

Validating Measures of Content Knowledge for Teaching Mathematics: A Validity Argument Approach

Deborah L. Ball

Carolyn Dean

Imani Masters Goffney

Stephen Schilling

Merrie L. Blunk

Seán Delaney

Heather Hill

Laurie Sleep

Chia-Ning Wang

University of Michigan

American Education Research Association Annual Meeting, Montréal

Session 63:039

April 14, 2005

<http://www.soe.umich.edu/lmt/>



Overview of Session

1. Brief description of Learning Mathematics for Teaching (LMT) Project
2. Define validation argument and lay out LMT validation work
3. Examples of validation work in LMT
 - Evidence from interviews with teachers
 - Evidence from interviews with mathematicians
 - Evidence from statistical analyses
4. Commentary by discussants: **Michael Kane, Mark Reckase, and Anna Sfard**

<http://www.soe.umich.edu/lmt/>



2

Validating Measures of Content Knowledge for Teaching Mathematics: A Validity Argument Approach

Merrie L. Blunk
Stephen G. Schilling

University of Michigan

<http://www.soe.umich.edu/lmt/>



3

Overview of Learning Mathematics for Teaching Project

Goal: Measure teachers' mathematical knowledge for teaching

- Multiple choice
- Particular content domains (number and operations; patterns, functions and algebra; geometry)
- Items meant to assess different kinds of mathematical knowledge:
 - Common and specialized content knowledge
 - Knowledge of students and content

<http://www.soe.umich.edu/lmt/>



4

Example:

Specialized Content Knowledge

Student A	Student B	Student C
$\begin{array}{r} 35 \\ \times 25 \\ \hline 125 \\ +75 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 175 \\ +700 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 25 \\ 150 \\ 100 \\ +600 \\ \hline 875 \end{array}$

Which of these students is using a method that could be used to multiply any two whole numbers?

Example:

Knowledge of Students and Content

Ms. Violeta was looking carefully at her students' papers, and she saw the following responses to the problem:

$$8 + 4 = \underline{\quad} + 5$$

- i) 12 ii) 17 iii) Can't do it iv) 1

Which of the following is the most likely explanation of the difficulty the students are having? (Mark ONE answer.)

- a) They do not know their basic addition facts.
- b) They cannot do multi-step problems.
- c) They do not know that addition is commutative.
- d) They do not understand the meaning of the equals sign.

Background on Validation

- In general, test validation focuses on two questions:
 - Do test scores supply the desired information?
 - Are the test scores helpful in making good decisions?
- Absence of validation guidelines has encouraged the use of selective evidence (Cronbach, 1988; Messick, 1975, 1981)
- Kane's argument-based approach attempts to provide a methodology for validation.

<http://www.soe.umich.edu/lmt/>



7

Kane's Argument-Based Approach

First Stage – Making explicit your desired interpretations of the scores

- Also: make explicit what decisions you intend to make from the test scores.

Second Stage – Seeking evidence for and against these as reasonable interpretations

- Focus on the most questionable assumptions.

<http://www.soe.umich.edu/lmt/>



8

Intended Use of Our Measures

- Uses of our measures
 - NOT used in high-stakes (licensing, certification, tenure)
 - For use in research and evaluation

<http://www.soe.umich.edu/lmt/>



9

The Design of Our Validation Work

- How do we want to be able to interpret teachers' performance on our questions?
 1. Teachers' scores capture teachers' mathematical knowledge
 2. Higher teacher scores are related to higher-quality mathematics instruction
 3. Higher teacher scores are related to improved student learning
 4. Teachers' scores reflect different dimensions of Content Knowledge for Teaching Mathematics
- Investigate interpretation with multiple sources of evidence

<http://www.soe.umich.edu/lmt/>



Multiple Sources of Evidence to Evaluate Our Claims

1. **Scores capture teachers' mathematical knowledge**
 - Cognitive interviews
2. **Higher scores are related to higher-quality mathematics instruction**
 - Videotape validation study
3. **Higher scores are related to improved student learning**
 - Study of Instructional Improvement student gains analysis
4. **Scores reflect different dimensions of “content knowledge for teaching mathematics”**
 - Mathematician and non-teacher interviews
 - Item response theory and factor analysis

<http://www.soe.umich.edu/lmt/>



Evidence for Claims

I. Scores capture mathematical knowledge

Evidence: Cognitive interviews

- Interviewed teachers
- Probed for reasoning behind their answers given
- Coded reasoning

Evidence for Claims

2. Higher scores are related to higher-quality mathematics instruction

Evidence: Videotape validation study

- Recruited 10 teachers, videotaped their mathematics teaching before and after professional development in mathematics
- Pretest and posttest
- Coding videotapes for teachers' mathematical knowledge and their use of mathematics with students
- Search for correspondence between teachers' scores and mathematical quality of their teaching

<http://www.soe.umich.edu/lmt/>



13

Evidence for Claims

3. Higher scores are related to improved student learning

Evidence: Study of Instructional Improvement student gains analysis

- Used students' gain scores (across on year)
- Compared students' achievement to teachers' scores
- Finding: Gains in student achievement were predicted by teacher scores

