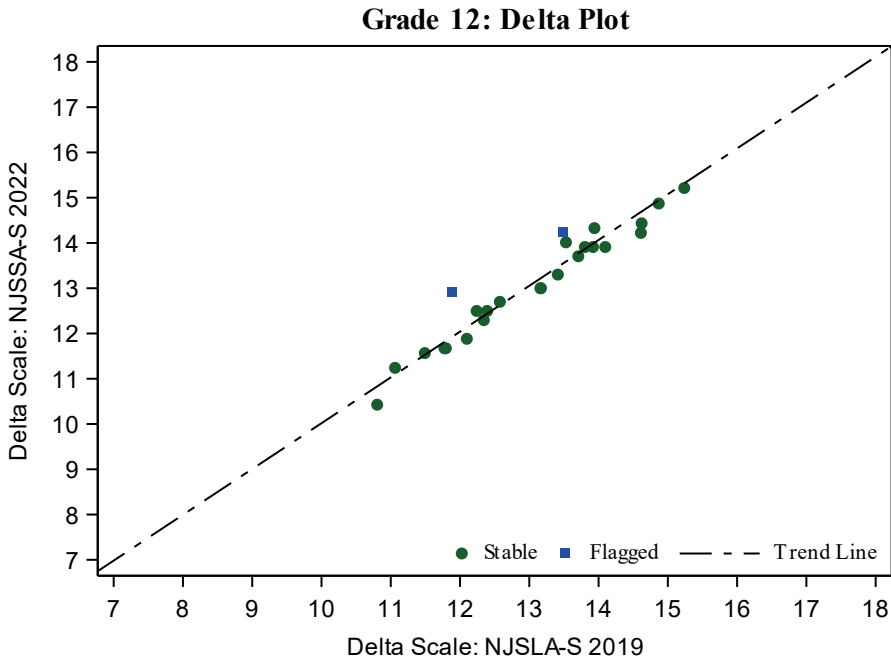


Figure 4.1.3. Grade 12 Delta Plot

4.2 0.3 Logits Absolute Difference Method

The 0.3 Logits Absolute Difference method was used to assess the stability of Rasch-based item difficulty parameters. The first step in the process was to recalibrate the NJSSA–S item difficulty parameters via Winsteps 3.74 (Linacre, 2012) using the 2022 test results. This calibration will hereafter be referenced as the unconstrained calibration, indicating that during the calibration process, the item difficulty parameters were allowed to freely converge regardless of previous results.

The second step was to use the 0.3 Logits Absolute Difference method to compare the results of the unconstrained calibration to the item difficulty parameters from the 2019 NJSLA–S. The latter serves as the basis for the NJSSA–S scale. To do so, an equating constant was calculated that represents the mean difference between the unconstrained item difficulty parameters and their 2019 NJSLA–S counterparts. Each unconstrained item difficulty parameter was then adjusted by adding to it the equating constant. If the absolute difference between an item’s adjusted, unconstrained item difficulty parameter and its 2019 NJSLA–S counterpart was greater than 0.3 logits, then the item with the greatest absolute difference was flagged. The process iterates until all items have adjusted, unconstrained item parameters within 0.3 logits of the constrained item difficulty parameters. Large percentages of flagged items would indicate that the item difficulty parameters were not stable from the 2019 NJSLA–S to the 2022 NJSSA–S administration.

The results of the 0.3 Logits Absolute Difference method were generally positive at all grade levels. Table 4.2.1 provides a summary of the results at each grade including the percentage of flagged items. Figures 4.2.1 to 4.2.3 are scatter plots with the constrained item difficulty parameters on the X-axis and the unconstrained item difficulty parameters on the Y-axis; items are demarcated as either stable or flagged, depending on the results of the 0.3 Logits Absolute Difference procedure. At grades 9 and 12, only two items ($\approx 8.0\%$) were flagged for both grades. They were also flagged by the delta plot method. At grade 6, five items (20%) were flagged. However, all of the flags were only slightly above the 0.3 Logits Absolute Difference threshold. Moreover, as was described in the previous section, only one item was also flagged by the delta plot method. For grade 9, the flagged items were due to the items being more difficult for the students than the model would have predicted, while for the two other grades, the pattern was indeterminate. [Appendix D](#) contains the detailed results of the 0.3 Logits Absolute Difference for each grade level.

Table 4.2.1: Summary 0.3 Logits Absolute Difference Method

Grade	Total Items	Flagged Items	Percent Flagged	Average Constrained Item Difficulty	Average Unconstrained Item Difficulty	Equating Constant
6	25	5	20.0%	–.286	.000	–.286
9	24	2	8.3%	–.151	.000	–.151
12	25	2	8.0%	–.282	.000	–.282

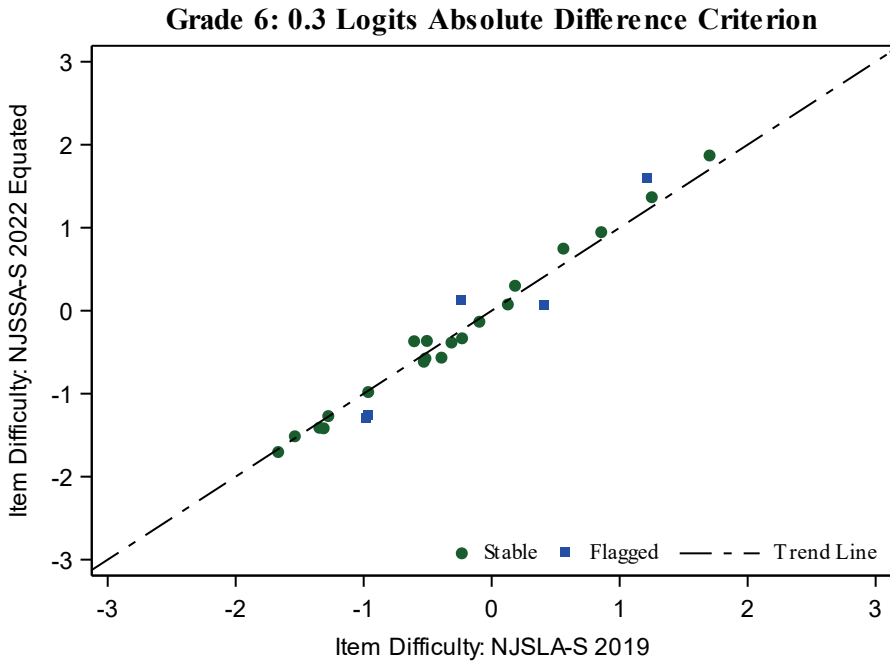


Figure 4.2.1. Grade 6 0.3 Logits Absolute Difference Criterion Item Difficulty Plot

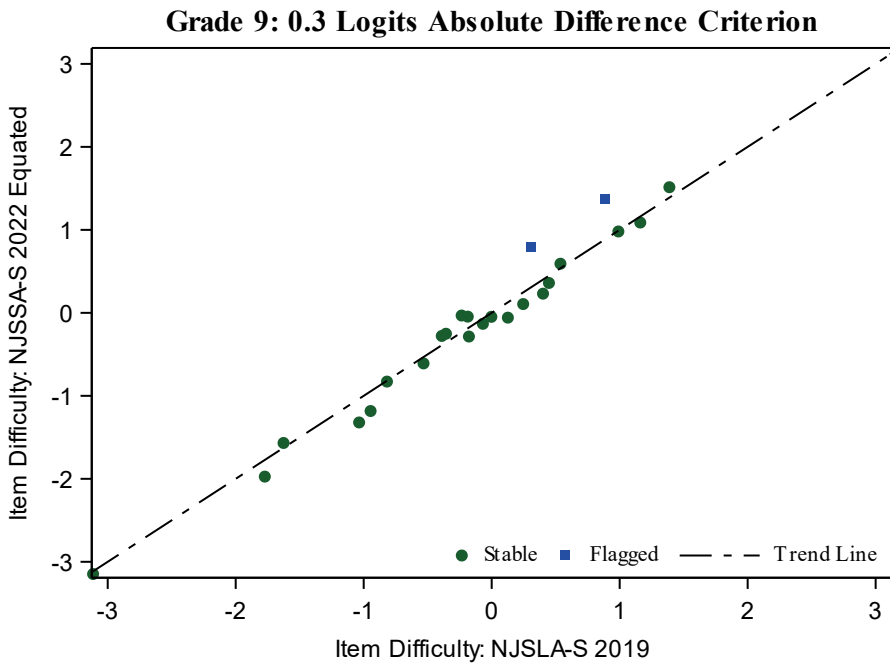


Figure 4.2.2. Grade 9 0.3 Logits Absolute Difference Criterion Item Difficulty Plot

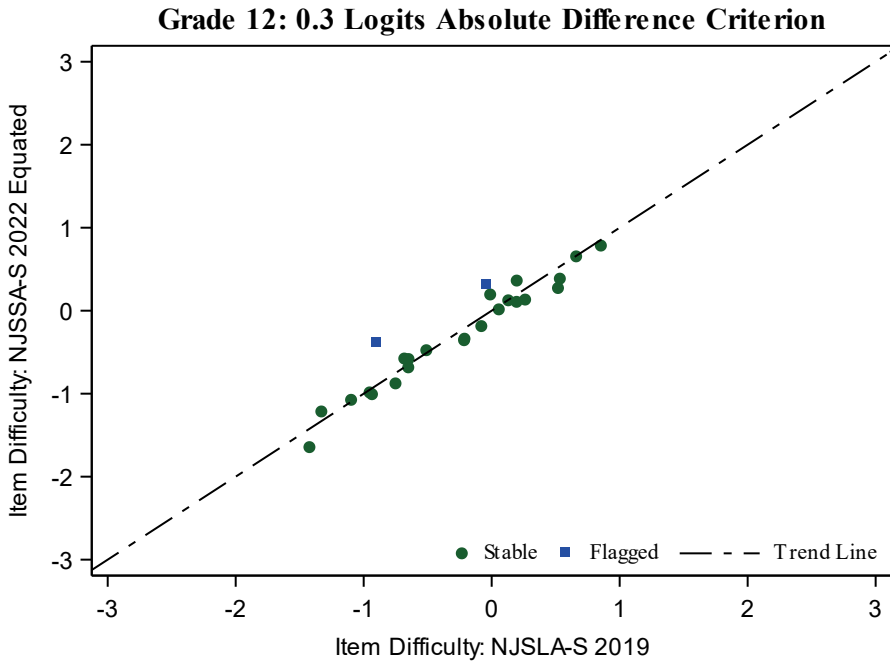


Figure 4.2.3. Grade 12 0.3 Logits Absolute Difference Criterion Item Difficulty Plot

Part 5: Reliability

Test reliability refers to the consistency of test scores. Ultimately, valid interpretations of test scores are dependent upon those scores being reliable. *Standard 2.0* states that “[a]ppropriate evidence of reliability/precision should be provided for the interpretation for each intended score use” (AERA, APA, NCME, 2014, p. 42). Examples of appropriate evidence include reliability coefficients, conditional standard errors of measurement (CSEM), test information functions, and decision consistency measures, amongst others. The following sections detail evidence supporting the reliability of the NJSSA–S test scores and subscores.

5.1 Classical Test Theory Reliability Estimates

This section describes the Classical Test Theory (CTT) reliability estimates calculated for the 2022 NJSSA–S. Section 5.1.1 describes the concept of reliability in the CTT framework, and Section 5.1.2 displays the results.

5.1.1 Reliability and Measurement Error

Student test scores are reliable when measurement error is minimized. Increasing reliability by minimizing measurement error is an important goal in the construction of any test. Under the assumptions of CTT, any observed measurement — such as a test score, X — is defined as a composite of true score, T , and its associated error:

$$X = T + \text{error} \quad \text{Equation 5.1}$$

Estimating the size of the measurement error associated with the true score is the key to estimating reliability. Errors in measurement can result from any of several factors, including environmental factors (e.g., testing conditions) and examinee factors (e.g., fatigue, stress). CTT provides a means for this quantification of examinee inconsistency (i.e., measurement error).

The definitions or assumptions in CTT lead to several important properties. For example, it can be demonstrated that

$$\sigma_x^2 = \sigma_t^2 + \sigma_e^2, \quad \text{Equation 5.2}$$

or observed score variance (σ_x^2) equals the sum of true score variance (σ_t^2) and error variance (σ_e^2). The relationships among the variance terms (i.e., σ_x^2 , σ_t^2 , σ_e^2) are critical to a more thorough understanding of important CTT concepts, including reliability and the standard error of measurement. For example, CTT reliability (ρ) is defined as the correlation between observed scores (x_1 , x_2) on parallel forms, which is equal to true score variance (σ_t^2) divided by observed score variance (σ_x^2):

$$\rho_{x_1x_2} = \sigma_t^2 / \sigma_x^2. \quad \text{Equation 5.3}$$

With just a few algebraic steps, the CTT definition of the standard error of measurement (SEM, σ_e) can be shown as:

$$\sigma_e = \sigma_x \sqrt{1 - \rho_{x_1 x_2}}. \quad \text{Equation 5.4}$$

Although the concepts of reliability and SEM are relatively straightforward, issues underlying the estimation of reliability are not. Reliability can be estimated via the correlation of scores on parallel forms or from test-retest data, or it can be estimated from a single test administration using any one of a variety of techniques (e.g., Brown, 1910; Cronbach, 1951; Kuder & Richardson, 1937).

For NJSSA–S, the consistency of individual student performance was estimated using Cronbach’s (1951) coefficient alpha. Coefficient alpha is conceptualized as the proportion of total raw score variance that may be attributed to a student’s true score variance. Ideally, more score variance should be attributable to true test scores than to measurement error.

Separate analyses were performed for each grade level. Scores from all item types were used in the computations. Coefficient alpha was estimated using the following formula:

$$\alpha_{\text{Cronbach}} = \frac{n}{n-1} \left[1 - \frac{\sum_{i=1}^n \sigma_{Y_i}^2}{\sigma_X^2} \right], \quad \text{Equation 5.5}$$

where n is the number of items, $\sigma_{Y_i}^2$ is the variance of item i , and σ_X^2 is the variance of observed total score, X . SEMs were calculated using the following formula:

$$SEM = S_X \sqrt{1 - \alpha_{\text{Cronbach}}}, \quad \text{Equation 5.6}$$

where S_X is the standard deviation of observed total scores.

5.1.2 Raw Score Internal Consistency

In order to accommodate the state’s diverse testing population, the NJSSA–S was delivered in multiple formats. Most students received the online computer-based test (CBT), the Spanish (SP) test, or the paper-based test (PBT). It should be noted that internal consistency reliability estimates, such as coefficient alpha, decrease when the students taking a given test form are more homogeneous in their test performance (i.e., do not show much variance in their test scores).

Table 5.1.1 displays the coefficient alpha and SEM for each form by grade. Overall, the reliability coefficients at each grade level indicate that students' raw scores were reliable. The results at grade 6 stand out as particularly exceptional given that there are only 25 points on the NJSSA–S. The grade 6 reliability coefficients ranged from .79 for the Spanish form to .88 for the CBT form. The most likely reason for the better results at grade 6, which had the same test length as the other grades, is that the grade 6 items were closely matched to the ability levels of the grade 6 students, thereby increasing the variance among test scores. At grade 9, where the distribution of test scores was heavily skewed toward the low end of the ability spectrum, reliability ranged from .49 on the Spanish form to .81 on the CBT form. Similarly, the grade 12 alpha coefficients ranged from .56 on the Spanish form to .84 on the CBT form. The relatively low-reliability measures for the Spanish forms were likely due to those populations showing limited variance in the Spanish form test scores.

Table 5.1.1: Coefficient Alpha and SEM, by Form

Grade	Form*	N-Count	Mean	SD	Alpha	SEM
6	CBT	93,918	13.64	5.97	.88	2.07
6	PBT	111	10.12	5.71	.86	2.10
6	SP	1,852	7.36	4.39	.79	2.02
9	CBT	97,417	10.02	4.85	.81	2.10
9	PBT	73	6.30	3.59	.69	2.00
9	SP	1,813	5.58	2.71	.49	1.94
12	CBT	87,081	12.32	5.50	.84	2.21
12	PBT	71	9.18	5.01	.81	2.16
12	SP	1,084	7.66	3.26	.56	2.18

* CBT: Computer-Based Test; PBT: Paper-Based Test; SP: Spanish

Table 5.1.2 summarizes the coefficient alpha and SEMs of the six reporting categories by grade. It should be noted that reliability coefficients are commonly low when based on small numbers of items (Traub & Rowley, 2008). Thus, reporting categories such as Critiquing and Investigating, which had fewer items, tended to have lower reliability measures. The highest subscore reliability of .81 was for Physical Science at grade 6; whereas the lowest subscore reliability of .46 was for Critiquing at grade 9. The reliability measures in Table 5.1.2 are based on all test takers at each grade level.

