

Raw Score to Scaled Score Conversions

Jon S. Twing, Ph.D.

Vice President, Psychometric Services

NCS Pearson - Iowa City

Personal Background

- **Doctorate in Educational Measurement and Statistics, University of Iowa.**
- **Responsible for psychometric and content development for NCS Pearson.**
- **Based in Iowa City - Measurement Services Division.**

Experience

- **WISC-III, Aprenda I, MAT-7, Stanford 9, WIAT, NNAT, CMEE, and others.**
- **Michigan (MEAP), Texas (TAAS, RPTE, SDAA), Minnesota (BST, MCA), Florida (FCAT).**
- **Advisory roles: Iowa DOE, Texas DOE, Virginia DOE, Federal Government.**

Raw Score to Scaled Score Conversions

Purpose

- What are “Scaled Scores”, why are they used and what are they good for?
- Gain a conceptual, if not specific understanding of scaled scores and augment understanding of assessment in general.

Goals and Agenda

- Everything starts with “item and test construction.”
- Field testing generates statistical data.
- Test form equating is used to ensure fairness in comparability.
- Standard setting conveys the meaning of what performance is required.
- Scaled scores help communicate student performance:
 - supports the reporting of performance standards across years;
 - facilitates equivalency of test forms across the years;
 - standardizes the meaning of performance across testing sessions;

Item Generation

- **Test questions (items) are generated to match various requirements as provided in the item development specifications:**
 - **What will be measured...the curriculum (Standards and Benchmarks):**
 - **Scope and sequence**
 - **Reporting objectives**
 - **New York State Learning Standards**
 - **How it will be measured...the specifications (Development Guidelines):**
 - **Format**
 - **Fairness**
 - **Is it Art or is it Science?**
 - **Expert judgement (measurement specialists)**
 - **Empirical evidence from field testing (psychometricians)**
 - **Educator review (teachers)**

Raw Score to Scaled Score Conversions

Field Testing

- **Newly developed test questions are field tested using real students to collect statistical information for a variety of purposes.**
 - **What are the psychometric / statistical properties of the items?**
 - **Difficulty, discrimination, response options.**
 - **Equating to existing item pools and/or test forms.**
 - **How can we fairly compare the statistical properties of all items in the pool?**
 - **How can we build test forms of comparable difficulty from year to year?**
 - **What evidence do we have that the items are fair?**
 - **Appropriate difficulty level.**
 - **Differential Item Functioning (DIF) analyses.**
 - **Educator review.**
 - **Ultimately the goal is to add items to the item pool.**

Raw Score to Scaled Score Conversions

Equating

- To better ensure that the items constructed, tested and used to build test forms this year are comparable to those used in previous years.
- Statistical test form equating is used to create a link between previous testing and current testing.
 - **Data collection designs (typically used during field testing):**
 - **Common item designs.**
 - **Randomly equivalent groups designs.**
 - **Usually implies a mathematical measurement model:**
 - **A mathematical function relating performance on items and tests to an underlying scale.**
 - **The Regents exams use the one-parameter Item-Response Theory measurement model known as the Rasch model.**
 - **Other multi-parameter models are also quite common.**

Raw Score to Scaled Score Conversions

Calibration

- Calibration is the process of relating student performance on test questions to statements about student ability.
 - This is done by describing the interaction between test items and student ability via a mathematical function or measurement model.
 - For the Regents examinations this is the Rasch measurement model.
 - Parameters of the mathematical function are estimated (usually via a computer) and used to obtain various derived scores (scaled scores) for items and students.

$$P_i(\theta) = \left[\frac{1}{1 + \exp^{-(\theta - b_i)}} \right]$$

- Derived scores facilitate test construction, equating, standard setting, fairness and longitudinal reporting.

An Example

- **Suppose we wanted to build an Algebra test.**
 - **Suppose you had test questions that were:**
 - **aligned to the curriculum;**
 - **reviewed for appropriateness and fairness.**
 - **Suppose further that a field testing had taken place and that the results of student performance on the field test were collected and placed into a data base for your use.**
 - **Finally, suppose we followed a statistical equating process such that all items in the pool were comparable to each other as well to other existing items already in the pool.**
- **How would you get started?**

Raw Score to Scaled Score Conversions

An Example Continued

- Typically in testing, a database of test questions is referred to as a test item bank:

The screenshot shows the 'Test Builder' application window. The title bar reads 'Test Builder'. The menu bar includes 'File', 'Search Items', 'Statistics', and 'Help'. The main window is divided into two panes: 'Available Items in Grade 09, Algebra' and 'Test Items'.