Evidence for Claims

4. Scores reflect different dimensions of Content Knowledge for Teaching Mathematics

Evidence: Mathematician and non-teacher interviews

- Interviewed non-teachers and mathematicians
- Probed for reasoning behind their answers given
- Coded answers and reasoning for mathematical correctness, consistency, and type of reasoning

Evidence: Item response theory and item factor analysis

<http://www.soe.umich.edu/lmt/>



15

Examples of our Validation Work

- Evidence from cognitive interviews
 - Imani Masters-Goffney
- Evidence from mathematicians
 - Laurie Sleep & Sean Delaney
- Evidence from item response theory and factor analysis
 - Stephen G. Schilling

<http://www.soe.umich.edu/lmt/>



Assessing “Content Knowledge for Teaching”: Data from Teachers, Non-teachers and Mathematicians

Heather Hill

Carolyn Dean

Imani Masters Goffney

University of Michigan

<http://www.soe.umich.edu/lmt/>



17

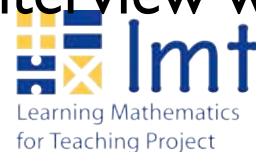
Rationale & Research Questions

- A. Are our measures a valid assessment of teachers' level of content knowledge for teaching mathematics?
- B. Using Kane to structure our investigation we analyzed two particular interpretations of scores:
- Do individuals' correct/incorrect responses capture their mathematical knowledge?
 - Do scores represent different dimensions of mathematical knowledge for teaching?

Sample Characteristics & Instrument Description

- Sample Characteristics
 - Teachers: $n = 26$
 - Non-teachers: $n = 21$
 - Mathematicians: $n = 18$
- Interview instrument
 - Standardized interview protocol
 - “Why did you choose (answer?) What process did you go through to decide? Did you consider any other answers? If so, why did you not select that answer?”
 - In the case of “knowledge of students and content” items, follow-up questions included a probe about whether the individual had ever experienced children making a similar error.
 - 17 items
 - Average length of interview was 50-75 minutes

<http://www.soe.umich.edu/lmt/>



19

Providing Evidence for Our Interpretations

1. Do individuals' correct/incorrect responses reflect their mathematical thinking?
 - Coded answers as consistent or inconsistent with thinking
2. Do scores represent different dimensions?
 - Reasoning process coded for all interview responses.

1. Mathematical justification
2. Mathematical reasoning
3. Definitions
4. Examples/ Counter examples
5. Pictures
6. Memorized rules

KSC

1. Knowledge of students

Non-mathematical/ Non-students

1. Test Taking Skills
2. Guessing
3. Other

Results I:

Consistency of Answers and Thinking

Consistency:

- Correct reasoning underlies correct answer
- Incorrect reasoning underlies incorrect reasoning

Inconsistency:

- Correct answer selected based on faulty mathematical reasoning
- Incorrect answer selected, but actually possesses correct understanding of the underlying mathematical concept

<http://www.soe.umich.edu/lmt/>



Table I: Consistency Rates

		Consistent	Inconsistent
CK items only	Teachers	94.1	5.9
	Mathematicians	95.8	4.2
	Non-teachers	91.2	8.8
	Average	93.7	6.3
KSC items only	Teachers	85.4	14.6
	Mathematicians	96.3	3.7
	Non-teachers	93.1	6.9
	Average	91.6	8.4

<http://www.soe.umich.edu/lmt/>



Flawed Item

Statement: A cube has eight edges.

(Always true, sometimes true, never true)

R: A cube has eight sides, I said that that was never true.

I: Ok, and how do you know that?

R: I counted them up and there was no way to get anything else but six.

Results II: Reasoning about items - CK

	Mathematical justification	Definitions	Examples	Mathematical reasoning	Pictures	Memorized rules	Guessing	Test-taking
Teachers	24.5	3.3	9.6	31.3	2.5	13.2	1.6	3.4
Mathematicians	35.8	6.4	6.4	33.8	1.2	7.4	0	2.4
Non-teachers	22.6	4.0	9.4	33.9	5.1	13.2	1.1	2.3
Average	26.9	4.3	8.7	32.7	2.9	11.7	1.0	2.8

(Results reported in percents)

<http://www.soe.umich.edu/lmt/>



Results II: Reasoning about Items - CK

	Mathematical justification	Definitions	Examples	Mathematical reasoning	Pictures	Memorized rules	Guessing	Test-taking
Teachers	24.5	3.3	9.6	31.3	2.5	13.2	1.6	3.4
Mathematicians	35.8	6.4	6.4	33.8	1.2	7.4	0	2.4
Non-teachers	22.6	4.0	9.4	33.9	5.1	13.2	1.1	2.3
Average	26.9	4.3	8.7	32.7	2.9	11.7	1.0	2.8

(Results reported in percents)

<http://www.soe.umich.edu/lmt/>



Results III: Reasoning about Items - KSC

	Students	Mathematical reasoning	Pictures	Guessing	Test taking
Teachers	40.9	40.9	4.9	0	16.3
Mathematicians	2.0	57.8	0	1.0	34.3
Non-teachers	15.9	59.5	0	2.4	21.9
Average	19.6	50.3	.5	.9	22.3

(Results reported in percents)

<http://www.soe.umich.edu/lmt/>



Overall Conclusions

1. Our interpretation of teachers scores on our assessment do measure some aspects of teachers' mathematical knowledge.
2. Evidence from analyses of KSC items suggests that, while weak, our interpretation of multidimensionality is valid.
3. Improving our future items will help in our efforts to more accurately measure the KSC construct.