Table 5.1.2: Coefficient Alpha and SEM by Reporting Category

Grade	Reporting Category	Total Items	MC Items	TE Items	TE2 Items	Max Points	Alpha	SEM
6	Total	25	7	18	0	25	.88	2.07
6	Earth and Space Science	8	4	4	0	8	.64	1.21
6	Life Science	7	1	6	0	7	.67	1.12
6	Physical Science	10	2	8	0	10	.81	1.23
6	Investigating	9	2	7	0	9	.75	1.27
6	Sensemaking	10	4	6	0	10	.72	1.30
6	Critiquing	6	1	5	0	6	.66	0.98
9	Total	24	8	15	1	25	.81	2.10
9	Earth and Space Science	8	4	4	0	8	.60	1.21
9	Life Science	7	1	6	0	7	.63	1.06
9	Physical Science	9	3	5	1	10	.57	1.34
9	Investigating	9	6	3	0	9	.58	1.26
9	Sensemaking	10	1	9	0	10	.69	1.32
9	Critiquing	5	1	3	1	6	.46	1.01
12	Total	25	10	15	0	25	.84	2.21
12	Earth and Space Science	10	5	5	0	10	.74	1.35
12	Life Science	7	1	6	0	7	.63	1.15
12	Physical Science	8	4	4	0	8	.51	1.31
12	Investigating	8	4	4	0	8	.63	1.26
12	Sensemaking	11	5	6	0	11	.68	1.49
12	Critiquing	6	1	5	0	6	.57	1.04

Table 5.1.3 shows the coefficient alpha and SEMs by demographic group. These calculations are based on the entire test. In general, the coefficient alphas are consistently high among the various demographic groups. At grade 6, the lowest value was .79, for English learner (EL) students, which is still very strong, given that the NJSSA–S only consisted of 25 points worth of items. At grade 9, the coefficient alphas were close to the .70 to .80 range, except for the EL students ($\alpha_{EL-Yes} = .50$). The same pattern was evident at grade 12, where all the coefficient alphas hovered close to the .75 to .85 range, except for EL students ($\alpha_{EL-Yes} = .58$)

Table 5.1.3: Coefficient Alpha and SEM by Demographic Group

Grade	Group	N	Mean	SD	Alpha	SEM
6	NJSSA–S	95,881	13.51	6.00	.88	2.07
6	Male	48,951	13.56	6.21	.89	2.05
6	Female	46,908	13.46	5.78	.87	2.08
6	Am. Indian	164	13.62	5.81	.87	2.08
6	Asian	10,324	17.65	4.96	.85	1.93
6	Black	13,428	10.71	5.66	.86	2.09
6	Hispanic	31,135	11.30	5.70	.86	2.10
6	Pacific Islander	190	14.66	5.67	.87	2.08
6	White	37,792	15.10	5.44	.86	2.06
6	EL–Yes	7,465	7.47	4.43	.79	2.03
6	EL–No	88,410	14.02	5.84	.87	2.07
6	EconDis–Yes	31,367	10.89	5.59	.86	2.10
6	EconDis–No	64,504	14.78	5.78	.87	2.05
6	SWD–Yes	19,541	10.10	5.94	.88	2.06
6	SWD–No	76,334	14.38	5.70	.87	2.07
9	NJSSA–S	99,303	9.94	4.85	.81	2.10
9	Male	50,525	10.19	5.11	.83	2.09
9	Female	48,653	9.67	4.56	.79	2.09
9	Am. Indian	169	10.02	4.96	.82	2.11
9	Asian	10,609	13.63	4.85	.81	2.11
9	Black	14,265	7.84	3.98	.73	2.05
9	Hispanic	32,176	8.11	4.13	.75	2.06
9	Pacific Islander	183	10.50	5.01	.82	2.11
9	White	39,419	11.11	4.72	.80	2.12
9	EL–Yes	6,062	5.48	2.73	.50	1.93
9	EL–No	93,240	10.23	4.82	.81	2.10
9	EconDis–Yes	29,728	7.90	4.01	.74	2.06
9	EconDis–No	69,574	10.81	4.92	.82	2.11
9	SWD–Yes	18,671	7.75	4.35	.78	2.04
9	SWD–No	80,630	10.44	4.83	.81	2.11
12	NJSSA–S	88,236	12.26	5.50	.84	2.21
12	Male	44,648	12.25	5.73	.85	2.18
12	Female	43,498	12.26	5.25	.82	2.24
12	Am. Indian	110	12.08	5.20	.81	2.25
12	Asian	9,732	15.96	5.35	.85	2.10
12	Black	11,482	9.83	4.69	.77	2.23
12	Hispanic	26,030	10.33	4.78	.78	2.24
12	Pacific Islander	196	13.06	5.05	.80	2.25
12	White	38,973	13.29	5.41	.83	2.21
12	EL–Yes	4,187	7.39	3.31	.58	2.14
12	EL–No	84,049	12.50	5.47	.84	2.21

Grade	Group	N	Mean	SD	Alpha	SEM
12	EconDis–Yes	23,477	10.26	4.76	.78	2.24
12	EconDis–No	64,758	12.98	5.57	.84	2.20
12	SWD–Yes	16,986	10.00	5.12	.82	2.20
12	SWD–No	71,247	12.80	5.45	.83	2.21

Table 5.1.4 displays coefficient alpha and SEM by the two main item types: multiple-choice (MC) and technology-enhanced (TE). As would be expected, as the number of points associated with a specific item type increases, so does the corresponding coefficient alpha. More than half of the points available on each test were associated with TE item types; thus, it is not surprising that the TE items displayed alphas from .73 to .85, and the MC items displayed alphas from .56 to .67 for each grade level.

Table 5.1.4: Coefficient Alpha and SEM by Item Type

Grade	Item Type	Items	Points	Mean	SD	Alpha	SEM
5	MC	7	7	4.63	1.88	.67	1.08
5	TE	18	18	8.88	4.48	.85	1.76
8	MC	8	8	3.67	2.01	.62	1.23
8	TE	16	17	6.26	3.27	.73	1.69
11	MC	10	10	5.19	2.20	.56	1.45
11	TE	15	15	7.07	3.76	.80	1.66

5.2 Item Response Theory Reliability

The reliability of the scores ascertained from the Partial Credit Model (PCM; Masters, 1982) was assessed in multiple ways. Test information functions (TIFs) and item maps were evaluated at each grade level. Overall, the 2022 NJSSA–S was reliable from the perspective of IRT and the PCM.

5.2.1 Test Information Functions

In IRT, the reliability of an assessment is conceptualized via the test information function (TIF; Hambleton & Swaminathan, 1985). Unlike coefficient alpha (Cronbach, 1951), the TIF is not uniform across the entire range of test scores. Instead, the TIF can assess test reliability across the full range of scores. This is particularly important to a criterion-referenced test such as the NJSSA–S because it allows for the reliability of the assessment to be evaluated specifically at the most important decision points (i.e., the “Some Support” and “Less Support” cut scores).

The TIF consists of the summation of all the item information functions (IIF) on a given test (Hambleton, 1989; Lord & Novick, 1968). An IIF is the probability of a correct response multiplied by the probability of an incorrect response. Item information functions (I_{ij}) for every

item (j) at every level of student ability (i) can be calculated for each item using the following equation:

$$I_{ij}(\theta_i, \delta_j) = P_{ij} * (1 - P_{ij}) \quad \text{Equation 5.7}$$

The total test information function is simply the sum of all the item information functions. Thus, each item contributes to the TIF, and proper selection of items during the test construction process will lead to TIFs that maximize information at important decision points.

Figures 5.2.1 to 5.2.3 illustrate, respectively, the TIFs for grades 6, 9, and 12 at person ability estimates ranging from -5 to $+5$. More information at a specific ability level implies less measurement error. Ideally, the peak of the information function would maximize information at both the “Some Support” (Level 2) and “Less Support” (Level 3) cut scores in order to minimize measurement error where the most important decisions are taking place. Within each figure, there are two vertical, dashed lines representing the cut scores.

At grade 6 the TIF peaked close to the “Some Support” cut score. There was a large drop in information at the “Less Support” cut. At grade 9 the TIF peaked almost directly in between the “Some Support” and “Less Support” cut scores. Overall, the grade 9 TIF is very close to being ideal. The grade 12 TIF was similar to its grade 6 counterpart. It peaked almost directly at the “Some Support” cut score, and there was a relatively large decrease in information at the “Less Support” cut score. Overall, the TIFs provide ample evidence that student ability estimates are reliable at the most important decision points. However, both grades 6 and 12 could be improved to gather more information at the “Less Support” cut.

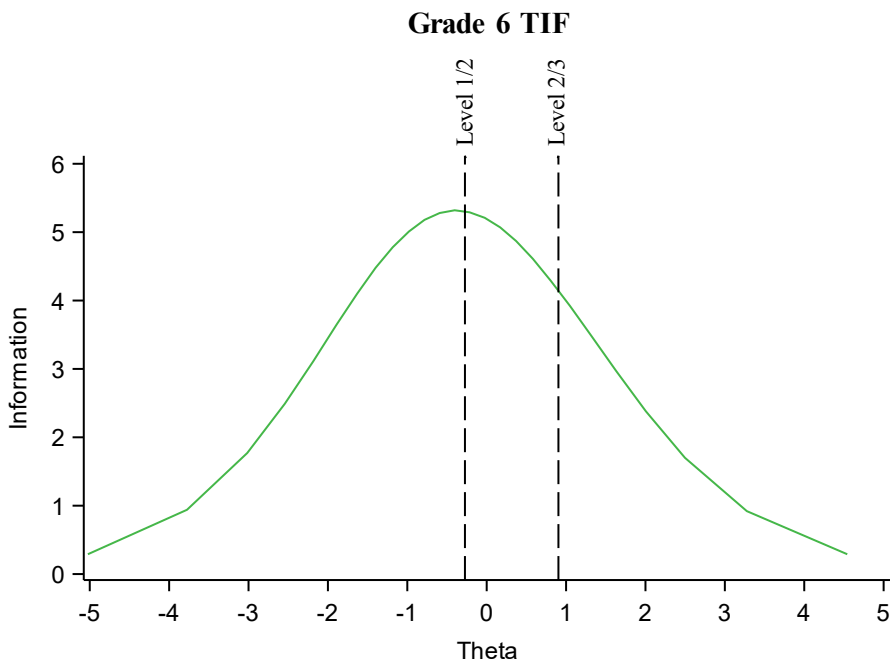


Figure 5.2.1. Grade 6 Test Information Function

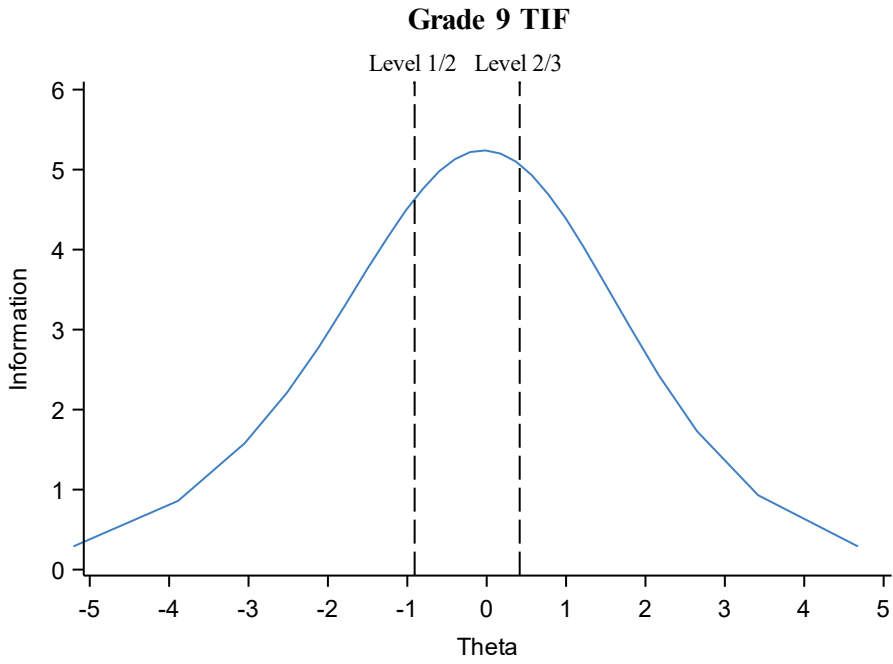


Figure 5.2.2. Grade 9 Test Information Function

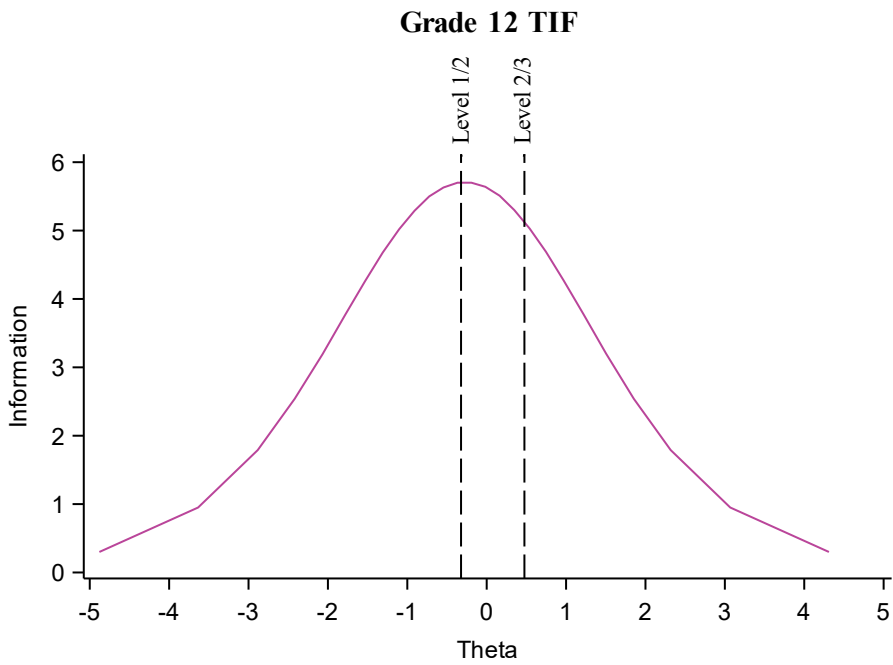


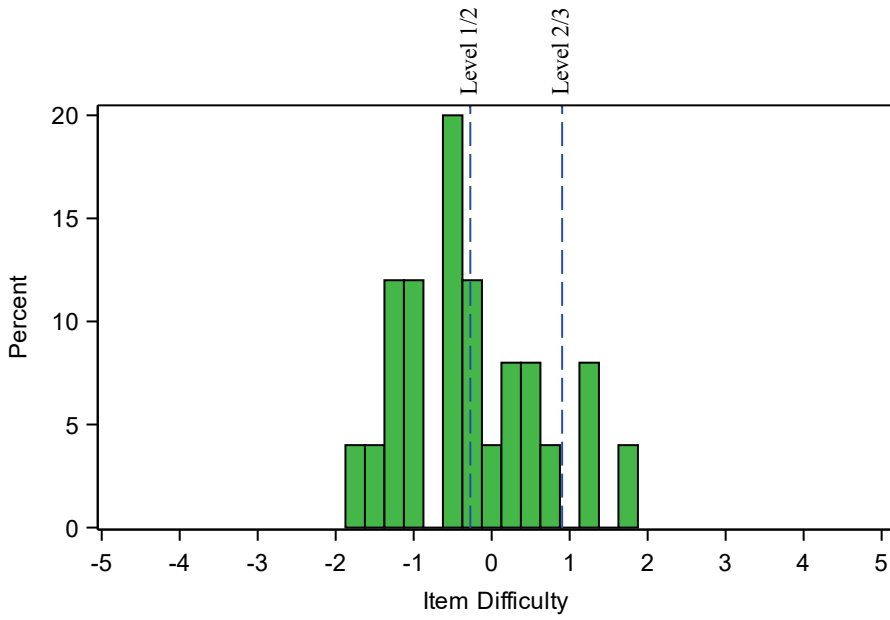
Figure 5.2.3. Grade 12 Test Information Function

5.2.2 Item Maps

Item maps indicate how well the item difficulties and person ability levels match. Items that are targeted to the ability levels of the students taking the test will result in more reliable measures of student ability. Figures 5.2.4 through 5.2.6 show the 2022 NJSSA–S item maps.

Unsurprisingly, given the TIFs, the grades 6 and 12 item difficulty distributions peak slightly below the “Some Support” (Level 1/2) cut score, while the grade 9 item difficulty distribution peaks directly in between the “Some Support” and “Less Support” cut scores. At grade 6, the theta distribution was normally distributed with student ability peaking close to the “Some Support” cut score. The item distribution peaked below the “Some Support” cut score. The theta distribution at grade 9 peaked between the “Some Support” and “Less Support” cut scores; however, it was more skewed towards the lower end of the ability spectrum, with large peaks of students below the “Some Support” cut. The grade 12 theta distribution peaked right at the “Some Support” cut score. The grade 6 item distribution contained more items below the ability levels of the students than would be ideal. The grade 9 item distribution, as the grade 9 TIF showed, matched the decision points on the scale very well. However, there were many students below the “Some Support” cut score and very few items along that part of the scale. The extremely tight grade 12 item distribution was lacking items at both the upper and lower parts of the scales in comparison to the ability levels of the students.