Available Items in Grade 09, Algebra

UIN	Pssg	M	W	N	R	D	P	R.Diff	P-Val	PBIS
09EA01AB2AZ00001								-1.767	88	0.45
09EA01AB2AZ00002					X			-1.441	84	0.50
09EA01AB2AZ00003						X		2.686	19	0.08
09EA01AB2AZ98001					X			-1.137	82	0.41
09EA01AB2AZ98002								-0.700	74	0.50
09EA01AB2AZ98003					X			-1.423	87	0.36
09EA01AB2AZ98004					X			-0.766	78	0.45
09EA01AB2AZ99003					X			-1.234	85	0.48
09EA01AB2BZ00007						X		1.612	36	0.28
09EA01AB2BZ00009								1.145	46	0.46
09EA01AB2BZ93321					X			-0.754	75	0.57
09EA01AB2BZ96603					X			-0.507	72	0.50
09EA01AB2BZ97702					X			-0.677	74	0.56
09EA01AB2BZ97706					X			0.201	57	0.54
09EA01AB2BZ97710					X			0.991	46	0.39
09EA01AB2BZ98005						X		1.743	29	0.26
09EA01AB2BZ98007						X		1.817	28	0.24
09EA01AB2BZ98008								1.747	30	0.42
09EA01AB2BZ98009								1.903	28	0.34
09EA01AB2BZ98010								1.724	29	0.38

Search Results

UIN	Pssg	M	W	N	R	D	P	R.Diff	P-Val	PBIS

Test Items

#	UIN	Pssg	R.Diff	P-Val	PBIS	Obj

At the bottom of the window, there is a summary area with the following data:

0	Items	Mean	0	0	0
		Sums	0	0	0

Buttons for 'Print', 'Search', 'Save', and 'Cancel' are visible at the bottom right, along with a trash icon.

Raw Score to Scaled Score Conversions

An Example Continued

- Suppose we select test items from the pool to build this Grade 9 Algebra Test from the resultant field test data:

The screenshot shows the 'Test Builder: Spring 1999' window. It displays two tables: 'Available Items in Grade 09, Algebra' and 'Test Items: 40 Records added'. The 'Available Items' table lists 20 items with columns for UIN, Pssg, M, W, N, R, D, P, R. Diff, P-val, and PBIS. The 'Test Items' table lists 40 selected items with columns for #, UIN, Pssg, R. Diff, P-val, PBIS, and Obj. A 'Search Results' table is also visible at the bottom left of the window.

UIN	Pssg	M	W	N	R	D	P	R. Diff	P-val	PBIS
09EA01AB2AZ00001								-1.767	88	0.45
09EA01AB2AZ00002					X			-1.441	84	0.50
09EA01AB2AZ00003						X		2.686	19	0.08
09EA01AB2AZ98001					X			-1.137	82	0.41
09EA01AB2AZ98002								-0.700	74	0.50
09EA01AB2AZ98003					X			-1.423	87	0.36
09EA01AB2AZ98004					X			-0.766	78	0.45
09EA01AB2AZ99003					X			-1.234	85	0.48
09EA01AB2BZ00007						X		1.612	36	0.28
09EA01AB2BZ00009								1.145	46	0.46
09EA01AB2BZ93321				X				-0.754	75	0.57
09EA01AB2BZ96603				X				-0.507	72	0.50
09EA01AB2BZ97702				X				-0.677	74	0.56
09EA01AB2BZ97706				X				0.201	57	0.54
09EA01AB2BZ97710				X				0.991	46	0.39
09EA01AB2BZ98005					X			1.743	29	0.26
09EA01AB2BZ98007					X			1.817	28	0.24
09EA01AB2BZ98008								1.747	30	0.42
09EA01AB2BZ98009								1.903	28	0.34
09EA01AB2BZ98010								1.724	29	0.38