Validating “Mathematics Knowledge for Teaching”: Evidence from Mathematicians’ Performance on Teacher Knowledge Items

Laurie Sleep

Seán Delaney

Carolyn Dean

Deborah Loewenberg Ball

Heather Hill

University of Michigan

<http://www.soe.umich.edu/lmt/>



28

The Multiple Sources of Evidence to Evaluate Our Claims

Claims about teachers' performance on our assessment:

1. Scores capture their mathematical knowledge
 - Cognitive interviews
2. Higher scores mean higher-quality mathematics instruction
 - Videotape validation study
3. Higher scores related to improved student learning
 - Study of Instructional Improvement student gains analysis
4. Scores represent different dimensions of Content Knowledge for Teaching Mathematics
 - Mathematician and non-teacher interviews
 - Item response theory and factor analysis

Data Collection & Analysis

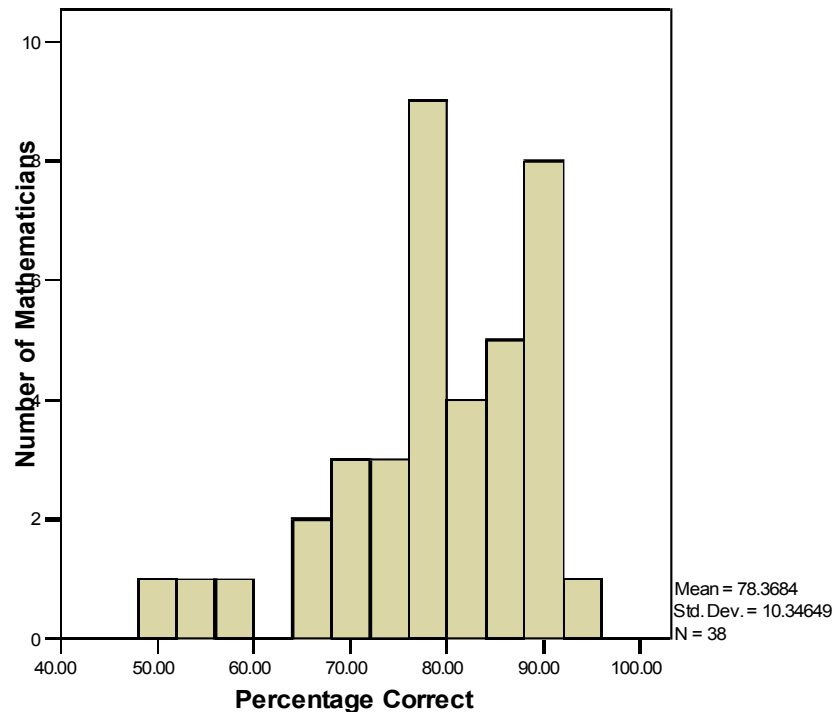
- Administered survey and interview to research mathematicians
 - 38 surveys, 18 interviews
- Coded interviews
 - right/wrong; consistent/inconsistent
 - type of explanations given
 - reasons for mathematicians' incorrect answers
- Ran descriptive statistics on survey data