Grade 6 Item Difficulty



Grade 6 Ability

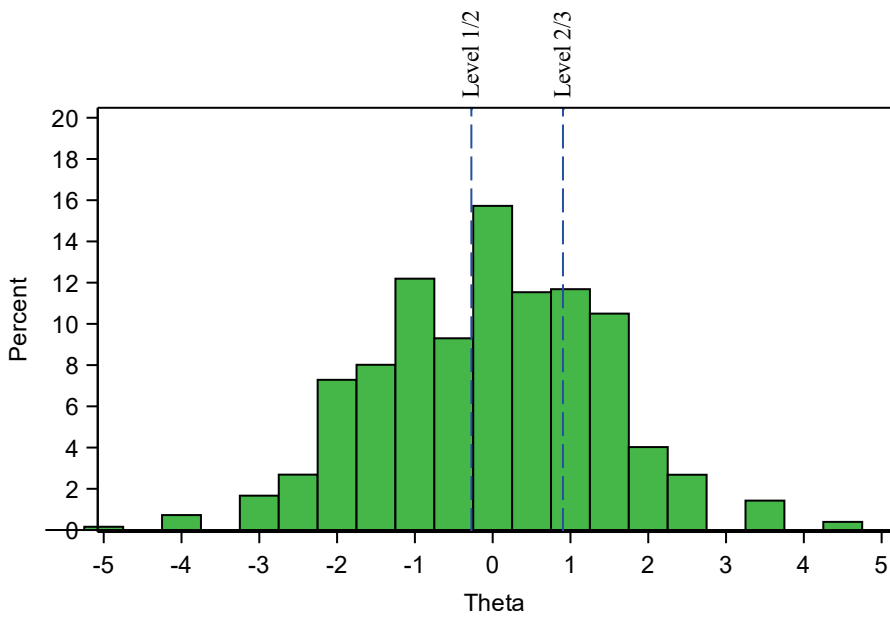


Figure 5.2.4. Grade 6 Item Difficulty and Student Ability Distributions

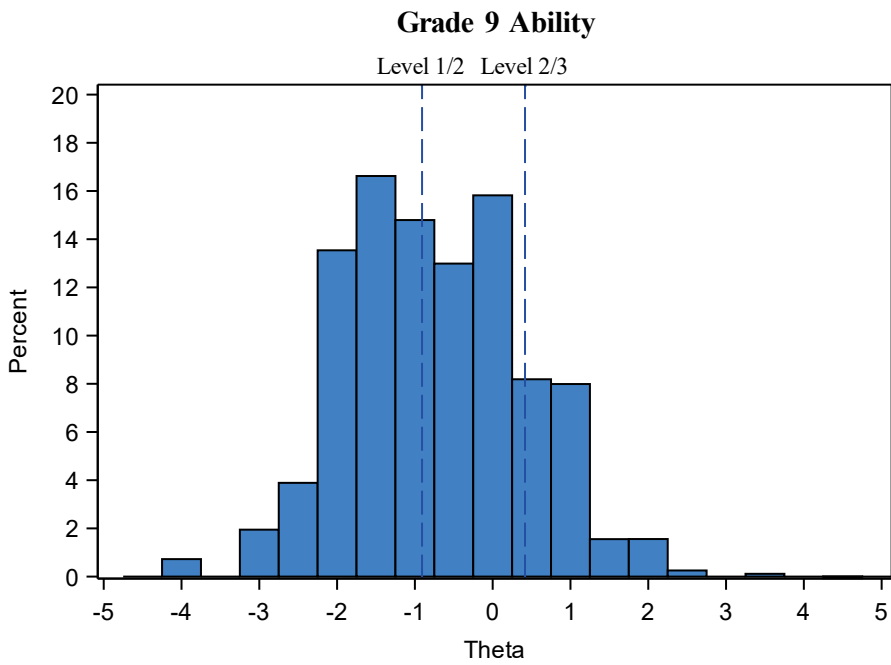
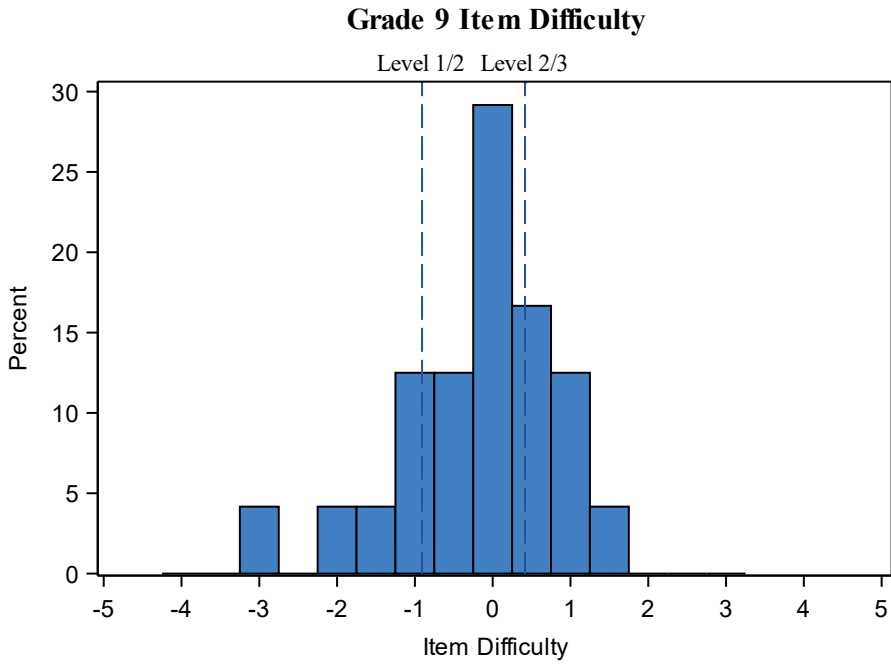


Figure 5.2.5. Grade 9 Item Difficulty and Student Ability Distributions

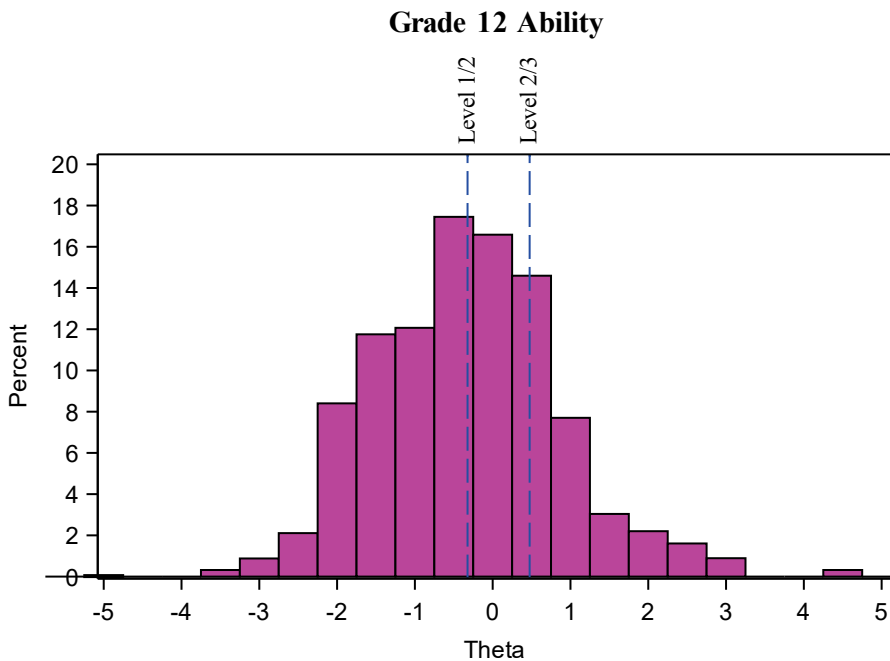
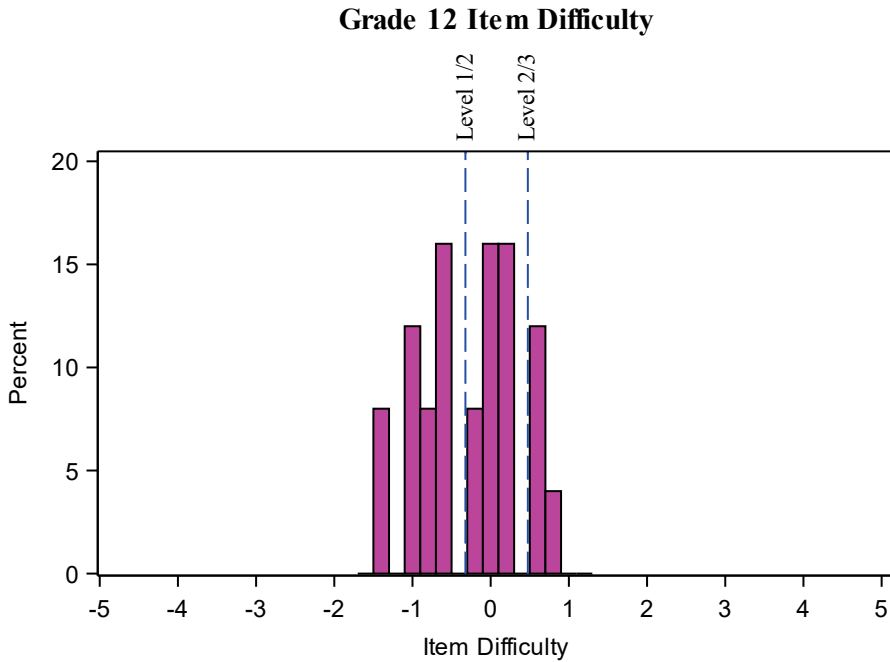


Figure 5.2.6. Grade 12 Item Difficulty and Student Ability Distributions

5.3 Reliability of Performance Classifications

The reliability of the performance level classifications was evaluated via two methods. First, error bands were placed around each cut score using the CSEM. Next, the BB-CLASS (Brennan, 2004) program was used to calculate performance-level classification consistency indices. The results of both methods indicate that the 2022 NJSSA–S performance level classifications were reliable.

5.3.1 Conditional Standard Error of Measurement at Each Cut-Score

Winsteps calculates the conditional standard error of measurement (CSEM) at each score point using the test information function. The equation for the standard error at each value of theta (ability) is given by:

$$SE(\hat{\theta}) = \frac{1}{\sqrt{I(\theta)}} \quad \text{Equation 5.8}$$

where $I(\theta)$ is the information function for a test at a score point (theta).

The 2022 NJSSA–S raw cut scores and the corresponding conditional standard error of measurement (CSEM) on the theta scale are summarized in Table 5.3.1. The theta scores and corresponding CSEMs for all raw scores are presented in [Appendix E](#). The lower and upper bound values in Table 5.3.1 have been placed on the raw score scale. Given that the CSEMs are the inverse of the TIF, their interpretations are similar. The test characteristic curves (TCCs) are graphical representations of the expected raw scores points a student would achieve at a given level of theta. The upper and lower bounds were defined by multiplying the CSEM for the theta cut score by two and either adding it to or subtracting it from the cut score’s theta value. Next, the upper and lower bound theta values were identified on the TCC to find their corresponding raw score point. Any overlap between the upper or lower bounds and one of the other cut scores could indicate reliability problems among the support level classifications. [Appendix F](#) contains both CSEM and TCC graphs for each grade level.

At grades 6 and 9, there was no overlap between the upper bound of the “Some Support” and the “Less Support” cut score, nor was there any overlap between the “Less Support” lower bound and the “Some Support” cut score. However, at grade 12, both cut scores showed overlap between the relevant upper or lower bound. The upper bound of the “Some Support” cut score was 17.4, which overlapped with the “Less Support” cut score of 17. Similarly, the “Less Support” lower bound of 12.6 overlapped with the “Some Support” cut score of 13. The overlap at grade 12 is due to the close proximity of the “Some Support” and “Less Support” cut scores on the theta scale and the relatively small number of items on the Start Strong in comparison to the NJSLA–S. If the Start Strong is continued with the same scale and cut scores, then a longer grade 12 test might be appropriate to address the overlap among the cut scores and the upper and lower bounds.

Table 5.3.1: Raw Cut Scores with Conditional Standard Error of Measurement

Grade	Level	Raw Cut score	CSEM*	Lower Bound	Upper Bound
6	Some Support	13	0.43	8.9	17.1
6	Less Support	19	0.51	14.9	23.1
9	Some Support	9	0.46	4.8	13.2
9	Less Support	16	0.45	11.8	20.2
12	Some Support	13	0.42	8.6	17.4
12	Less Support	17	0.45	12.6	21.4

*CSEM placed on the theta scale

5.3.2 Classification Consistency Indices

The reliability index for proficiency classifications (coefficient kappa; Cohen, 1960) is an estimate of how reliably the test classifies students into the support level categorizations (i.e., Strong Support, Some Support, Less Support). Kappa was computed with the BB-CLASS program (Brennan, 2004) based on the beta-binomial model. Coefficient kappa (K) is given by:

$$\kappa = \frac{\varphi - \varphi_c}{1 - \varphi_c}, \quad \text{Equation 5.9}$$

where φ is the probability of a consistent classification and φ_c is the probability of a consistent classification by chance. A classification consistency index can be regarded as the percentage of examinees that would hypothetically be assigned to the same achievement level if the same test was administered a second time or an equivalent test was administered under the same conditions.

Table 5.3.2 displays the results from BB-CLASS (Brennan, 2004) using the Livingston and Lewis (1995) approach to estimate classification consistency. At each grade level, the classification consistency ranged from .72 to .74. Thus, if the NJSSA–S had been administered a second time, approximately 72–74% of the students would have been classified at exactly the same performance level. The classification consistency is similar to the much longer NJSLA–S (NJDOE, 2020). Overall, the NJSSA–S support level classifications should be interpreted as being consistent.

Table 5.3.2: Support Level Classification Reliability

Grade	Alpha	SEM	Some Support Cut	Less Support Cut	Kappa	Classification Consistency
6	.88	2.07	13	19	.59	.74
9	.81	2.10	9	16	.54	.72
12	.84	2.21	13	17	.56	.73

Part 6: Validity

The *Standards* state that “[v]alidity is a unitary concept. It is the degree to which all the accumulated evidence supports the intended interpretation of test scores for the proposed use” (AERA, APA, NCME, 2014, p. 14). If there is ample evidence to support reasonable interpretations and test uses, then they are considered to possess high validity (Kane, 2013).

The primary purpose of the NJSSA–S is to provide instructional information to classroom teachers and school/district educators about student needs for additional support upon returning to school after the COVID pandemic hiatus. The NJSSA–S produces the resources used to locally evaluate the needs of students. The assessment provides an initial indication of conceptual or skill gaps that might exist in a student’s understanding of NJSLS–S and the level of support students may need to inform instruction. The information provided by the assessment is only one piece of the puzzle used to holistically understand a student’s academic performance. The data should be used with other supporting evidence (e.g., assessments, homework, etc.) in guiding instruction. It is important to note that the NJSSA–S does not assess all the learning standards on the summative assessment. The NJSSA–S is not a replacement for the NJSLS–S. Nonetheless, the NJSSA–S provides applicable information. This information can be evaluated in terms of its validity evidence.

Conceptually, Kane (2006) labeled the process of evaluating validity evidence as validation. As such, test validation is an ongoing, ever-evolving process that extends through the duration of an assessment program. Every component within this technical brief, from test development to score reporting, is evidence both for and against the valid interpretation and uses of test scores.

The *Standards* (AERA, APA, NCME, 2014) categorize validity evidence into five sections:

- Evidence based on test content.
- Evidence based on response processes.
- Evidence based on internal structure.
- Evidence based on relation to other variables.
- Evidence based on the consequences of testing.

The following sections detail what evidence exists both for and against those five categories of validity evidence. Next, a section describes other validity evidence that was collected. Finally, the validity evidence pertaining to the intended interpretations is summarized. Overall, the evidence suggests that the 2022 NJSSA–S fosters valid interpretations and uses of test scores as they pertain to the overall classifications of students into support levels.

6.1 Evidence Based on Test Content

Validity evidence based on test content refers to the relevance of the content of the test to the construct the test is purporting to measure. *Standard 1.11* states that

[w]hen the rationale for test score interpretation for a given use rests in part on the appropriateness of test content, the procedures followed in specifying and generating content should be described and justified with reference to the intended population to be tested and the construct the test is intended to measure or the domain it is intended to represent (AERA, APA, NCME, 2014, p. 26).

The content-related evidence of validity includes the extent to which the test items represent the specified content domains and cognitive dimensions. Adequacy of the content representation of the NJSSA–S is critical because the tests purport to provide an indication of the support students need to progress toward achieving the knowledge, skills, and abilities (KSAs) identified in the NJSLS–S.

Adequate representation of the content domains defined in the NJSLS–S is assured through the use of a test blueprint and a responsible test construction process as was described in Part 2. The NJSLS–S was taken into consideration in the writing of all NJSSA–S items. In accordance with the test blueprint, the test construction process attempts to balance the six reporting categories and to ensure that the NJSSA–S contains an adequate representation of each content domain and scientific practice. Furthermore, all DCIs, SEPs, and CCCs are represented on the test. Section 2.4 provides a summary of test construction in comparison to the goals established in the test blueprint.

The test content was well-balanced at the content domain level (i.e., Earth and Space Science, Life Science, and Physical Science). At each grade level, the content domains were all within three points of being perfectly balanced. The scientific and engineering practices (i.e., Investigating, Sensemaking, and Critiquing) were less balanced. At each grade level, the Sensemaking scientific practice was over-represented and Critiquing was under-represented.

On a more granular level, the length of the test and the need for the items to be grouped as clusters made it impossible to effectively sample all the DCIs, SEPs, and CCCs. At grade 6, only six out of eleven DCIs were tested, including only one of three Earth and Space Science DCIs. Grade 9 had all eleven DCIs multiple DCIs represented on the test; however, three DCIs were only aligned to one item. At grade 12, three out of the eleven DCIs were not tested. The SEPs at both grades 6 and 9 did not include OEI items. Moreover, some CCCs were not represented at a given grade level.

Overall, the content domains and the range of DCIs, SEPs, and CCCs provide evidence that the test is adequately measuring the KSAs defined by the NJSLS–S. However, the relative lack of balance in the scientific and engineering practices and individual DCIs, SEPs and CCCs provides

evidence that the scale may be over-represented by certain components within the NJSLS–S, which could affect interpretations of test scores at both the overall and subscore levels.