#	UIN	Pssg	R. Diff	P-val	PBIS	Obj
1	09EA08AC4CZ97797		-1.519	82	0.40	08
2	09EA08AC3CZ95507		-1.485	80	0.42	08
3	09EA09AB2DZ97704		-1.124	76	0.40	09
4	09EA09AC2GZ97711		-1.116	78	0.44	09
5	09EA06AB4AZ97717		-0.801	75	0.52	06
6	09EA04AC4AZ97705		-0.666	72	0.53	04
7	09EA01AC2CZ96607		-0.413	62	0.51	01
8	09EA04AC4AZ94401		-0.338	59	0.51	04
9	09EA02AB2BZ97712		-0.247	66	0.53	02
10	09EA01AB2BZ97702		-0.161	62	0.57	01
11	09EA01AC2CZ97711		-0.003	58	0.56	01
12	09EA02AB2BZ97713		0.025	51	0.56	02
13	09EA01AB2BZ93321		0.054	51	0.51	01
14	09EA03AC1CZ95506		0.073	52	0.56	03
15	09EA06AB4AZ97706		0.202	56	0.60	06
16	09EA06AB4AZ97716		0.370	57	0.60	06
17	09EA03AC1CZ93310		0.467	49	0.59	03
18	09EA03AC1CZ96625		0.358	51	0.60	03
19	09EA07AD2AZ97704		0.381	53	0.54	07
20	09EA02AC2DZ97717		0.388	57	0.52	02
21	09EA05AD2AZ97717		0.397	49	0.55	05
22	09EA04AC4AZ97709		0.414	59	0.51	04
23	09EA05AD2AZ93353		0.453	42	0.55	05
24	09EA02AB2CZ95516		0.511	41	0.40	02
25	09EA08AC4CZ97707		0.513	54	0.58	08
26	09EA07AD2AZ96603		0.449	52	0.68	07
27	09EA08AC4CZ95503		0.507	42	0.59	08
28	09EA03AC1CZ97727		-0.003	63	0.57	03
40	Items		Mean	0.019	56.2	0.51
			Sums	0.763	2246	20.39

An Example Continued

- How did we do it?
 - Start with one objective at a time.
 - Use the last “live” assessment for “target” statistical values.
- Don’t forget, content is more important than the statistical parameters when selecting items to place on a test.
 - “...content must rule the day.”
- This is what is often referred to as the combination of “art and science” in test construction.
 - Obviously, this requires highly skilled individuals and is a painstaking process.

Test Statistics			
Test Statistics by Objective			
Obj	Rasch Dif	Pval	PBIS
01	-0.131	58	0.54
02	0.169	54	0.50
03	0.224	54	0.58
04	0.330	55	0.50
05	0.105	51	0.49
06	-0.092	63	0.55
07	-0.130	57	0.57
08	0.088	53	0.46
09	-0.458	66	0.45
All	0.019	56	0.51

UIN	Rasch Dif	Pval	PBIS
09EA01AC2CZ96607	-0.413	62	0.51
09EA01AB2BZ97702	-0.161	62	0.57
09EA01AC2CZ97711	-0.003	58	0.56
09EA01AB2BZ93321	0.054	51	0.51

Mean:	-0.131	58.3	0.54
Sum:	-0.523	233	2.15

An Example Continued

- The results of the field-testing can be very extensive depending upon the requirements:
 - DIF
 - History
 - Multiple Forms
 - Reviewer Comments
 - Actual Item Image
- As such, each builder must consider a lot of statistical and evaluative information.

Item Statistical Details

Item UIN 09EA01AB2AZ98002 Grade 09 Subject A

Objective: 01 Skill: B2 Test date: Passage: Do Not Use Flags:
 Released (R) Passage (P) Data Review (D) Misc (M) Note (N) Warning (W)

Expectations: A Sub: Z Sub-Sub: Sub-Sub-Sub: Comments: Restricted Comments:

Statistical History

#	Type	Test Date	Form	V.Const	R.Diff	P-val	PBIS	Mantel-Haenszel Alpha			ANS
								A-W	H-W	F-M	
22	PT	05/1999	04	0	-0.700	74	0.500	1.2	1.0	0.9	3

Current Statistics

Rasch Diff: -0.700
P-value: 74
PBIS: 0.50

Demographical Statistical History

Group	A,F	B,G	C,H	D,J	E,K	Other	Omit	P-Val	PBIS	R.Diff	R.Fit
Composite	9	4	74	8	5	0	0	74	0.50	-0.700	1.078
AfricanAmerican	17	8	53	13	8	1	0	53	0.52	-0.702	0.905
Female	9	3	75	8	4	0	0	75	0.47	-0.830	1.164
Hispanic	10	4	71	10	6	0	0	71	0.47	-0.901	0.919
Male	9	4	73	9	5	0	0	73	0.52	-0.795	0.984
White	7	3	80	6	3	0	0	80	0.45	-0.794	1.152