<http://www.soe.umich.edu/lmt/>



Results

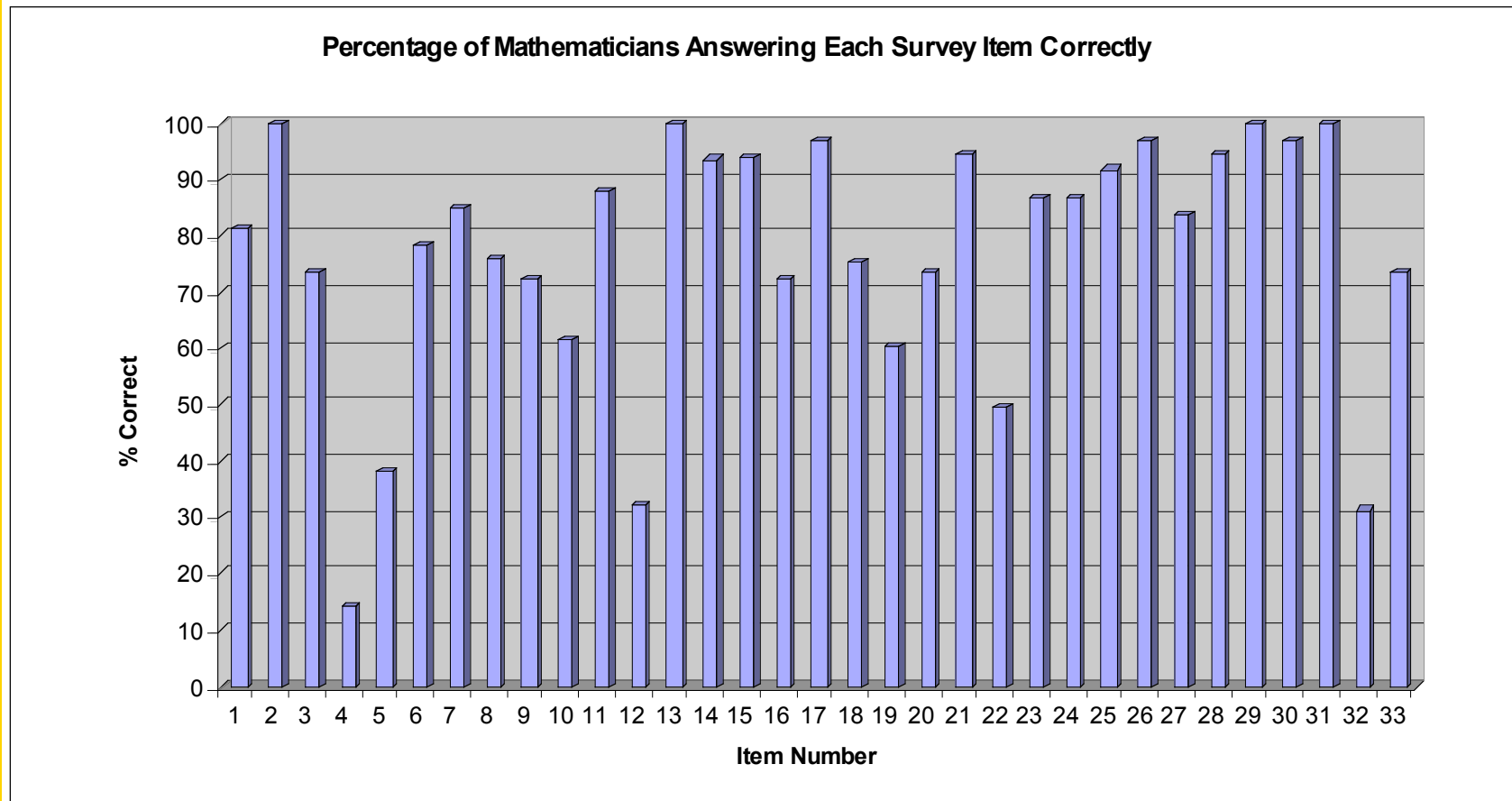
Mathematicians' Scores on Mathematical Knowledge for Teaching Number Items



Overall, mathematicians did well on the items:

- percent correct ranged from 48% to 94%
- mean score of 78%

Percentage of Mathematicians Answering Each Survey Item Correctly



Items Containing Flaws

Mr. Lewis was not surprised by this as he had seen students do this before. What did Mr. Lewis know?

“OK, so this one I answered (e), ‘I'm not sure’, because I don't know him, and the question says, ‘What did he know?’”

(Morin, interview, #3)

Items Requiring Knowledge of Students

Ms. Violeta was looking carefully at her students' papers, and she saw the following responses to the problem:

$$8 + 4 = \underline{\quad} + 5$$

- a) 12 b) 17 c) Can't do it d) 1

Which of the following is the most likely explanation of the difficulty the students are having? (Mark ONE answer.)

- a) They do not know their basic addition facts.
- b) They cannot do multi-step problems.
- c) They do not know that addition is commutative.
- d) They do not understand the meaning of the equals sign.

Items Requiring Mathematical Knowledge Unique to the Work of Teaching

Which of these students is using a method that could be used to multiply any two whole numbers?

Student A	Student B	Student C
$\begin{array}{r} 35 \\ \times 25 \\ \hline 125 \\ +75 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 175 \\ +700 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 25 \\ 150 \\ 100 \\ +600 \\ \hline 875 \end{array}$

Lack of Flexibility with Non-Standard Approaches

“Because there is one standard procedure for multiplication or division or subtraction done with a pen and a paper and if a person does something different first of all for me it is incorrect to begin with. So there is one way, the person should know this way, that's how it should be done, and also it's really difficult to figure out what the person was thinking.”

(Manchester, interview, #7)

Student A	Student B	Student C
$\begin{array}{r} 35 \\ \times 25 \\ \hline 125 \\ +75 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 175 \\ +700 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 25 \\ 150 \\ 100 \\ +600 \\ \hline 875 \end{array}$

Lack of Flexibility with Non-Standard Approaches

“Well, Method A, I didn’t know what ... was going on. Method B looks great, clearly, and Method C, I just wasn’t sure either. I didn’t know what was going on there too. 25, 50—where ... did he get all this? How could any student do a thing like that? I mean it seems absurd to me.”