6.2 Evidence Based on Response Processes

Standard 1.12 states that “[i]f the rationale for a test score interpretation for a given use depends on premises about the psychological processes or cognitive operations of test takers, then theoretical or empirical evidence in support of those premises should be provided” (AERA, APA, NCME, 2014, p. 26). Evidence based on response processes is complementary to evidence based on test content; it can come from several sources including response times, eye-tracking, think-aloud protocols, interviews, and/or focus groups. This complementary evidence is different from content evidence because its source is not content experts or teachers, but rather the actual student test takers. Padilla and Benitez (2014) noted that “validation studies aimed at obtaining evidence from response processes are scant” (p. 139), and at present time the NJSSA–S evidence based on response processes is limited to judgments from the NJSAC and content specialists.

The alignment of each item to the NJSLS–S Range PLDs provides limited evidence of the cognitive processes theoretically being assessed by the NJSSA–S. The Support Level Descriptors were based on the NJSSA–S Range PLDs, which were created in a collaborative effort by NJDOE, the NJSAC, content specialists, and psychometricians; they are based upon the NJSLS–S content standards. A more detailed description of the NJSLS–S PLD development process can be found in the 2019 NJSLS–S Technical Report. The detailed test maps presented in [Appendix A](#) display the NJSLS–S Range PLD alignment for each NJSSA–S item.

The NJSSA–S program does not currently have its own Range PLDs. It instead uses NJSLS–S Range PLDs as the theoretical cognitive structure underlying all current NJSSA–S item and test development. The NJSLS–S Range PLDs contain detailed descriptions of the KSAs that a student needs to display in order to be classified at a given support level. Each item on the NJSSA–S was aligned to two Range PLDs: one based on the DCI, and one based on the SEP. Those alignments were verified by the NJSAC. The alignment of each item to the Range PLDs offers a theoretical link from the NJSSA–S’s underlying cognitive structure to the student responses, which provides limited validity evidence based on response processes.

Table 6.2.1 shows the distributions of the support levels associated with each item by grade level and by DCI and SEP. The DCI distribution of items at grade 6 had more items aligned to “Strong Support” than were necessary, whereas the grade 9 DCI distribution had too many items aligned to “Less Support” and too few items aligned to “Strong Support.” Both grades had SEP distributions that were close to ideal. Grade 12’s SEP alignment had too many items aligned to “Strong Support,” but its DCI alignment distribution was close to ideal. These support level alignment distributions largely correspond to the item difficulty distributions illustrated in Figures 5.2.4 through 5.2.6.

Table 6.2.1: Support Level Alignment by DCI, SEP, and Grade Level

Grade	Domain/Practice	Strong Support	Some Support	Less Support
6	DCI	10	12	3
6	SEP	7	11	7
9	DCI	2	12	10
9	SEP	5	14	5
12	DCI	5	13	7
12	SEP	10	10	5

6.3 Evidence Based on Internal Structure

According to the *Standards*, “[a]nalysis of the internal structure of a test can indicate the degree to which the relationships among test items and test components conform to the construct on which the proposed test score interpretations are based” (AERA, APA, NCME, 2014, p. 16). The NJSSA–S was constructed as a unidimensional test. However, it also assesses student performance in several content clusters. It is important to study the pattern of relationships among the content clusters and testing methods. Therefore, this section addresses evidence based on responses and internal structure. Overall, the evidence supports the notion that the internal structure of the NJSSA–S is unidimensional, and that its items are measuring the same construct. However, at the subscore level, unexpected patterns of correlations provide evidence that the internal structure was not performing as intended.

6.3.1 Intercorrelations

One method of providing evidence supporting the inferences made from test scores is to evaluate the correlations among the total test score and its subscores. If the subscores are highly correlated, then that provides evidence that the test is unidimensional. Section 3.2.1.1 of this document summarizes correlation coefficients among test content domains and clusters by grade level. The intercorrelations of the NJSSA–S provide clear evidence that the NJSSA–S is unidimensional.

While the NJSSA–S was designed to be unidimensional, it was also designed to assess specific, distinct content domains (Earth and Space Science, Life Science, and Physical Science) and scientific practices (Critiquing, Investigating, and Sensemaking). The lowest correlation among all content subscores at all grade levels was .54 at Grade 12 between the Life and Physical Science content domain categories. Among the scientific practices, the lowest correlation observed was .54 between the Critiquing and Investigating practices. The moderate to high positive correlations between the content domains and between the scientific practices provide evidence for the interpretations of these domains and practices. The correlations suggest that each domain and practice is part of an overall unidimensional measure.

6.3.2 Other Internal Structure Evidence

Evidence of the internal structure of the NJSSA–S was also presented via a principal component analysis (PCA). Its results are presented in Section 3.2.1.2. These scree plots show further evidence that the variability in the NJSSA–S test scores is due to a single dimension. No secondary factors at any grade level practically contributed to explaining the variation in the test scores.

Part 5 of this technical brief provides ample evidence to support NJSSA–S reliability. Reliability is a measure of internal consistency that provides a sign as to whether the internal structure of the NJSSA–S is unidimensional. The grade-level reliability coefficients presented in Section 5.1 were relatively strong, ranging from .81 to .88. At the subscore level, the reliability coefficients were adequate given the dearth of points available for many subscores, with only grade 9 Critiquing falling below .50.

6.4 Evidence Based on Relationships to Other Variables

Evidence based on relationships to other variables takes the form of relationships between test scores and other variables that are external to the test (AERA, APA, NCME, 2014). This evidence can come from investigating the relationships among tests that measure similar constructs, tests that measure different constructs, or other outcomes that a test purports to predict. No evidence based on relationships of the NJSSA–S to other variables currently exists.

6.5 Evidence Based on the Consequences of Testing

Standard 1.25 states that “[w]hen unintended consequences result from test use, an attempt should be made to investigate whether such consequences arise from the test’s sensitivity to characteristics other than those it is intended to assess or from the test’s failure to fully represent the intended construct” (AERA, APA, NCME, 2014, p. 30). Lane and Stone (2002, p. 24) list the following types of evidence that can be collected to evaluate the consequences of a large-scale statewide accountability assessment program:

- Student, teacher, and administrator motivation and effort.
- Curriculum and instructional content and strategies.
- Content and format of classroom assessments.
- Improved learning for all students.
- Professional development support.
- Use and nature of test preparation activities.
- Student, teacher, administrator, and public awareness and beliefs about the assessment and criteria for judging performance and the use of assessment results.

No NJSSA–S validity evidence based on the consequences of testing exists at the moment.

6.6 Other Validity Evidence

Each section within this technical brief contributes evidence relevant to validity. The following is a summary of evidence within each section that is specific to the NJSSA–S:

Part 1: Introduction — This section describes the purpose of the assessment including:

- Intended inferences and uses of test scores.
- The relationship between the NJSLS–S and NJSSA–S.

Part 2: Test Development — This section describes the processes used to design and develop the NJSSA–S including:

- The steps taken to link test development to the NJSSA–S’ intended inferences and uses.
- The training and QC procedures implemented in the NJSLS–S item development process.
- The use of NJDOE, the NJSAC, and the Sensitivity committee during the initial creation of the items for the NJSLS–S to ensure the work of item writers and content specialists was aligned to the NJSLS–S.
- The steps taken to ensure the test construction process matched the NJSSA–S blueprint and statistical constraints.

Part 3: Item and Test Statistics — This section describes the battery of statistics that were used to evaluate the NJSSA–S at the test, item, and person levels including:

- Summaries of item performance across grade level, content domain, scientific practice, and item type to verify that the items are appropriate.
- Measures of test speededness to assess whether students could finish the test in the allotted time.
- Confirming the test items were not disadvantaging large subgroups of students via DIF statistics.
- Descriptive statistics of raw scores by test form and subgroups of students to evaluate how appropriate the test is for portions of the population.
- Evaluating the IRT assumptions of the PCM to ensure it is appropriate for modeling student ability estimates.
- Evaluating IRT person-fit statistics by subgroups of students.

Part 4: Scale Stability — This section describes the methods used to test the stability of the NJSSA–S scale including:

- Statistical and visual inspections of the stability of CTT and IRT item difficulty measures.
- Monitoring the IRT-based equating constant.

Part 5: Reliability — This section describes several reliability statistics that were calculated to verify the consistency of the NJSSA–S test scores including:

- Verifying the reliability at the total score, form, subscore, item type, and subgroup levels.
- Evaluating graphical displays of IRT reliability such as TIFs.
- Assessing the consistency of student support level classifications.

The following is a summary of validity evidence presented within this document that is specific to the entire NJSSA:

Section 2.5: Test Administration — This section describes the care that was taken to implement standardized test administration procedures including:

- Documents produced to communicate NJSSA–S test administration procedures for all versions of the test.
- Steps taken to ensure testing materials were handled using safe and secure procedures.
- Accommodations and accessibility features that were used during the test administration to provide all NJSSA–S test-takers with equal opportunities on the test.

Section 2.9: Scores and Score Reports — This section describes the procedures that were implemented to verify the accuracy of scoring student responses including:

- Confirming all computer-scored answer keys for both MC and TE item types.
- Verifying that student raw scores and subscores were calculated accurately.

6.7 Summary of Validity Evidence

Messick (1989) defined validity as “an integrated evaluative judgment of the degree to which empirical evidence and theoretical rationales support the adequacy and appropriateness of inferences and actions based on test scores and other modes of assessment” (p. 13). Making an integrated evaluative judgment with such a diverse assortment of evidence is challenging given that the validity process is ongoing and exists throughout the duration of the testing program. Overall, there is ample evidence that the NJSSA–S will foster valid inferences and uses. However, the NJSSA–S validity argument requires continuing attention, as well as an iterative process of identifying its weakest components, making modifications and then re-evaluating their effectiveness is needed. As Cronbach (1980) said, “The job of validation is not to support an interpretation, but to find out what might be wrong with it. A proposition deserves some degree of trust only when it has survived serious attempts to falsify it” (p. 103). The following sections set forth the pros and cons of the NJSSA–S validity evidence by the primary inferences and uses of the test.

The most important inferences made from the NJSSA–S involve the student support-level classifications. Students are classified as needing “Strong Support,” “Some Support,” or “Less

Support.” All interpretations based on NJSSA–S support-level classifications should be validated for evaluating student support as it pertains to the KSAs defined in the NJSLS–S.

Validity evidence in support of the proposed support-level classification interpretations has been presented throughout this document and within the validity section. The NJSSA–S was developed and constructed by well-trained experts with assistance from NJDOE and the NJSAC to specifically measure the wide range of KSAs defined in the NJSLS–S. It was administered under standardized processes and procedures. The accuracy of the scoring of all NJSLS–S items was verified. After the test administration, the items were statistically reviewed to ensure they met the assumptions of the proposed IRT model. Finally, both the overall scale and the support-level classifications were verified as being internally consistent.

There are some areas in which the validity evidence in support of the support-level classification inferences could be improved. The validity section on response processes contained limited evidence. Without a degree of evidence that student responses to test items are indeed measuring what the test is intended to measure, the validity argument is incomplete. Even if content experts and the NJSAC say an item measures a specific skill, that claim should be verified with evidence from the students who actually have to answer the item. The validity section on consequences also has no evidence, which is somewhat expected due to the challenge of integrating consequential validity evidence into a coherent validity argument (Cizek, 2016), as well as the difficulty of identifying the long-term consequences of a testing program after its first year of large-scale operational use. More pressingly, while there is ample validity evidence presented in both Part 5 of this document and in this validity section that the support-level classifications were consistent, the creation of the support levels themselves does not conform to best practices. There was no standard-setting or standards validation for the NJSSA–S. The support levels map directly to cut scores associated with the NJSLA–S, and the Reporting PLDs used on the score reports are based on the KSAs detailed in the NJSLA–S Range PLDs. If the NJSSA–S is continued in the future, it would behoove NJDOE to validate the use of the NJSLA–S performance standards as the basis for their support levels.

Overall, the evidence in favor of the valid interpretations of support-level classification outweighs the areas in which evidence is lacking or nonexistent. The NJSSA–S is a standards-based assessment; thus, the content validity evidence linking the test scores and interpretations to the NJSLS–S and the test blueprint is of chief importance (Sireci et al., 2008). However, there is clearly a need to study the issues noted above to enhance the validity evidence.

Appendix A: Detailed Test Maps

Table A.1: Grade 6 Test Map — Metadata

UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	NJSLA–S Range PLD Levels
518043_01	1	TE	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = B1; SEP = B3
518043_03	1	MC	CEDS	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = B1; SEP = B3
518043_05	1	MC	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = B1; SEP = B3
518008_01	1	MC	PACI	PS2	C and E	Physical Science	Investigating	DCI = B2; SEP = D2
518008_02	1	TE	CEDS	PS2	C and E	Physical Science	Sensemaking	DCI = B2; SEP = B1
518008_06	1	TE	EAE	PS2	C and E	Physical Science	Critiquing	DCI = B2; SEP = E2
518010_01	1	MC	AID	PS1	SC	Physical Science	Sensemaking	DCI = B2; SEP = B3
518010_03	1	TE	UMCT	PS1	PAT	Physical Science	Investigating	DCI = B2; SEP = B3
518010_05	1	TE	PACI	PS1	PAT	Physical Science	Investigating	DCI = B2; SEP = B3
518012_07	1	TE	EAE	LS4	SF	Life Science	Critiquing	DCI = C2; SEP = A2
518012_04	1	MC	EAE	LS4	SF	Life Science	Critiquing	DCI = A2; SEP = C2
518012_02	1	TE	EAE	LS4	SF	Life Science	Critiquing	DCI = A2; SEP = C2
519003_01a	1	TE	AQDP	LS1	SF	Life Science	Investigating	DCI = B1; SEP = A1
519003_02a	1	TE	PACI	LS1	PAT	Life Science	Investigating	DCI = B1; SEP = A1
519003_04a	1	TE	UMCT	LS1	S, P, and Q	Life Science	Investigating	DCI = A1; SEP = D1
519003_05a	1	TE	CEDS	LS1	C and E	Life Science	Sensemaking	DCI = B1; SEP = B2
518011_06	1	MC	UMCT	ESS2	S, P, and Q	Earth and Space Science	Investigating	DCI = A1; SEP = B2
518011_09	1	TE	EAE	ESS2	C and E	Earth and Space Science	Critiquing	DCI = A1; SEP = D2
518060_01	1	MC	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = C2; SEP = A1
518060_02	1	TE	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = C1; SEP = A1
518060_03	1	TE	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = C3; SEP = A1
519001_01a	1	TE	AQDP	PS3	E & M	Physical Science	Investigating	DCI = A2; SEP = C2
519001_07b	1	TE	DUM	PS3	E & M	Physical Science	Sensemaking	DCI = A3; SEP = E2
519001_08b	1	TE	EAE	PS3	E & M	Physical Science	Critiquing	DCI = B2; SEP = D2
519001_10b	1	TE	PACI	PS3	E & M	Physical Science	Investigating	DCI = A3; SEP = D3

Table A.2: Grade 6 Test Map — Item Statistics 2019

UIN	Points	Item Type	Rasch	P-Value	RPB	Median Time
518043_01	1	TE	-0.5302	.67	.51	174
518043_03	1	MC	-1.3127	.79	.38	79
518043_05	1	MC	-1.3463	.81	.45	54
518008_01	1	MC	-1.5379	.81	.31	67
518008_02	1	TE	-0.9556	.73	.44	32
518008_06	1	TE	1.7023	.24	.36	86
518010_01	1	MC	-0.9646	.73	.54	94
518010_03	1	TE	-0.3141	.63	.56	103
518010_05	1	TE	0.4156	.49	.58	73
518012_07	1	TE	0.1263	.53	.46	117
518012_04	1	MC	-0.9766	.74	.49	46
518012_02	1	TE	-0.5176	.67	.55	46
519003_01a	1	TE	-0.6058	.67	.37	71
519003_02a	1	TE	0.1824	.53	.42	45
519003_04a	1	TE	-0.2313	.60	.25	68
519003_05a	1	TE	-1.6682	.83	.31	40
518011_06	1	MC	-0.2299	.60	.28	43
518011_09	1	TE	-0.5062	.66	.37	82
518060_01	1	MC	-0.3923	.63	.30	90
518060_02	1	TE	1.2508	.31	.32	78
518060_03	1	TE	0.8560	.39	.30	97
519001_01a	1	TE	-0.0962	.58	.47	78
519001_07b	1	TE	-1.2766	.78	.38	40
519001_08b	1	TE	1.2180	.33	.53	64
519001_10b	1	TE	0.5602	.46	.59	53