Technical Comments: Close

Raw Score to Scaled Score Conversions

An Example Continued

- Because we are using the Rasch measurement model, the item difficulty estimates (R.Diff) and the person ability estimates (theta or θ) are on the same scale (logistic scale).
 - This makes direct comparisons between student skills and test items possible.
 - A student with a theta value = 1.00 has a probability of 0.50 of answering a test item with a difficulty of 1.00 correctly.
 - As the student ability increases and/or the test item difficulty decreases, this probability will go up.
 - Similarly, as the student ability decreases and/or the test item difficulty increases, this probability will go down.

Raw	Ability	Scale	PR	TLI
0	-4.922	0920	0	0
1	-3.903	1020	0	0
2	-3.168	1100	0	0
3	-2.721	1140	0	0
4	-2.390	1170	0	0
5	-2.125	1200	0	0
6	-1.899	1220	0	0
7	-1.701	1240	0	0
8	-1.523	1260	0	0
9	-1.360	1270	0	0
10	-1.208	1290	0	0
11	-1.066	1300	0	0
12	-0.931	1320	0	0
13	-0.802	1330	0	0
14	-0.677	1340	0	0
15	-0.556	1350	0	0
16	-0.439	1360	0	0
17	-0.323	1380	0	0
18	-0.209	1390	0	0

An Example Continued

- In fact, this is the probability returned in the mathematical conversion formula provided by the Rasch model as already discussed:

$$P_i(\theta) = \left[\frac{1}{1 + \exp^{-(\theta - b_i)}} \right]$$

- The calibration process provides both the required theta values (θ) as well as the item difficulty values (δ)
 - Because a statistical equating process was followed, these values are comparable not only to each other, but to other items (and students) who have been tested before.
 - This means that the test we just constructed is comparable to those tests previously used...that is, this test has been pre-equated.

Raw Score to Scaled Score Conversions

Something Exciting has Just Happened!

- So what...what does all this mean to me? Think about and summarized what has just transpired:
 - Items were constructed measuring the New York State Learning Standards.
 - These items were reviewed for fairness and appropriateness.
 - These newly constructed items were included in a “small” field test.
 - Many thousands of students...but not all...and each student was not required to take all items.
 - This field test included link items or used some other methodology to statistically equate the newly constructed items as well as the student performances on these items to previous tests that actually counted.
 - A statistical equating process was followed allowing for comparability across the years.
 - We know what the passing standard was for these previously used tests.
 - Since the new items have been equated to these previously used tests we can also know what the passing standard is on the new form we just constructed:
 - » **BUT REMEMBER...NO ONE HAS TAKEN THE TEST YET!!!**

Equivalent Test Forms via Pre-Equating

- The test we constructed is equivalent in difficulty to the test forms previously used.
 - The same passing standard established on the previous tests will be the passing standard on the current test.
- Suppose we know that the previous passing standard (the one established on the previous test form) was 65 percent of the items.
 - This process usually comes about from a formal “Standard Setting”
 - The Regents exams use “item mapping” standard setting methodology.
 - Assume that this 65 percent was 24 items out of 40 and that a Rasch ability value of 1.50 was associated with this value of 24.
- Even though we have constructed this new form to be equivalent, random error may mean that the best we can do is to get within one raw score point. So, assume that our test is one raw score harder and that this is the best we can do.
 - We would know this...and actually do know it...as a result of the equating.

Raw Score to Scaled Score Conversions

A Derived Score (Scaled Score) is Needed

- This means that the equivalent passing standard on our test (which is one score point harder than last year's) is 23.
 - In other words, a 23 / 40 on our test (which is one point harder than last year's test) is equivalent to last year's passing standard (on a test one item easier) which was 24 / 40.
 - How many of you would like to try to explain this to your Board of Education or a parent...that 24 / 40 is equivalent to 23 / 40?
 - Yet, we know that 24 / 40 last year was equal to a Rasch theta value of 1.50. And, if you believe me when I say that we know the test we constructed is one item more difficult, then 23 / 40 on our test will have a theta value of 1.50!