(Marshall, interview, #7)

Student A	Student B	Student C
$\begin{array}{r} 35 \\ \times 25 \\ \hline 125 \\ +75 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 175 \\ +700 \\ \hline 875 \end{array}$	$\begin{array}{r} 35 \\ \times 25 \\ \hline 25 \\ 150 \\ 100 \\ +600 \\ \hline 875 \end{array}$

Compressed Mathematical Knowledge

Which of the following story problems could be used to illustrate $1 \frac{1}{4}$ divided by $\frac{1}{2}$?

b) You have \$1.25 and may soon double your money. How much money would you end up with?

“It's reflecting the fact that it's going to be half of something else, meaning you're going to double it. So it's just reflecting the fact that division by a half is multiplication by two, so that's why that one is correct.” (McIntyre, interview, #6b)

Conclusions

- Mathematicians' correct answers help validate that items are measuring mathematical knowledge.
- There is mathematical knowledge specialized to the work of teaching.
- Mathematicians also correctly answered many of these items.
 - Raises questions about both item design and the extent to which expert knowledge and reasoning abilities can compensate for unfamiliarity with the work of teaching

Validating Content Knowledge for Teaching Mathematics Using Unidimensional and Multidimensional Item Response Theory

Stephen Schilling
Chia-Ning Wang

University of Michigan

<http://www.soe.umich.edu/lmt/>



40

Piloting Items

- Three forms: A, B, C
- Content knowledge (CK) and knowledge of students and content (KSC)
- Linking items across forms

<http://www.soe.umich.edu/lmt/>



Interpretive Claims

1. CKT-M consists of two distinct constructs – CK and KSC – and these constructs are measured by the items as organized in these two separate domains.
2. Teachers can be reliably distinguished by IRT scores reflecting this organization by types of knowledge.
3. IRT scores do not vary across different samples of items reflective of these constructs.

Note: Claims are interconnected and mutually supporting.

<http://www.soe.umich.edu/lmt/>

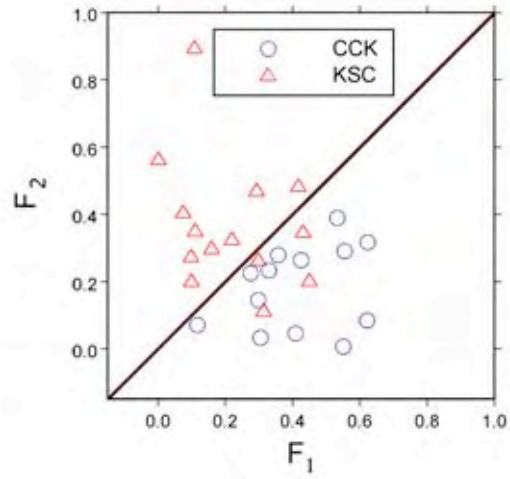


Method for Evaluating Claim #1

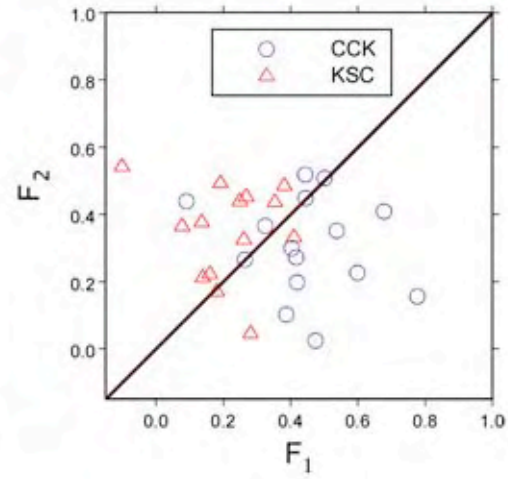
- I. CKT-M consists of two distinct constructs –CK and KSC – and these constructs are measured by the items as organized in these two separate domains.