Table A.3: Grade 9 Test Map — Metadata

UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	NJSLA–S Range PLD Levels
818077	1	TE	EAE	ESS2	PAT	Earth and Space Science	Critiquing	DCI = A3; SEP = C1
818307_01	1	TE	CEDS	ESS3	SF	Earth and Space Science	Sensemaking	DCI = A3; SEP = A2
818283	1	MC	EAE	PS2	C and E	Physical Science	Critiquing	DCI = B2; SEP = C1
818033_02	1	MC	AQDP	PS4	E & M	Physical Science	Investigating	DCI = B2; SEP = A2
818055_02	1	TE	DUM	LS2	E & M	Life Science	Sensemaking	DCI = A2; SEP = E2
818055_01	1	TE	CEDS	LS2	C and E	Life Science	Sensemaking	DCI = C3; SEP = B2
818055_03	1	TE	DUM	LS2	SC	Life Science	Sensemaking	DCI = B3; SEP = E2
818095_01	1	MC	CEDS	ESS3	C and E	Earth and Space Science	Sensemaking	DCI = A2; SEP = B2
818300_01	1	TE	UMCT	ESS1	S, P, and Q	Earth and Space Science	Investigating	DCI = B2; SEP = C2
818306_01	1	MC	AQDP	ESS3	SF	Earth and Space Science	Investigating	DCI = A2; SEP = A1
818302	1	MC	AQDP	ESS1	S, P, and Q	Earth and Space Science	Investigating	DCI = A2; SEP = C2
818267	1	MC	AQDP	ESS2	S & SM	Earth and Space Science	Investigating	DCI = A3; SEP = B2
818271	1	TE	EAE	ESS2	PAT	Earth and Space Science	Critiquing	DCI = A3; SEP = C3
818003_02a	1	TE	AID	PS3	PAT	Physical Science	Sensemaking	DCI = A3; SEP = A2
818003_01a	1	TE	DUM	PS3	E & M	Physical Science	Sensemaking	DCI = B3; SEP = D2
818003_03a	2	TE	EAE	PS3	E & M	Physical Science	Critiquing	DCI = A3; SEP = B3
818109	1	TE	AID	LS4	C and E	Life Science	Sensemaking	DCI = B2; SEP = D2
818296_02	1	TE	DUM	LS3	PAT	Life Science	Sensemaking	DCI = A2; SEP = E2
818065	1	TE	EAE	LS1	SF	Life Science	Critiquing	DCI = C3; SEP = C3
818062	1	MC	PACI	LS1	E & M	Life Science	Investigating	DCI = A1; SEP = A3
818250	1	TE	PACI	PS4	SF	Physical Science	Investigating	DCI = B2; SEP = B4
818285	1	TE	UMCT	PS2	S & SM	Physical Science	Investigating	DCI = B2; SEP = D2
818089_01	1	MC	AQDP	PS2	C and E	Physical Science	Investigating	DCI = A1; SEP = A1
818028	1	TE	AID	PS1	SF	Physical Science	Sensemaking	DCI = C2; SEP = D1

Table A.4: Grade 9 Test Map — Item Statistics 2019

UIN	Points	Item Type	Rasch	P-Value	RPB	Median Time
818077	1	TE	-0.5325	.52	.30	100
818307_01	1	TE	-0.9462	.61	.41	85
818283	1	MC	-0.8181	.59	.54	55
818033_02	1	MC	0.4011	.34	.47	60
818055_02	1	TE	-3.1142	.92	.34	60
818055_01	1	TE	-0.1868	.45	.46	95
818055_03	1	TE	-0.3573	.49	.49	59
818095_01	1	MC	-1.7725	.76	.49	49
818300_01	1	TE	1.1610	.20	.25	105
818306_01	1	MC	-1.0357	.63	.47	59
818302	1	MC	0.1264	.38	.39	77
818267	1	MC	0.2461	.36	.41	53
818271	1	TE	0.4473	.33	.42	64
818003_02a	1	TE	-0.1778	.45	.35	82
818003_01a	1	TE	-0.0686	.42	.27	44
818003_03a	2	TE	0.9899	.49	.34	94
818109	1	TE	-1.6263	.74	.43	67
818296_02	1	TE	-0.2348	.46	.49	98
818065	1	TE	0.8934	.25	.44	53
818062	1	MC	-0.0026	.40	.26	80
818250	1	TE	1.3884	.17	.19	74
818285	1	TE	0.3165	.36	.56	95
818089_01	1	MC	-0.3900	.49	.31	54
818028	1	TE	0.5373	.31	.32	58

Table A.5: Grade 12 Test Map — Metadata

UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	NJSLA–S Range PLD Levels
HS18038_02	1	TE	AID	LS4	S, P, and Q	Life Science	Sensemaking	DCI = A1; SEP = A3
HS18038_10	1	TE	AID	LS4	S & SM	Life Science	Sensemaking	DCI = A1; SEP = E2
HS18038_12	1	TE	EAE	LS4	S, P, and Q	Life Science	Critiquing	DCI = A2; SEP = B2
HS18038_16	1	TE	OECI	LS4	S & SM	Life Science	Critiquing	DCI = A2; SEP = B2
HS18004_01	1	MC	OECI	PS1	PAT	Physical Science	Critiquing	DCI = B2; SEP = A1
HS18004_04	1	TE	DUM	PS1	PAT	Physical Science	Sensemaking	DCI = B2; SEP = C2
HS18004_05	1	MC	DUM	PS1	S & SM	Physical Science	Sensemaking	DCI = B2; SEP = F1
HS18069_01	1	TE	AID	LS2	S, P, and Q	Life Science	Sensemaking	DCI = A1; SEP = A1
HS18069_04	1	MC	AID	LS2	SC	Life Science	Sensemaking	DCI = C1; SEP = A1
HS18069_07	1	TE	OECI	LS2	SC	Life Science	Critiquing	DCI = C2; SEP = A1
HS18013_01	1	MC	UMCT	ESS1	S & SM	Earth and Space Science	Investigating	DCI = A3; SEP = D3
HS18013_03	1	MC	UMCT	ESS1	S & SM	Earth and Space Science	Investigating	DCI = A4; SEP = F3
HS18013_05	1	MC	DUM	ESS1	S, P, and Q	Earth and Space Science	Sensemaking	DCI = B1; SEP = C2
HS18040_01	1	MC	AQDP	PS3	E & M	Physical Science	Investigating	DCI = C3; SEP = A2
HS18040_03	1	TE	PACI	PS3	C and E	Physical Science	Investigating	DCI = C4; SEP = E1
HS18040_04	1	TE	PACI	PS3	C and E	Physical Science	Investigating	DCI = C3; SEP = E1
HS19004_01a	1	TE	AQD	ESS3	S & SM	Earth and Space Science	Investigating	DCI = B2; SEP = D4
HS19004_03a	1	TE	CEDS	ESS3	S & SM	Earth and Space Science	Sensemaking	DCI = A2; SEP = E1
HS19004_06b	1	TE	OECI	ESS3	C and E	Earth and Space Science	Critiquing	DCI = A2; SEP = A1
HS19004_07a	1	TE	EAE	ESS3	S & SM	Earth and Space Science	Critiquing	DCI = B3; SEP = E3
HS18001_01	1	TE	UMCT	PS2	C and E	Physical Science	Investigating	DCI = A2; SEP = F2
HS18001_07	1	MC	UMCT	PS2	SC	Physical Science	Investigating	DCI = A3; SEP = F2
HS18071_01	1	MC	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = D2; SEP = A2
HS18071_04	1	MC	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = D2; SEP = A1
HS18071_05	1	TE	AID	ESS2	PAT	Earth and Space Science	Sensemaking	DCI = D2; SEP = A2

Table A.6: Grade 12 Test Map — Item Statistics 2019

UIN	Points	Item Type	Rasch	P-Value	RPB	Median Time
HS18038_02	1	TE	-0.6828	.58	.35	74
HS18038_10	1	TE	-0.2127	.48	.35	67
HS18038_12	1	TE	-1.4240	.71	.47	50
HS18038_16	1	TE	-0.7515	.59	.31	76
HS18004_01	1	MC	-0.0331	.45	.38	41
HS18004_04	1	TE	0.1956	.41	.25	39
HS18004_05	1	MC	-0.0117	.45	.27	38
HS18069_01	1	TE	-0.6516	.57	.55	59
HS18069_04	1	MC	-0.9354	.62	.42	36
HS18069_07	1	TE	0.2613	.39	.59	79
HS18013_01	1	MC	-1.0985	.65	.52	65
HS18013_03	1	MC	-1.3309	.69	.46	37
HS18013_05	1	MC	-0.5107	.54	.41	34
HS18040_01	1	MC	0.1944	.41	.27	32
HS18040_03	1	TE	0.5180	.34	.40	26
HS18040_04	1	TE	0.1295	.42	.59	35
HS19004_01a	1	TE	-0.2159	.48	.53	80
HS19004_03a	1	TE	0.5325	.34	.51	36
HS19004_06b	1	TE	0.6599	.32	.46	43
HS19004_07a	1	TE	0.8530	.29	.42	55
HS18001_01	1	TE	-0.8978	.61	.62	56
HS18001_07	1	MC	-0.0810	.46	.19	53
HS18071_01	1	MC	0.0555	.43	.43	47
HS18071_04	1	MC	-0.9538	.62	.52	38
HS18071_05	1	TE	-0.6498	.56	.55	21

Appendix B: Raw Score Cumulative Frequency Distributions

Table B.1: Grade 6 — Raw Score Cumulative Frequency Distribution — Gender

Raw Score	All Cum. #	All Cum. %	Female Cum. %	Male Cum. %
0	153	0.16	0.12	0.19
1	848	0.88	0.68	1.07
2	2,448	2.55	2.06	3.02
3	5,026	5.24	4.29	6.15
4	8,329	8.68	7.41	9.90
5	12,023	12.53	11.15	13.86
6	15,783	16.45	15.09	17.76
7	19,712	20.55	19.34	21.71
8	23,483	24.48	23.44	25.48
9	27,380	28.54	27.90	29.16
10	31,413	32.75	32.36	33.12
11	35,765	37.28	37.32	37.26
12	40,335	42.05	42.38	41.74
13	45,008	46.92	47.58	46.30
14	50,081	52.21	53.23	51.24
15	55,414	57.77	59.16	56.45
16	60,875	63.46	65.15	61.85
17	66,478	69.30	71.06	67.62
18	72,056	75.12	76.85	73.46
19	77,682	80.98	82.61	79.42
20	83,002	86.53	88.01	85.10
21	87,751	91.48	92.59	90.41
22	91,609	95.50	96.30	94.73
23	94,180	98.18	98.54	97.84
24	95,548	99.61	99.71	99.51
25	95,926	100.00	100.00	100.00

Table B.2: Grade 6 — Raw Score Cumulative Frequency Distribution — Ethnicity

Raw Score	All Cum. #	All Cum. %	Am. Indian Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	Pacific Islander Cum. %	White Cum. %
0	153	0.16	0.61	0.09	0.25	0.22	0.00	0.10
1	848	0.88	0.61	0.21	1.63	1.32	0.00	0.44
2	2,448	2.55	1.22	0.61	4.64	3.90	0.53	1.27
3	5,026	5.24	2.44	1.16	9.54	8.08	2.63	2.61
4	8,329	8.68	4.88	1.89	15.46	13.53	5.26	4.34
5	12,023	12.53	7.93	2.72	22.18	19.34	8.95	6.46
6	15,783	16.45	14.63	3.92	28.55	25.10	11.05	8.80
7	19,712	20.55	17.68	5.25	34.93	30.90	15.26	11.53
8	23,483	24.48	25.00	6.65	40.75	36.39	17.37	14.28
9	27,380	28.54	28.66	8.40	46.58	41.58	21.05	17.42
10	31,413	32.75	34.15	10.37	51.67	46.82	23.68	21.15
11	35,765	37.28	39.02	12.61	56.73	52.19	30.53	25.44
12	40,335	42.05	44.51	15.65	61.84	57.57	35.26	30.11
13	45,008	46.92	48.17	18.76	66.87	62.80	39.47	35.12
14	50,081	52.21	53.05	22.80	72.03	68.11	45.26	40.76
15	55,414	57.77	59.15	27.35	76.66	73.55	49.47	47.03
16	60,875	63.46	63.41	32.64	81.00	78.60	53.68	53.86
17	66,478	69.30	67.07	39.60	85.53	83.14	62.11	60.91
18	72,056	75.12	71.34	47.06	89.34	87.42	68.42	68.21
19	77,682	80.98	80.49	56.49	92.53	91.07	76.32	75.76
20	83,002	86.53	87.80	66.54	95.21	94.17	83.68	82.96
21	87,751	91.48	92.07	76.64	97.21	96.67	91.58	89.45
22	91,609	95.50	95.12	86.27	98.62	98.36	95.26	94.69
23	94,180	98.18	99.39	93.66	99.52	99.40	97.37	98.00
24	95,548	99.61	100.00	98.35	99.90	99.89	100.00	99.63
25	95,926	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table B.3: Grade 6 — Raw Score Cumulative Frequency Distribution — Other Demographics

Raw Score	All Cum. #	All Cum. %	EL – Yes Cum. %	EL – No Cum. %	Econ. Dis – Yes Cum. %	Econ. Dis – No Cum. %	SWD – Yes Cum. %	SWD – No Cum. %
0	153	0.16	0.64	0.12	0.20	0.14	0.36	0.11
1	848	0.88	3.30	0.68	1.35	0.65	2.37	0.50
2	2,448	2.55	9.39	1.97	4.21	1.74	6.64	1.50
3	5,026	5.24	18.69	4.10	8.70	3.55	13.11	3.22
4	8,329	8.68	29.75	6.90	14.50	5.85	20.68	5.61
5	12,023	12.53	40.24	10.19	20.74	8.54	28.04	8.56
6	15,783	16.45	49.69	13.64	26.84	11.40	34.95	11.71
7	19,712	20.55	58.91	17.30	33.12	14.43	41.56	15.16
8	23,483	24.48	65.84	20.98	38.89	17.47	47.17	18.67
9	27,380	28.54	71.44	24.91	44.40	20.82	51.98	22.54
10	31,413	32.75	76.56	29.04	49.70	24.50	56.89	26.56
11	35,765	37.28	80.90	33.59	55.09	28.62	61.50	31.08
12	40,335	42.05	84.99	38.41	60.52	33.06	66.16	35.87
13	45,008	46.92	88.27	43.42	65.73	37.77	70.27	40.94
14	50,081	52.21	91.33	48.90	71.09	43.02	74.45	46.51
15	55,414	57.77	93.87	54.71	76.36	48.72	78.35	52.49
16	60,875	63.46	95.78	60.72	81.16	54.85	81.89	58.74
17	66,478	69.30	97.07	66.95	85.51	61.41	85.39	65.18
18	72,056	75.12	98.10	73.17	89.29	68.22	88.63	71.65
19	77,682	80.98	98.97	79.46	92.61	75.32	91.58	78.26
20	83,002	86.53	99.45	85.43	95.27	82.27	94.11	84.58
21	87,751	91.48	99.71	90.78	97.34	88.63	96.41	90.21
22	91,609	95.50	99.91	95.13	98.69	93.95	98.20	94.81
23	94,180	98.18	100.00	98.03	99.54	97.52	99.36	97.88
24	95,548	99.61	100.00	99.57	99.92	99.45	99.89	99.53
25	95,926	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table B.4: Grade 9 — Raw Score Cumulative Frequency Distribution — Gender

Raw Score	All Cum. #	All Cum. %	Female Cum. %	Male Cum. %
0	145	0.15	0.11	0.18
1	864	0.87	0.74	1.00
2	2,796	2.82	2.42	3.20
3	6,653	6.70	6.12	7.27
4	12,565	12.65	12.02	13.28
5	20,078	20.22	19.84	20.61
6	28,275	28.47	28.59	28.39
7	36,560	36.82	37.52	36.17
8	44,055	44.36	45.66	43.16
9	51,230	51.59	53.48	49.82
10	57,883	58.29	60.52	56.18
11	64,109	64.56	67.19	62.07
12	69,788	70.28	73.04	67.65
13	75,073	75.60	78.43	72.91
14	79,797	80.36	83.17	77.68
15	84,198	84.79	87.43	82.25
16	87,916	88.53	90.96	86.21
17	91,141	91.78	93.92	89.73
18	93,773	94.43	96.10	92.83
19	95,839	96.51	97.78	95.30
20	97,380	98.06	98.85	97.31
21	98,349	99.04	99.48	98.61
22	98,925	99.62	99.81	99.43
23	99,177	99.87	99.93	99.81
24	99,291	99.99	99.99	99.98
25	99,305	100.00	100.00	100.00

Table B.5: Grade 9 — Raw Score Cumulative Frequency Distribution — Ethnicity

Raw Score	All Cum. #	All Cum. %	Am. Indian Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	Pacific Islander Cum. %	White Cum. %
0	145	0.15	0.59	0.01	0.27	0.25	0.00	0.06
1	864	0.87	1.18	0.14	1.56	1.38	0.55	0.43
2	2,796	2.82	4.14	0.54	4.84	4.41	4.37	1.48
3	6,653	6.70	5.92	1.59	11.06	10.41	6.01	3.68
4	12,565	12.65	12.43	3.40	20.52	19.15	12.57	7.33
5	20,078	20.22	19.53	5.78	31.43	29.97	16.39	12.55
6	28,275	28.47	27.81	8.73	43.22	41.37	25.68	18.48
7	36,560	36.82	37.28	12.62	54.62	51.68	34.43	25.35
8	44,055	44.36	41.42	16.64	63.27	60.40	40.44	32.50
9	51,230	51.59	48.52	21.84	70.82	67.97	47.54	39.89
10	57,883	58.29	57.99	27.36	76.68	74.35	53.01	47.43
11	64,109	64.56	65.09	33.23	81.87	79.78	59.02	54.87
12	69,788	70.28	72.19	39.57	86.39	84.41	62.30	61.73
13	75,073	75.60	75.74	46.54	90.12	88.30	67.76	68.36
14	79,797	80.36	80.47	53.94	92.87	91.30	74.32	74.51
15	84,198	84.79	85.80	61.52	95.01	93.80	78.69	80.38
16	87,916	88.53	87.57	69.00	96.64	95.76	86.34	85.27
17	91,141	91.78	91.72	76.27	97.85	97.26	91.26	89.52
18	93,773	94.43	94.67	82.74	98.67	98.30	95.08	93.03
19	95,839	96.51	96.47	88.51	99.21	99.08	96.72	95.72
20	97,380	98.06	95.86	92.96	99.68	99.55	99.45	97.69
21	98,349	99.04	98.22	96.31	99.87	99.80	99.45	98.89
22	98,925	99.62	99.41	98.36	99.98	99.93	99.45	99.60
23	99,177	99.87	100.00	99.39	99.99	99.98	99.45	99.88
24	99,291	99.99	100.00	99.93	100.00	100.00	100.00	99.98
25	99,305	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table B.6: Grade 9 — Raw Score Cumulative Frequency Distribution — Other Demographics