Last Year's Test		Our Test	
Total Score at 65% Passing Standard	Rasch Theta Value	Rasch Theta Value	Equivalent 65% Total Score Passing Standard
24 / 40	1.50	1.50	23 / 40

A Derived Score (Scaled Score) is Needed

- Clearly, if we could communicate the 65% passing standard in terms of the Rasch theta value then we are likely to be more successful since this value is the same from year to year (and will always be the true passing standard).
 - Unfortunately, the decimilized signed metric of the theta scale is not very convenient...no one likes decimals and negative numbers!
 - So, we usually do a convenience scaling, which is nothing more than another mathematical conversion to a scale that is easier to use.
 - For the Regents exam, this transformation is:

$$\text{Scaled Score} = (a) * (x ** 3) + (b) * (x ** 2) + (c) * (x) + d$$

where

x is the theta value of a student and a, b, c, and d are parameters found by solving for simultaneous equations where 0 is the scaled score associated with the lowest theta value, 65 is the scaled score associated with the “passing” standard, 85 is the scaled score associated with the “passing with distinction” standard and 100 is the scaled score associated with the highest theta value.

Raw Score to Scaled Score Conversions

Fun with Scaled Scores

- Do you remember the conversion table we generated from our test...based on the pre-equating?
- If we had used the Regents conversion we would have known what the passing scores were and which raw score was associated with each one (0, 65, 85 and 100).
- Furthermore, we could go back and add or remove items to construct a test to get the “correct” raw score associated with each scaled score.
 - For example, we could go back and remove a hard item and add an easy item to see if we could build a test where a scaled score of 65 was equal to 24 and not 23.

Test Item : Pre-equating

Pre-Equating for Test: Spring 1999

Raw	Ability	Scale	PR	TLI
0	-4.922	0920	0	0
1	-3.903	1020	0	0
2	-3.168	1100	0	0
3	-2.721	1140	0	0
4	-2.390	1170	0	0
5	-2.125	1200	0	0
6	-1.899	1220	0	0
7	-1.701	1240	0	0
8	-1.523	1260	0	0
9	-1.360	1270	0	0
10	-1.208	1290	0	0
11	-1.066	1300	0	0
12	-0.931	1320	0	0
13	-0.802	1330	0	0
14	-0.677	1340	0	0
15	-0.556	1350	0	0
16	-0.439	1360	0	0
17	-0.323	1380	0	0
18	-0.209	1390	0	0

Export Print Close

Raw Score to Scaled Score Conversions

Fun with Scaled Scores

- What would happen without such scaled scores?
- Consider the table below:

Last Year's Test		Our Test	
Total Score at 65% Passing Standard	Rasch Theta Value	Rasch Theta Value	Equivalent 65% Total Score Passing Standard
.	.	.	.
.	.	.	.
.	.	.	.
23 / 40	1.45	1.47	22 / 40
24 / 40	1.50	1.50	23 / 40
25 / 40	1.55	1.56	24 / 40
.	.	.	.
.	.	.	.
.	.	.	.

- If we mistakenly used 24 as the passing standard for both years we can see that the passing standard would be a theta value of 1.50 one year and 1.56 the next year
 - In other words...it would be harder to pass the second year without the scale!

Fun with Scaled Scores

- **Measurement standards and best practice tell us some things:**
 - **Validity of score use is completely tied to the meaning and interpretability of scores generated from a testing occasion.**
 - **Builders of assessments are burdened with facilitating the valid interpretation and use of scores resulting from an assessment.**
- **Experience tells us some things:**
 - **Teachers, parents and the public alike think they know and understand simple scores like total scores, percent correct and perhaps percentile ranks...when they don't.**
 - **Teachers, parents and the public rarely claim they know and understand how to interpret scaled scores...but they need to.**
 - **Most large-scaled assessment systems use scaled scores and provide training around the use and interpretation of such scores.**
 - **Still, these scores are not embraced by the majority of the public.**
 - **We, as leaders in education, must be diligent in our explanations of derived scores such as scaled scores.**

Conclusions:

- Scaled scores result from one step of a multi-step test construction process.
- When used with thoughtful item development, standard setting, field testing and statistical test form equating, scaled scores:
 - Facilitate fair communication of student performance relative to established passing standards;
 - Allow for easy to understand comparability of student performance from year to year;
 - Provide for additional steps in the test development process to build forms parallel in both content and statistical parameters.
- Scaled scores suffer from the same misinterpretations as do raw scores but have the perception of being more difficult to understand.
 - People misinterpret raw scores but fail to recognize it...most people recognize, that they fail to understand scaled score interpretations.