Method: Full-information item factor analysis

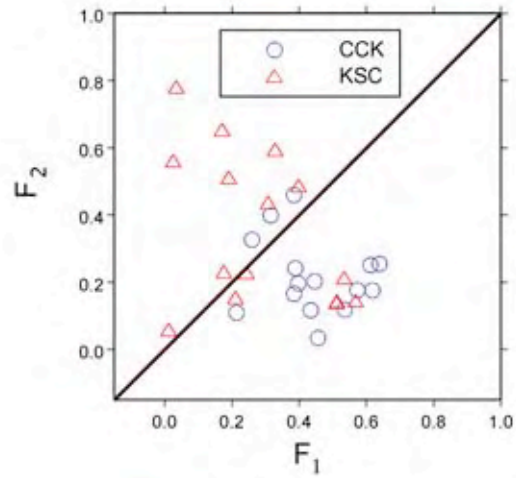
Form A Factor Loadings



Form B Factor Loadings



Form C Factor Loadings



Method for Evaluating Claim #2

2. Teachers can be reliably distinguished by IRT scores reflecting this organization by types of knowledge.

Method: Unidimensional IRT models

- marginal reliability of each scale
- test information plots

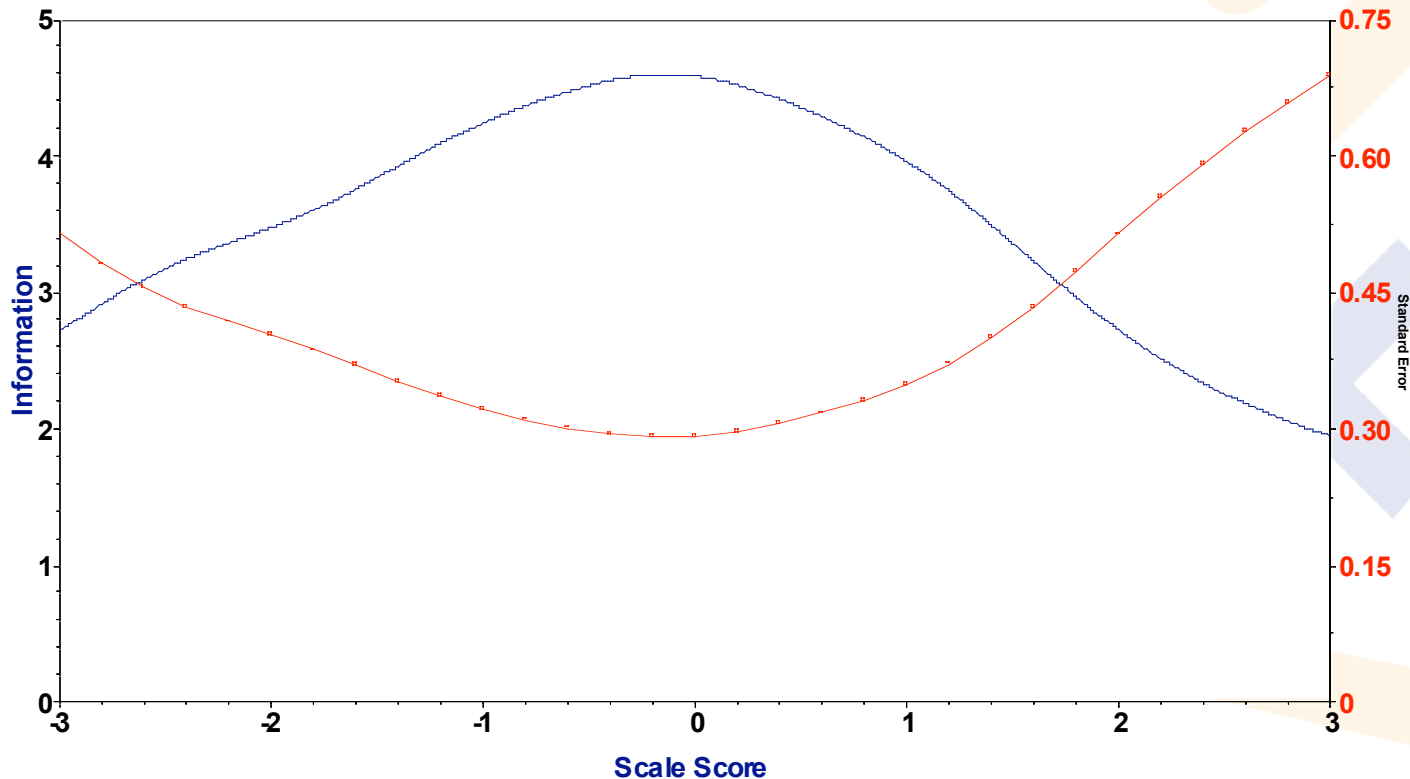
<http://www.soe.umich.edu/lmt/>



Test Information - Form A

Content Knowledge

Test Information and Measurement Error - Form A Content Knowledge



Reliability = 0.81

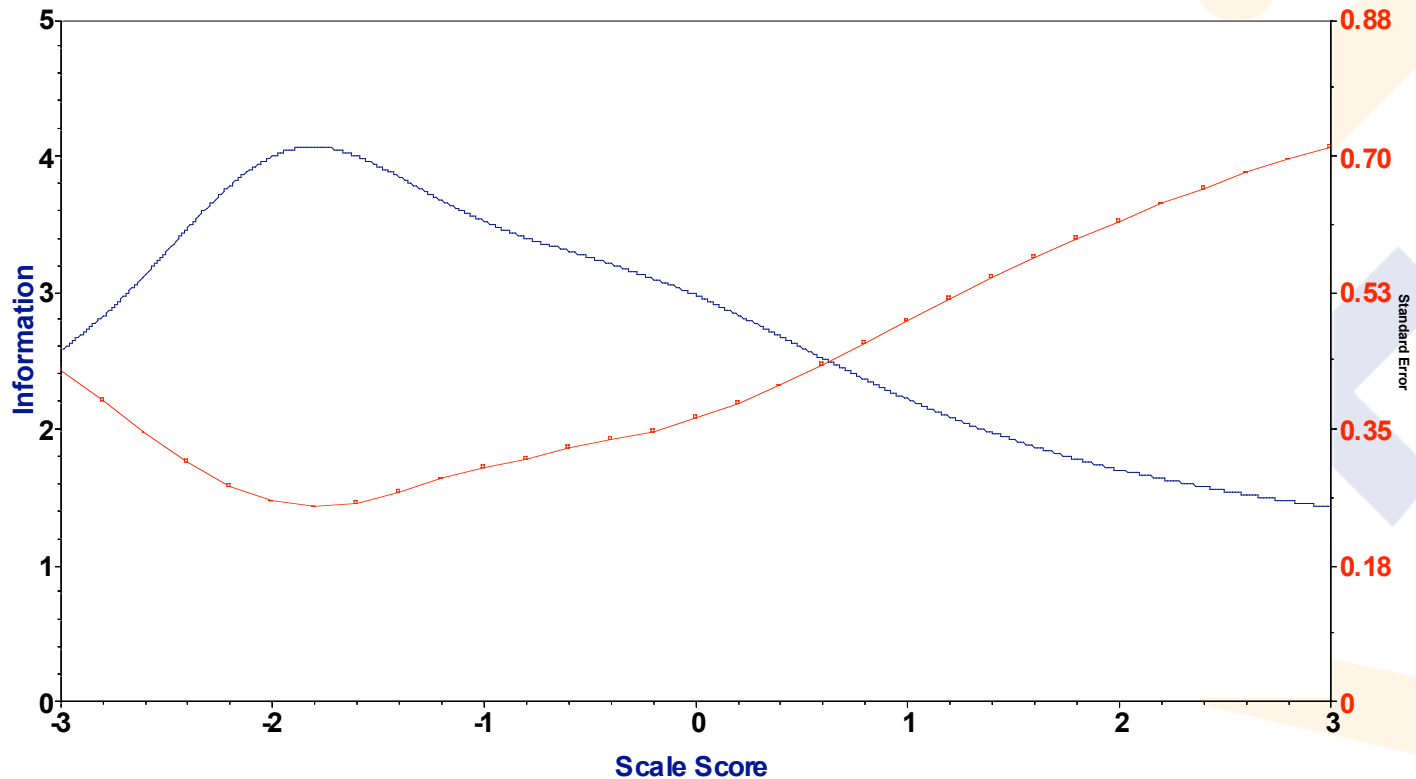
<http://www.soe.umich.edu/lmt/>



Test Information - Form A