Raw Score	All Cum. #	All Cum. %	EL – Yes Cum. %	EL – No Cum. %	Econ. Dis – Yes Cum. %	Econ. Dis – No Cum. %	SWD – Yes Cum. %	SWD – No Cum. %
0	145	0.15	0.68	0.11	0.22	0.11	0.27	0.12
1	864	0.87	3.74	0.68	1.33	0.67	1.88	0.64
2	2,796	2.82	11.08	2.28	4.62	2.05	6.00	2.08
3	6,653	6.70	23.57	5.60	11.10	4.82	13.32	5.17
4	12,565	12.65	39.78	10.89	20.35	9.36	23.87	10.05
5	20,078	20.22	55.85	17.90	31.50	15.40	36.17	16.52
6	28,275	28.47	69.75	25.79	43.08	22.23	47.47	24.07
7	36,560	36.82	80.65	33.97	53.68	29.61	57.33	32.06
8	44,055	44.36	87.37	41.57	62.45	36.64	65.32	39.51
9	51,230	51.59	91.97	48.96	69.89	43.77	71.94	46.88
10	57,883	58.29	95.22	55.89	76.30	50.59	77.10	53.93
11	64,109	64.56	96.97	62.45	81.51	57.32	81.65	60.60
12	69,788	70.28	97.89	68.48	85.98	63.57	85.15	66.83
13	75,073	75.60	98.85	74.09	89.55	69.64	88.32	72.65
14	79,797	80.36	99.34	79.12	92.39	75.22	90.73	77.95
15	84,198	84.79	99.60	83.82	94.68	80.56	92.96	82.90
16	87,916	88.53	99.79	87.80	96.46	85.14	94.72	87.10
17	91,141	91.78	99.84	91.25	97.74	89.23	96.40	90.71
18	93,773	94.43	99.93	94.07	98.69	92.61	97.62	93.69
19	95,839	96.51	99.97	96.28	99.29	95.32	98.49	96.05
20	97,380	98.06	100.00	97.94	99.67	97.37	99.16	97.81
21	98,349	99.04	100.00	98.97	99.85	98.69	99.54	98.92
22	98,925	99.62	100.00	99.59	99.92	99.49	99.82	99.57
23	99,177	99.87	100.00	99.86	99.96	99.83	99.94	99.85
24	99,291	99.99	100.00	99.98	100.00	99.98	99.99	99.98
25	99,305	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table B.7: Grade 12 — Raw Score Cumulative Frequency Distribution — Gender

Raw Score	All Cum. #	All Cum. %	Female Cum. %	Male Cum. %
0	64	0.07	0.06	0.08
1	346	0.39	0.33	0.45
2	1122	1.27	1.00	1.53
3	2,983	3.38	2.80	3.96
4	6,098	6.91	5.90	7.90
5	10,400	11.79	10.26	13.29
6	15,434	17.49	15.55	19.40
7	20,770	23.54	21.53	25.51
8	26,137	29.62	27.78	31.43
9	31,418	35.60	34.10	37.10
10	36,684	41.57	40.60	42.55
11	41,809	47.38	46.96	47.84
12	46,818	53.06	53.01	53.15
13	51,867	58.78	59.26	58.36
14	56,667	64.22	65.14	63.36
15	61,454	69.64	70.90	68.45
16	65,994	74.79	76.21	73.44
17	70,336	79.71	81.16	78.33
18	74,333	84.24	85.73	82.81
19	77,925	88.31	89.68	86.99
20	81,129	91.94	93.12	90.80
21	83,811	94.98	95.95	94.04
22	85,754	97.18	97.89	96.48
23	87,172	98.79	99.16	98.42
24	87,959	99.68	99.78	99.59
25	88,241	100.00	100.00	100.00

Table B.8: Grade 12 — Raw Score Cumulative Frequency Distribution — Ethnicity

Raw Score	All Cum. #	All Cum. %	Am. Indian Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	Pacific Islander Cum. %	White Cum. %
0	64	0.07	0.00	0.02	0.16	0.09	0.00	0.05
1	346	0.39	0.00	0.14	0.59	0.54	0.00	0.30
2	1122	1.27	0.00	0.45	2.07	1.76	0.00	0.94
3	2,983	3.38	3.64	1.19	5.55	4.63	1.53	2.48
4	6,098	6.91	5.45	2.22	11.17	9.53	3.06	5.16
5	10,400	11.79	10.00	3.77	18.54	16.30	7.14	8.90
6	15,434	17.49	15.45	5.89	27.50	24.13	9.69	13.20
7	20,770	23.54	20.91	8.57	36.52	32.58	14.80	17.67
8	26,137	29.62	26.36	11.16	45.30	40.68	19.90	22.55
9	31,418	35.60	32.73	14.30	53.54	48.38	28.57	27.48
10	36,684	41.57	43.64	18.12	60.49	55.75	36.73	32.74
11	41,809	47.38	52.73	22.09	66.99	62.55	41.84	38.19
12	46,818	53.06	56.36	26.20	72.45	68.83	44.90	43.95
13	51,867	58.78	62.73	31.07	77.61	74.60	54.59	49.97
14	56,667	64.22	68.18	36.51	82.40	79.40	59.18	56.03
15	61,454	69.64	73.64	42.06	86.11	84.04	68.37	62.38
16	65,994	74.79	75.45	47.81	89.96	87.74	72.96	68.70
17	70,336	79.71	80.91	55.13	92.59	90.92	76.53	74.83
18	74,333	84.24	88.18	62.16	94.83	93.87	81.12	80.41
19	77,925	88.31	91.82	69.38	96.54	95.90	88.27	85.72
20	81,129	91.94	92.73	76.98	97.87	97.44	92.35	90.44
21	83,811	94.98	95.45	84.57	98.82	98.56	96.43	94.21
22	85,754	97.18	98.18	90.69	99.40	99.31	97.45	96.81
23	87,172	98.79	98.18	95.39	99.78	99.79	98.98	98.71
24	87,959	99.68	99.09	98.61	99.95	99.97	100.00	99.69
25	88,241	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table B.9: Grade 12 — Raw Score Cumulative Frequency Distribution — Other Demographics

Raw Score	All Cum. #	All Cum. %	EL – Yes Cum. %	EL – No Cum. %	Econ. Dis. – Yes Cum. %	Econ. Dis. – No Cum. %	SWD – Yes Cum. %	SWD – No Cum. %
0	64	0.07	0.17	0.07	0.13	0.05	0.14	0.06
1	346	0.39	0.93	0.37	0.54	0.34	0.78	0.30
2	1122	1.27	3.44	1.16	1.70	1.11	2.47	0.99
3	2,983	3.38	9.31	3.08	4.86	2.84	6.17	2.71
4	6,098	6.91	18.76	6.32	9.67	5.91	12.30	5.63
5	10,400	11.79	31.08	10.82	16.56	10.06	20.54	9.70
6	15,434	17.49	44.59	16.14	24.59	14.92	29.48	14.63
7	20,770	23.54	57.20	21.86	32.97	20.12	38.47	19.98
8	26,137	29.62	67.82	27.72	41.10	25.46	46.45	25.61
9	31,418	35.60	76.75	33.55	49.01	30.74	53.79	31.27
10	36,684	41.57	83.62	39.48	56.41	36.19	60.17	37.14
11	41,809	47.38	88.21	45.35	63.07	41.69	65.69	43.01
12	46,818	53.06	92.19	51.11	69.19	47.21	70.74	48.84
13	51,867	58.78	94.70	56.99	74.99	52.90	75.34	54.83
14	56,667	64.22	96.28	62.62	79.94	58.52	79.41	60.60
15	61,454	69.64	97.76	68.24	84.38	64.30	83.21	66.41
16	65,994	74.79	98.76	73.59	88.28	69.90	86.63	71.96
17	70,336	79.71	99.47	78.72	91.25	75.53	89.63	77.34
18	74,333	84.24	99.71	83.47	93.99	80.70	92.24	82.33
19	77,925	88.31	99.90	87.73	96.06	85.50	94.27	86.89
20	81,129	91.94	99.90	91.54	97.45	89.94	96.26	90.91
21	83,811	94.98	99.95	94.73	98.68	93.64	97.76	94.32
22	85,754	97.18	100.00	97.04	99.38	96.39	98.80	96.80
23	87,172	98.79	100.00	98.73	99.78	98.43	99.56	98.60
24	87,959	99.68	100.00	99.66	99.97	99.58	99.88	99.63
25	88,241	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Appendix C: Item Parameter Estimates and Model Fit Tables

Table C.1: Grade 6 — IRT Item Parameter Estimates and Fit Statistics

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower	Item Mean
518008_01	-1.53790	1.06	1.12	.39	0.91	.00	.74
518008_02	-0.95565	0.94	0.86	.45	1.12	.10	.71
518008_06	1.70235	0.98	1.18	.31	0.99	.00	.20
518010_01	-0.96460	0.82	0.72	.58	1.31	.00	.66
518010_03	-0.31411	0.79	0.72	.62	1.44	.00	.56
518010_05	0.41557	0.87	0.82	.59	1.26	.00	.48
518011_06	-0.22986	1.09	1.18	.41	0.78	.00	.47
518011_09	-0.50623	1.07	1.11	.42	0.85	.00	.56
518012_02	-0.51762	0.85	0.80	.57	1.30	.00	.59
518012_04	-0.97657	0.87	0.79	.50	1.23	.04	.71
518012_07	0.12630	0.96	0.95	.49	1.09	.00	.48
518043_01	-0.53023	0.85	0.79	.56	1.29	.00	.60
518043_03	-1.31270	1.03	1.04	.39	0.95	.07	.73
518043_05	-1.34630	0.85	0.79	.52	1.22	.00	.73
518060_01	-0.39235	1.22	1.36	.29	0.54	.18	.59
518060_02	1.25079	1.09	1.41	.28	0.80	.03	.26
518060_03	0.85600	1.16	1.41	.28	0.65	.06	.33
519001_01a	-0.09625	0.91	0.87	.53	1.19	.00	.51
519001_07b	-1.27660	0.95	0.89	.48	1.09	.00	.71
519001_08b	1.21804	0.77	0.74	.47	1.35	.00	.23
519001_10b	0.56024	0.79	0.78	.56	1.39	.00	.36
519003_01a	-0.60579	1.13	1.14	.40	0.77	.01	.56
519003_02a	0.18245	1.01	0.99	.45	1.00	.00	.44
519003_04a	-0.23129	1.23	1.38	.29	0.49	.08	.55
519003_05a	-1.66820	1.02	1.14	.38	0.95	.00	.77

Table C.2: Grade 9 — IRT Item Parameter Estimates and Fit Statistics

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower	Item Mean
818003_01a	-0.06862	1.16	1.22	.22	0.63	.07	.38
818003_02a	-0.17777	1.07	1.08	.32	0.84	.05	.41
818003_03a	0.98992	1.17	1.28	.29	0.77	.07	.45
818028	0.53728	1.02	1.10	.30	0.95	.01	.26
818033_02	0.40108	0.98	0.98	.42	1.03	.00	.32
818055_01	-0.18681	0.95	0.96	.39	1.10	.00	.37
818055_02	-3.11420	0.85	0.59	.35	1.15	.00	.88
818055_03	-0.35735	0.88	0.83	.47	1.32	.00	.41
818062	-0.00256	1.11	1.22	.26	0.73	.06	.37
818065	0.89336	0.76	0.67	.37	1.29	.00	.16
818077	-0.53246	1.13	1.19	.25	0.63	.09	.48
818089_01	-0.39005	1.18	1.27	.19	0.51	.11	.41
818095_01	-1.77250	0.79	0.67	.46	1.41	.00	.73
818109	-1.62630	0.87	0.82	.46	1.27	.00	.66
818250	1.38845	1.09	1.52	.11	0.86	.02	.14
818267	0.24607	1.03	1.12	.36	0.91	.03	.34
818271	0.44726	0.98	0.93	.40	1.04	.00	.30
818283	-0.81809	0.85	0.80	.50	1.43	.00	.52
818285	0.31647	0.72	0.64	.51	1.49	.00	.23
818296_02	-0.23481	0.85	0.81	.48	1.36	.00	.37
818300_01	1.16099	1.11	1.36	.23	0.84	.03	.19
818302	0.12636	1.04	1.07	.36	0.91	.03	.37
818306_01	-1.03570	0.89	0.85	.43	1.30	.03	.61
818307_01	-0.94620	0.97	0.95	.37	1.09	.03	.59

Table C.3: Grade 12 — IRT Item Parameter Estimates and Fit Statistics

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower	Item Mean
HS18001_01	-0.89782	0.89	0.84	.57	1.28	.00	.51
HS18001_07	-0.08102	1.36	1.49	.08	0.07	.23	.47
HS18004_01	-0.03307	1.09	1.10	.29	0.79	.00	.38
HS18004_04	0.19559	1.21	1.31	.18	0.51	.08	.37
HS18004_05	-0.01170	1.23	1.33	.16	0.41	.12	.40
HS18013_01	-1.09850	0.99	0.94	.38	1.04	.03	.64
HS18013_03	-1.33090	0.97	0.98	.40	1.06	.00	.67
HS18013_05	-0.51066	1.00	0.99	.38	1.00	.03	.53
HS18038_02	-0.68282	1.07	1.08	.32	0.81	.02	.55
HS18038_10	-0.21273	1.08	1.10	.32	0.79	.10	.50
HS18038_12	-1.42390	0.82	0.72	.44	1.35	.00	.74
HS18038_16	-0.75146	1.08	1.11	.29	0.79	.13	.61
HS18040_01	0.19436	1.18	1.27	.24	0.57	.10	.41
HS18040_03	0.51802	1.07	1.05	.39	0.88	.03	.38
HS18040_04	0.12950	0.84	0.78	.54	1.38	.00	.41
HS18069_01	-0.65162	0.86	0.82	.50	1.37	.00	.57
HS18069_04	-0.93538	0.98	0.97	.37	1.04	.01	.63
HS18069_07	0.26127	0.87	0.81	.53	1.30	.00	.41
HS18071_01	0.05550	1.02	1.06	.37	0.93	.04	.43
HS18071_04	-0.95385	0.91	0.84	.45	1.24	.00	.63
HS18071_05	-0.64984	0.87	0.82	.50	1.35	.00	.55
HS19004_01a	-0.21588	0.91	0.87	.48	1.26	.00	.50
HS19004_03a	0.53253	0.92	0.86	.49	1.16	.00	.36
HS19004_06b	0.65995	0.97	0.93	.41	1.06	.00	.32
HS19004_07a	0.85302	1.02	0.94	.39	1.00	.00	.29

Appendix D: Scale Stability Results Tables

Table D.1: Grade 6 — Delta Plot Results Table

UIN	Item Type	Points	P-Value 2019	P-Value 2022	Delta 2019	Delta 2022	Distance	Decision
518008_01	MC	1	0.808	0.742	9.515	10.427	0.107	Stable
518008_02	TE	1	0.731	0.706	10.532	10.786	0.355	Stable
518008_06	TE	1	0.244	0.195	15.775	16.366	0.095	Stable
518010_01	MC	1	0.735	0.660	10.490	11.350	0.074	Stable
518010_03	TE	1	0.633	0.559	11.639	12.396	0.006	Stable
518010_05	TE	1	0.486	0.482	13.142	13.201	0.484	Flagged
518011_06	MC	1	0.601	0.470	11.975	13.301	0.411	Stable
518011_09	TE	1	0.655	0.556	11.403	12.396	0.172	Stable
518012_02	TE	1	0.667	0.593	11.272	12.090	0.047	Stable
518012_04	MC	1	0.739	0.712	10.440	10.786	0.290	Stable
518012_07	TE	1	0.528	0.477	12.715	13.201	0.182	Stable
518043_01	TE	1	0.669	0.599	11.250	11.987	0.010	Stable
518043_03	MC	1	0.791	0.728	9.762	10.549	0.020	Stable
518043_05	MC	1	0.805	0.727	9.562	10.549	0.161	Stable
518060_01	MC	1	0.625	0.591	11.722	12.090	0.270	Stable
518060_02	TE	1	0.310	0.263	14.980	15.573	0.097	Stable
518060_03	TE	1	0.386	0.327	14.155	14.760	0.092	Stable
519001_01a	TE	1	0.580	0.515	12.194	12.900	0.029	Stable
519001_07b	TE	1	0.777	0.706	9.946	10.786	0.058	Stable
519001_08b	TE	1	0.330	0.232	14.764	15.955	0.326	Stable
519001_10b	TE	1	0.459	0.360	13.417	14.434	0.197	Stable
519003_01a	TE	1	0.671	0.556	11.233	12.396	0.293	Stable
519003_02a	TE	1	0.531	0.437	12.694	13.604	0.119	Stable
519003_04a	TE	1	0.598	0.550	12.004	12.497	0.180	Stable
519003_05a	TE	1	0.832	0.768	9.155	10.045	0.090	Stable