Knowledge of Students and Content

Test Information and Measurement Error - Form A KSC



Reliability = 0.74

<http://www.soe.umich.edu/lmt/>



Reliability and Point of Maximum Information - All Forms

	Content Knowledge		Knowledge of Students and Content	
	IRT reliability	Point of maximum information	IRT reliability	Point of maximum information
Form A	0.81	-0.1	0.74	-1.8
Form B	0.83	-1.1	0.74	-1.4
Form C	0.81	-0.3	0.74	-1.3

<http://www.soe.umich.edu/lmt/>



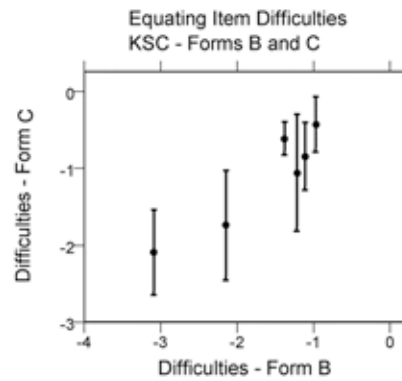
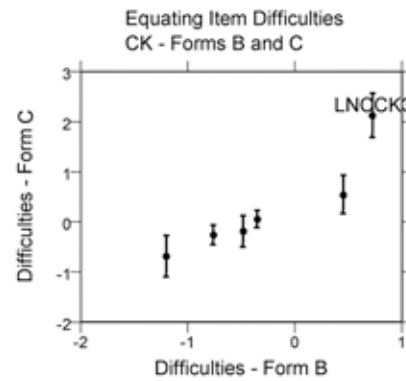
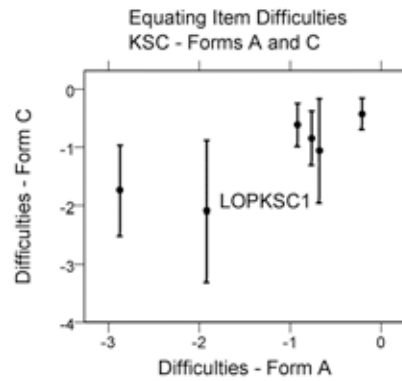
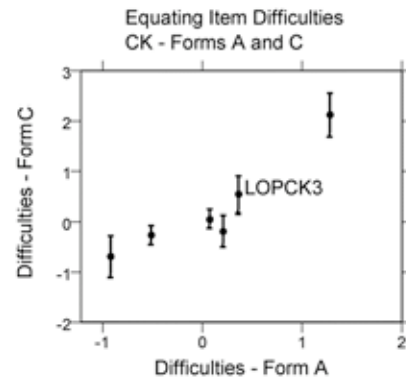
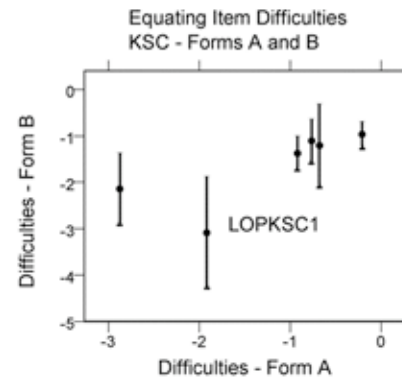
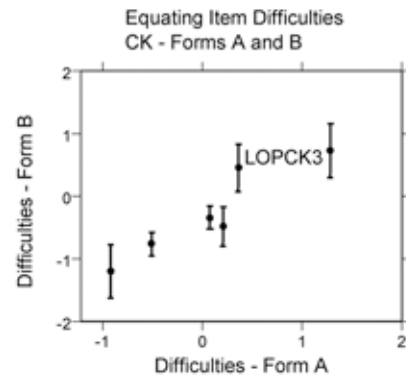
Method for Evaluating Claim #3