Table D.2: Grade 9 — Delta Plot Results Table

UIN	Item Type	Points	P-Value 2019	P-Value 2022	Delta 2019	Delta 2022	Distance	Decision
818003_01a	TE	1	0.423	0.384	13.778	14.222	0.084	Stable
818003_02a	TE	1	0.448	0.413	13.528	13.910	0.131	Stable
818003_03a	TE	2	0.247	0.223	15.740	16.022	0.181	Stable
818028	TE	1	0.306	0.259	15.031	15.573	0.003	Stable
818033_02	MC	1	0.338	0.319	14.670	14.871	0.249	Stable
818055_01	TE	1	0.453	0.368	13.473	14.327	0.205	Stable
818055_02	TE	1	0.921	0.879	7.342	8.300	0.222	Stable
818055_03	TE	1	0.493	0.407	13.069	13.910	0.192	Stable
818062	MC	1	0.403	0.369	13.982	14.327	0.153	Stable
818065	TE	1	0.252	0.156	15.669	16.978	0.549	Flagged
818077	TE	1	0.516	0.476	12.842	13.201	0.153	Stable
818089_01	MC	1	0.485	0.412	13.150	13.910	0.135	Stable
818095_01	MC	1	0.761	0.730	10.168	10.549	0.163	Stable
818109	TE	1	0.743	0.661	10.386	11.350	0.254	Stable
818250	TE	1	0.171	0.140	16.798	17.321	0.001	Stable
818267	MC	1	0.362	0.341	14.417	14.650	0.228	Stable
818271	TE	1	0.328	0.297	14.785	15.098	0.168	Stable
818283	MC	1	0.592	0.519	12.073	12.799	0.101	Stable
818285	TE	1	0.356	0.231	14.481	15.955	0.656	Flagged
818296_02	TE	1	0.464	0.366	13.357	14.327	0.286	Stable
818300_01	TE	1	0.204	0.189	16.311	16.512	0.234	Stable
818302	MC	1	0.384	0.370	14.182	14.327	0.293	Stable
818306_01	MC	1	0.630	0.615	11.672	11.883	0.270	Stable
818307_01	TE	1	0.607	0.589	11.910	12.090	0.290	Stable

Table D.3: Grade 12 — Delta Plot Results Table

UIN	Item Type	Points	P-Value 2019	P-Value 2022	Delta 2019	Delta 2022	Distance	Decision
HS18001_01	TE	1	0.609	0.511	11.897	12.900	0.676	Flagged
HS18001_07	MC	1	0.459	0.471	13.412	13.301	0.119	Stable
HS18004_01	MC	1	0.450	0.377	13.501	14.222	0.465	Flagged
HS18004_04	TE	1	0.407	0.367	13.940	14.327	0.227	Stable
HS18004_05	MC	1	0.447	0.398	13.529	14.013	0.299	Stable
HS18013_01	MC	1	0.647	0.643	11.491	11.566	0.028	Stable
HS18013_03	MC	1	0.686	0.668	11.062	11.240	0.104	Stable
HS18013_05	MC	1	0.542	0.528	12.575	12.699	0.053	Stable
HS18038_02	TE	1	0.575	0.548	12.239	12.497	0.150	Stable
HS18038_10	TE	1	0.484	0.501	13.160	13.000	0.152	Stable
HS18038_12	TE	1	0.709	0.741	10.804	10.427	0.284	Stable
HS18038_16	TE	1	0.589	0.606	12.098	11.883	0.181	Stable
HS18040_01	MC	1	0.409	0.415	13.919	13.910	0.052	Stable
HS18040_03	TE	1	0.344	0.383	14.611	14.222	0.324	Stable
HS18040_04	TE	1	0.421	0.411	13.802	13.910	0.032	Stable
HS18069_01	TE	1	0.565	0.569	12.343	12.295	0.066	Stable
HS18069_04	MC	1	0.620	0.630	11.776	11.673	0.100	Stable
HS18069_07	TE	1	0.392	0.409	14.096	13.910	0.177	Stable
HS18071_01	MC	1	0.430	0.432	13.708	13.705	0.045	Stable
HS18071_04	MC	1	0.618	0.626	11.795	11.673	0.113	Stable
HS18071_05	TE	1	0.560	0.549	12.392	12.497	0.042	Stable
HS19004_01a	TE	1	0.483	0.504	13.167	13.000	0.157	Stable
HS19004_03a	TE	1	0.343	0.362	14.622	14.434	0.183	Stable
HS19004_06b	TE	1	0.320	0.315	14.869	14.871	0.051	Stable
HS19004_07a	TE	1	0.288	0.293	15.237	15.214	0.072	Stable

Table D.4: Grade 6 — 0.3 Logits Absolute Difference Results Table

UIN	Item Type	Points	Rasch 2019	Rasch 2022	EQK	Adjusted Difficulty	Difficulty Difference	Absolute Difference	Decision
518008_01	MC	1	-1.538	-1.228	-.286	-1.514	-0.024	0.024	Stable
518008_02	TE	1	-0.956	-0.981	-.286	-1.267	0.312	0.312	Flagged
518008_06	TE	1	1.702	2.156	-.286	1.870	-0.168	0.168	Stable
518010_01	MC	1	-0.965	-0.695	-.286	-0.981	0.016	0.016	Stable
518010_03	TE	1	-0.314	-0.098	-.286	-0.384	0.070	0.070	Stable
518010_05	TE	1	0.416	0.338	-.286	0.052	0.364	0.364	Flagged
518011_06	MC	1	-0.230	0.404	-.286	0.118	-0.348	0.348	Flagged
518011_09	TE	1	-0.506	-0.079	-.286	-0.365	-0.142	0.142	Stable
518012_02	TE	1	-0.518	-0.290	-.286	-0.576	0.058	0.058	Stable
518012_04	MC	1	-0.977	-1.025	-.286	-1.311	0.334	0.334	Flagged
518012_07	TE	1	0.126	0.361	-.286	0.075	0.051	0.051	Stable
518043_01	TE	1	-0.530	-0.329	-.286	-0.615	0.085	0.085	Stable
518043_03	MC	1	-1.313	-1.131	-.286	-1.417	0.105	0.105	Stable
518043_05	MC	1	-1.346	-1.125	-.286	-1.411	0.064	0.064	Stable
518060_01	MC	1	-0.392	-0.280	-.286	-0.566	0.174	0.174	Stable
518060_02	TE	1	1.251	1.653	-.286	1.367	-0.117	0.117	Stable
518060_03	TE	1	0.856	1.232	-.286	0.946	-0.090	0.090	Stable
519001_01a	TE	1	-0.096	0.153	-.286	-0.133	0.037	0.037	Stable
519001_07b	TE	1	-1.277	-0.984	-.286	-1.270	-0.006	0.006	Stable
519001_08b	TE	1	1.218	1.873	-.286	1.587	-0.369	0.369	Flagged
519001_10b	TE	1	0.560	1.034	-.286	0.748	-0.188	0.188	Stable
519003_01a	TE	1	-0.606	-0.082	-.286	-0.368	-0.238	0.238	Stable
519003_02a	TE	1	0.182	0.586	-.286	0.300	-0.118	0.118	Stable
519003_04a	TE	1	-0.231	-0.047	-.286	-0.333	0.102	0.102	Stable
519003_05a	TE	1	-1.668	-1.417	-.286	-1.703	0.035	0.035	Stable

Table D.5: Grade 9 — 0.3 Logits Absolute Difference Results Table

UIN	Item Type	Points	Rasch 2019	Rasch 2022	EQK	Adjusted Difficulty	Difficulty Difference	Absolute Difference	Decision
818003_01a	TE	1	-0.069	0.066	-.198	-0.132	0.064	0.064	Stable
818003_02a	TE	1	-0.178	-0.088	-.198	-0.285	0.108	0.108	Stable
818003_03a	TE	2	0.990	1.178	-.198	0.980	0.010	0.010	Stable
818028	TE	1	0.537	0.790	-.198	0.592	-0.055	0.055	Stable
818033_02	MC	1	0.401	0.429	-.198	0.231	0.170	0.170	Stable
818055_01	TE	1	-0.187	0.152	-.198	-0.046	-0.141	0.141	Stable
818055_02	TE	1	-3.114	-2.950	-.198	-3.147	0.033	0.033	Stable
818055_03	TE	1	-0.357	-0.054	-.198	-0.252	-0.106	0.106	Stable
818062	MC	1	-0.003	0.150	-.198	-0.048	0.046	0.046	Stable
818065	TE	1	0.893	1.560	-.198	1.363	-0.469	0.469	Flagged
818077	TE	1	-0.532	-0.412	-.198	-0.609	0.077	0.077	Stable
818089_01	MC	1	-0.390	-0.081	-.198	-0.278	-0.112	0.112	Stable
818095_01	MC	1	-1.773	-1.776	-.198	-1.974	0.201	0.201	Stable
818109	TE	1	-1.626	-1.370	-.198	-1.568	-0.058	0.058	Stable
818250	TE	1	1.388	1.712	-.198	1.514	-0.126	0.126	Stable
818267	MC	1	0.246	0.304	-.198	0.106	0.140	0.140	Stable
818271	TE	1	0.447	0.558	-.198	0.360	0.087	0.087	Stable
818283	MC	1	-0.818	-0.631	-.198	-0.828	0.010	0.010	Stable
818285	TE	1	0.316	0.979	-.198	0.782	-0.465	0.465	Flagged
818296_02	TE	1	-0.235	0.165	-.198	-0.032	-0.202	0.202	Stable
818300_01	TE	1	1.161	1.286	-.198	1.088	0.073	0.073	Stable
818302	MC	1	0.126	0.140	-.198	-0.058	0.184	0.184	Stable
818306_01	MC	1	-1.036	-1.123	-.198	-1.321	0.285	0.285	Stable
818307_01	TE	1	-0.946	-0.985	-.198	-1.183	0.237	0.237	Stable

Table D.6: Grade 12 — 0.3 Logits Absolute Difference Results Table

UIN	Item Type	Points	Rasch 2019	Rasch 2022	EQK	Adjusted Difficulty	Difficulty Difference	Absolute Difference	Decision
HS18001_01	TE	1	-0.898	-0.110	-.282	-0.391	-0.506	0.506	Flagged
HS18001_07	MC	1	-0.081	0.096	-.282	-0.186	0.105	0.105	Stable
HS18004_01	MC	1	-0.033	0.589	-.282	0.307	-0.340	0.340	Flagged
HS18004_04	TE	1	0.196	0.645	-.282	0.363	-0.167	0.167	Stable
HS18004_05	MC	1	-0.012	0.477	-.282	0.195	-0.207	0.207	Stable
HS18013_01	MC	1	-1.099	-0.792	-.282	-1.074	-0.025	0.025	Stable
HS18013_03	MC	1	-1.331	-0.933	-.282	-1.214	-0.117	0.117	Stable
HS18013_05	MC	1	-0.511	-0.195	-.282	-0.476	-0.035	0.035	Stable
HS18038_02	TE	1	-0.683	-0.295	-.282	-0.577	-0.106	0.106	Stable
HS18038_10	TE	1	-0.213	-0.056	-.282	-0.338	0.125	0.125	Stable
HS18038_12	TE	1	-1.424	-1.362	-.282	-1.644	0.220	0.220	Stable
HS18038_16	TE	1	-0.751	-0.594	-.282	-0.876	0.125	0.125	Stable
HS18040_01	MC	1	0.194	0.389	-.282	0.107	0.088	0.088	Stable
HS18040_03	TE	1	0.518	0.555	-.282	0.273	0.245	0.245	Stable
HS18040_04	TE	1	0.130	0.407	-.282	0.125	0.005	0.005	Stable
HS18069_01	TE	1	-0.652	-0.401	-.282	-0.683	0.031	0.031	Stable
HS18069_04	MC	1	-0.935	-0.725	-.282	-1.006	0.071	0.071	Stable
HS18069_07	TE	1	0.261	0.416	-.282	0.134	0.127	0.127	Stable
HS18071_01	MC	1	0.056	0.298	-.282	0.016	0.040	0.040	Stable
HS18071_04	MC	1	-0.954	-0.703	-.282	-0.985	0.031	0.031	Stable
HS18071_05	TE	1	-0.650	-0.301	-.282	-0.583	-0.067	0.067	Stable
HS19004_01a	TE	1	-0.216	-0.074	-.282	-0.356	0.140	0.140	Stable
HS19004_03a	TE	1	0.533	0.668	-.282	0.386	0.146	0.146	Stable
HS19004_06b	TE	1	0.660	0.936	-.282	0.655	0.005	0.005	Stable
HS19004_07a	TE	1	0.853	1.066	-.282	0.785	0.068	0.068	Stable

Appendix E: Raw-to-Theta Score Tables

Table E.1: Grade 6 — Raw-to-Theta Score Table

Raw Score	Theta	CSEM	Information	Support Level
0	-5.0236	1.8429	0.29	Strong Support
1	-3.7750	1.0317	0.94	Strong Support
2	-3.0164	0.7527	1.77	Strong Support
3	-2.5432	0.6342	2.49	Strong Support
4	-2.1855	0.5669	3.11	Strong Support
5	-1.8897	0.5235	3.65	Strong Support
6	-1.6318	0.4937	4.10	Strong Support
7	-1.3989	0.4725	4.48	Strong Support
8	-1.1831	0.4574	4.78	Strong Support
9	-0.9791	0.4467	5.01	Strong Support
10	-0.7830	0.4395	5.18	Strong Support
11	-0.5919	0.4353	5.28	Strong Support
12	-0.4033	0.4337	5.32	Strong Support
13	-0.2149	0.4347	5.29	Some Support
14	-0.0246	0.4382	5.21	Some Support
15	0.1697	0.4442	5.07	Some Support
16	0.3709	0.4533	4.87	Some Support
17	0.5818	0.4658	4.61	Some Support
18	0.8063	0.4826	4.29	Some Support
19	1.0498	0.5053	3.92	Less Support
20	1.3202	0.5364	3.48	Less Support
21	1.6307	0.5806	2.97	Less Support
22	2.0053	0.6481	2.38	Less Support
23	2.4974	0.7659	1.70	Less Support
24	3.2772	1.0424	0.92	Less Support
25	4.5419	1.8494	0.29	Less Support

Table E.2: Grade 9 — Raw-to-Theta Score Table

Raw Score	Theta	CSEM	Information	Support Level
0	-5.2107	1.8725	0.29	Strong Support
1	-3.8920	1.0747	0.87	Strong Support
2	-3.0554	0.7957	1.58	Strong Support
3	-2.5249	0.6719	2.22	Strong Support
4	-2.1244	0.5987	2.79	Strong Support
5	-1.7961	0.5500	3.31	Strong Support
6	-1.5131	0.5156	3.76	Strong Support
7	-1.2607	0.4904	4.16	Strong Support
8	-1.0297	0.4718	4.49	Strong Support
9	-0.8138	0.4582	4.76	Some Support
10	-0.6086	0.4484	4.97	Some Support
11	-0.4106	0.4420	5.12	Some Support
12	-0.2170	0.4385	5.20	Some Support
13	-0.0252	0.4378	5.22	Some Support
14	0.1670	0.4398	5.17	Some Support
15	0.3624	0.4447	5.06	Some Support
16	0.5634	0.4527	4.88	Less Support
17	0.7734	0.4644	4.64	Less Support
18	0.9964	0.4808	4.33	Less Support
19	1.2380	0.5033	3.95	Less Support
20	1.5064	0.5345	3.50	Less Support
21	1.8150	0.5791	2.98	Less Support
22	2.1882	0.6475	2.39	Less Support
23	2.6801	0.7663	1.70	Less Support
24	3.4614	1.0437	0.92	Less Support
25	4.7286	1.8506	0.29	Less Support

Table E.3: Grade 12 — Raw-to-Theta Score Table

Raw Score	Theta	CSEM	Information	Support Level
0	-4.8788	1.8410	0.30	Strong Support
1	-3.6351	1.0282	0.95	Strong Support
2	-2.8839	0.7477	1.79	Strong Support
3	-2.4184	0.6279	2.54	Strong Support
4	-2.0687	0.5596	3.19	Strong Support
5	-1.7813	0.5152	3.77	Strong Support
6	-1.5322	0.4845	4.26	Strong Support
7	-1.3086	0.4624	4.68	Strong Support
8	-1.1026	0.4463	5.02	Strong Support
9	-0.9089	0.4346	5.29	Strong Support
10	-0.7237	0.4265	5.50	Strong Support
11	-0.5442	0.4213	5.63	Strong Support
12	-0.3679	0.4188	5.70	Strong Support
13	-0.1927	0.4187	5.70	Some Support
14	-0.0165	0.4211	5.64	Some Support
15	0.1627	0.4262	5.51	Some Support
16	0.3475	0.4342	5.30	Some Support
17	0.5408	0.4457	5.03	Less Support
18	0.7463	0.4618	4.69	Less Support
19	0.9693	0.4838	4.27	Less Support
20	1.2177	0.5145	3.78	Less Support
21	1.5043	0.5588	3.20	Less Support
22	1.8531	0.6272	2.54	Less Support
23	2.3177	0.7471	1.79	Less Support
24	3.0680	1.0277	0.95	Less Support
25	4.3109	1.8408	0.30	Less Support

Appendix F: Conditional Standard Error of Measurement and Test Characteristic Curve Graphs

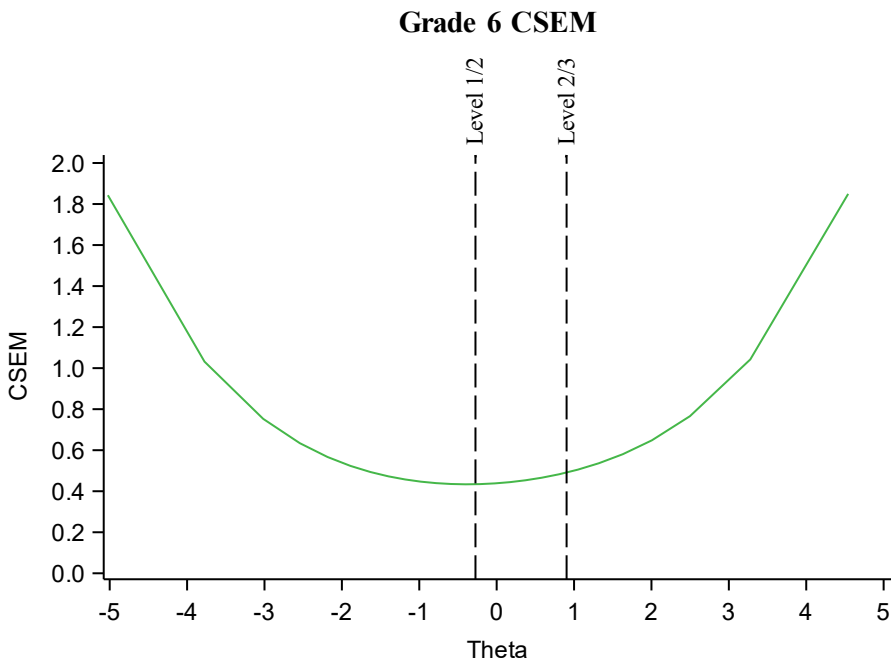


Figure F.1. Grade 6 Conditional Standard Error of Measurement

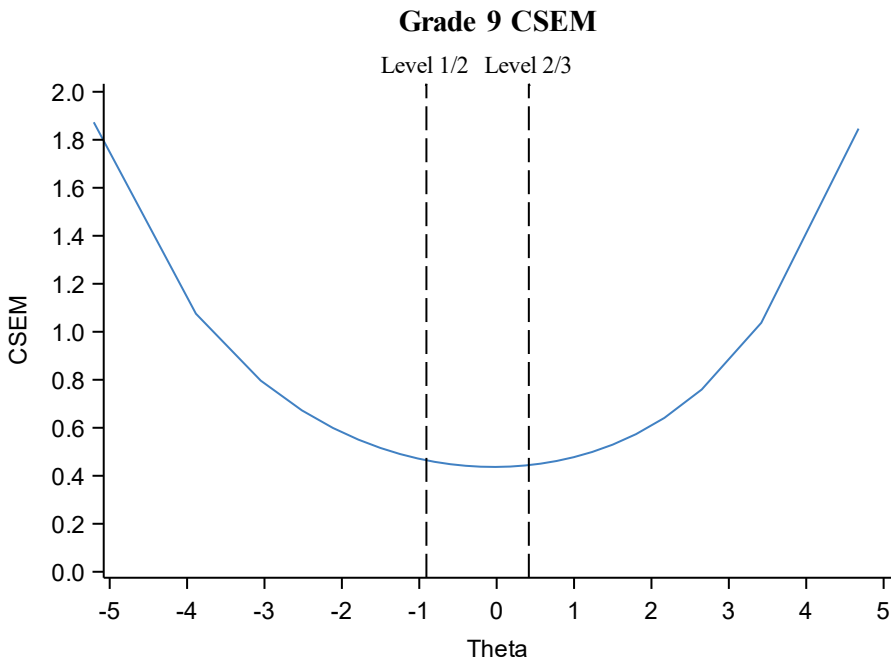


Figure F.2. Grade 9 Conditional Standard Error of Measurement

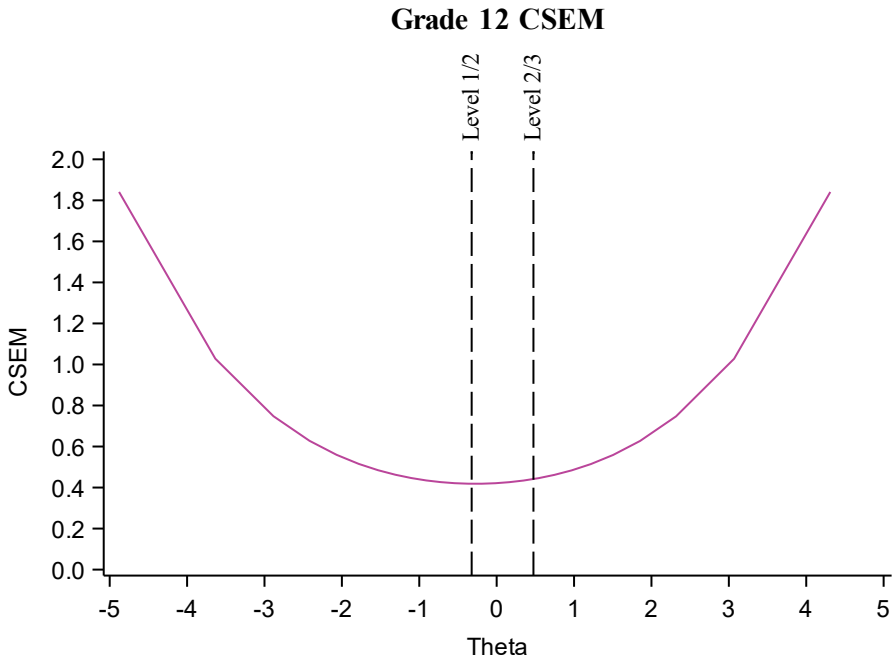


Figure F.3. Grade 12 Conditional Standard Error of Measurement

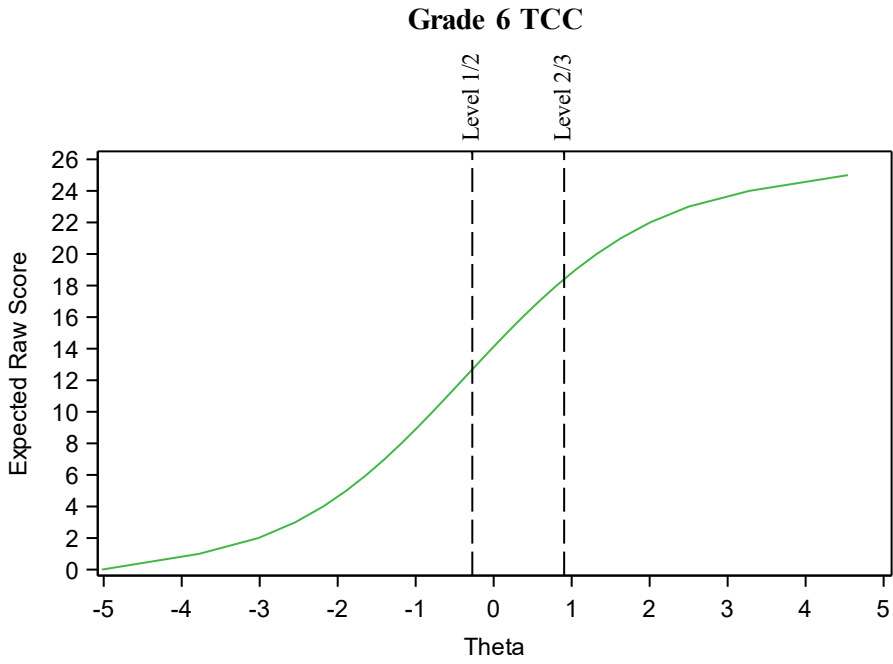


Figure F.4. Grade 6 Test Characteristic Curve

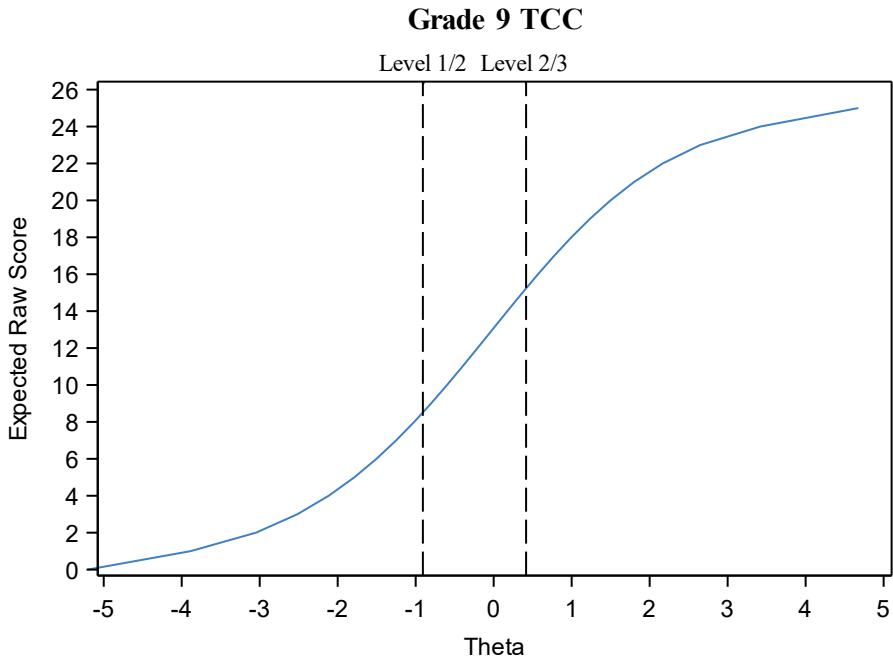


Figure F.5. Grade 9 Test Characteristic Curve

Grade 12 TCC

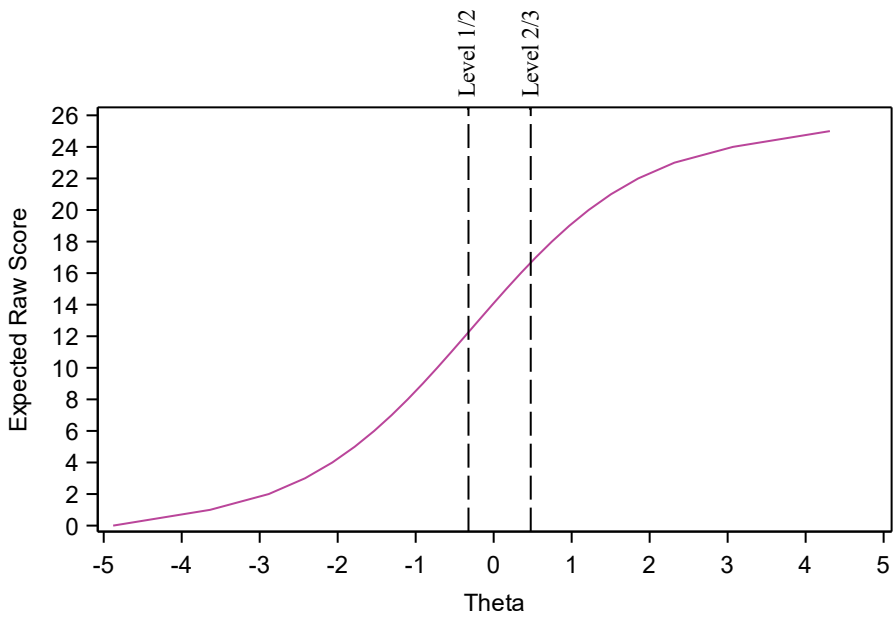


Figure F.6. Grade 12 Test Characteristic Curve

References

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for Educational and Psychological Testing*. Washington, D.C.: American Educational Research Association.
- Angoff, W. H. (1972, September). *A technique for the investigation of cultural differences*. Paper presented at the annual meeting of the American Psychological Association, Honolulu. (ERIC Document Reproduction Service No. ED 069 686)
- Angoff, W. H. (1982). Use of difficulty and discrimination indices for detecting item bias. In R. A. Berk (Ed.), *Handbook of methods for detecting test bias* (pp. 96–116). Baltimore: Johns Hopkins University Press.
- Brennan, R. L. (2004). Manual for BB-CLASS: A computer program that uses the beta-binomial model for classification consistency and accuracy (version 1). CASMA Research Report 9. Iowa City, IA.
- Brown, W. (1910). Some experimental results in the correlation of mental abilities. *British Journal of Psychology*, 3, 296–322.
- Brown, T. A. (2006). *Confirmatory Factor Analysis for Applied Research*. New York, NY: The Guilford Press.
- Camilli, G., & Shepard, L. A. (1994). *Methods for identifying biased test items*. Newbury Park, CA: Sage.
- Chen, W. H., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, 22(3), 265–289.
- Cizek, G. J. (2016). Validating test score meaning and defending test score use: different aims, different methods. *Assessment in Education: Principles, Policy, & Practice*, 23:2, 212–225.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1), 37-46.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Cronbach, L. J. (1980). Validity on parole: How can we go straight? *New Directions for Testing and Measurement: Measuring Achievement over a Decade*, 5, 99–108.
- Dorans, N. J., & Holland, P. W. (1993). *DIF detection and description: Mantel-Haenszel and standardization*. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). Hillsdale, NJ: Lawrence Erlbaum.
- Every Student Succeeds Act, 20 U.S.C. § 6301 (2015). <https://www.congress.gov/bill/114th-congress/senate-bill/1177>

- Gorsuch, R. L. (1983). *Factor Analysis*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Hambleton, R. K. (1989). Principles and selected applications of item response theory. In R. L. Linn (Ed.), *Educational measurement* (3rd ed.). New York: American Council on Education and Macmillan.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Norwell, MA: Kluwer Academic Publishers.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage.
- Hambleton, R. K., & van der Linden, W. J. (1982). Advanced in item response theory and applications: An introduction. *Applied Psychological Measurement*, 6, 373–378.
- Kane, M. T. (2006). Validation. In R. L. Brennan (Ed.), *Educational measurement* (pp. 17–64). Washington, DC: American Council on Education/Praeger.
- Kane, M. T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50, 1–73.
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating: Methods and practice*. NY: Springer.
- Kuder, G. F., & Richardson, M. W. (1937). The theory of the estimation of test reliability. *Psychometrika*, 2, 151–160.
- Lane, S., & Stone, C. A. (2002). Strategies for examining the consequences of assessment and accountability programs. *Educational Measurement: Issues and Practice*, 21, 23–30.
- Linacre, J. M. (2012). *A User's Guide to WINSTEPS MINISTEP Rasch-Model Computer Programs*. Chicago, IL.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading MA: Addison-Wesley.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute*, 22, 719–748.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47, 149–174.
- McNeill, K. L., Katsh-Singer, R., & Pelletier, P. (2015). Assessing science practices: Moving your class along a continuum. *Science Scope*, 39, 21–28.
- Messick S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 13–104). New York, NY: American Council on Education.
- Miller, G. E., Rotou, O., & Twing, J. S. (2004). Evaluation of the 0.3 logits screening criterion in common item equating. *Journal of Applied Measurement*, 5(2), 172–177.

- National Research Council (2012). *A framework for K-12 science education: Practices, crosscutting concepts, and core ideas*. Washington, DC: National Academies Press.
- New Jersey Department of Education (Winter 2022). New Jersey Student Learning Assessment (NJSLA) accessibility features and accommodations manual: Guidance for districts and decision-making teams to ensure that NJSLA summative assessments produce valid results for all students.
[NJSLA NJGPA Accessibility Features and Accommodations 11th Edition 011223 V5 Final \(1\).pdf \(mypearsonsupport.com\)](#),
<https://nj.mypearsonsupport.com/resources/manuals/NJSLASpring2019AFA.pdf>
- New Jersey Department of Education (2020). *New Jersey Student Learning Assessment – Science (NJSLA–S): Technical report grades 5, 8, and 11 2019* [Technical report]. NJDOE.
- NGSS Lead States (2013). *Next Generation Science Standards: For states, by states*. The National Academies Press.
- Ostini, R., & Nering, M. L. (2010). New perspectives and applications. In M. L. Nering & R. Ostini (Ed.), *Handbook of Polytomous Item Response Models* (pp. 3–20). New York, NY: Routledge.
- Padilla, J. L., & Benitez, I. (2014). Validity evidence based on response processes. *Psicothema*, 26, 136–144.
- Penfield, R. (2007). An approach for categorizing DIF in polytomous items. *Applied Measurement in Education*, 20, 335–355.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen: Danish Institute for Educational Research.
- Sireci, S. G., Baldwin, P., Martone, A., Zenisky, A. L., Kaira, L., Lam, W., Shea, C. L., Han, K., Deng, N., Delton, J., & Hambleton, R. K. (2008). *Massachusetts Adult Proficiency Tests Technical Manual*. Boston, MA: Massachusetts Department of Elementary and Secondary Education.
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal Design Applied to Large Scale Assessments*. Synthesis Report.
- Traub, R. E., & Rowley, G. L. (2008). *Understanding reliability*. Instructional topics in educational measurement. Madison, WI: National Council on Measurement and Education 176–177.
- Wright B.D., & Linacre, J. M. (1994). Reasonable mean-square fit values. *Rasch Measurement Transactions*, 8:3, 370.
- Wright B.D., & Masters, G. N. (1982). *Rating Scale Analysis*, Chicago: MESA Press.
- Yen, W. M. (1984). Effects of local item dependence on the fit and equating performance of the three-parameter logistic model. *Applied Psychological Measurement*, 8, 125–145.

Zieky, M. (1993). Practical questions in the use of DIF statistics in item development. In P. W. Holland, & H. Wainer (Eds.), *Differential item functioning* (pp. 337–347). Hillsdale, NJ: Erlbaum.