3. IRT scores do not vary across different samples of items reflective of these constructs.

Method: Unidimensional IRT models
- evaluating equating of scales

<http://www.soe.umich.edu/lmt/>





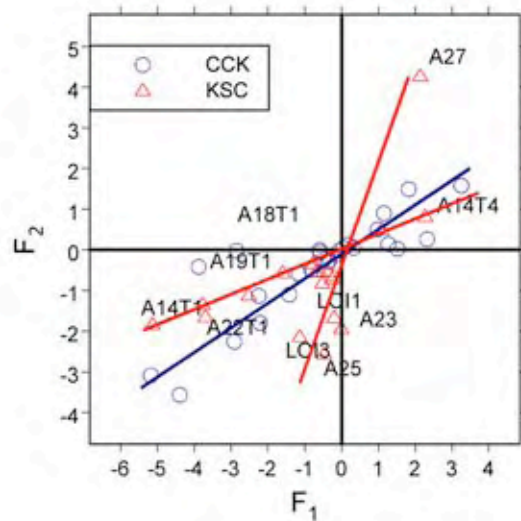
Examining the Failure of Item Equating – Two Dimensional Item Difficulties



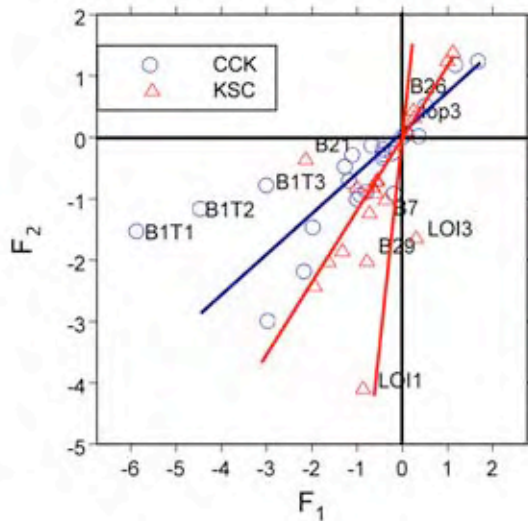
<http://www.soe.umich.edu/lmt/>



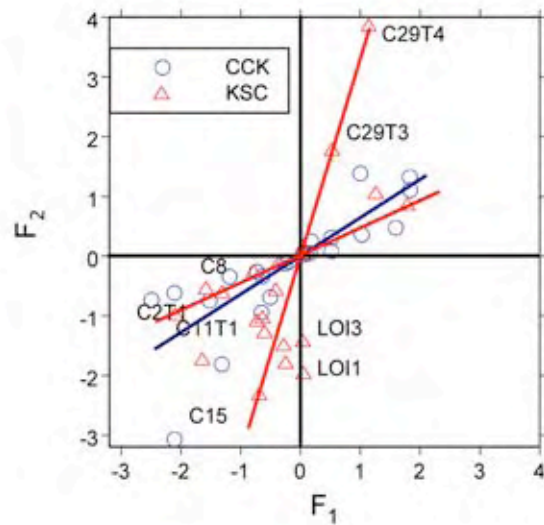
Two Dimensional Form A Item Difficulties



Two Dimensional Form B Item Difficulties



Two Dimensional Form C Item Difficulties



Reformulating the Interpretive Argument for KSC: Two Approaches

1. View KSC from a generalizability interpretation
 - Item universe is two dimensional and we are sampling items from a two dimensional item universe. Leads to a Rasch model.
2. View KSC as those items close to the F_2 axis as measuring the KSC construct

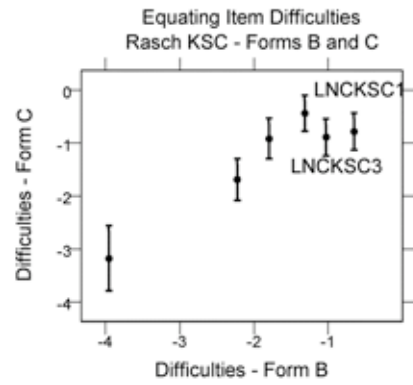
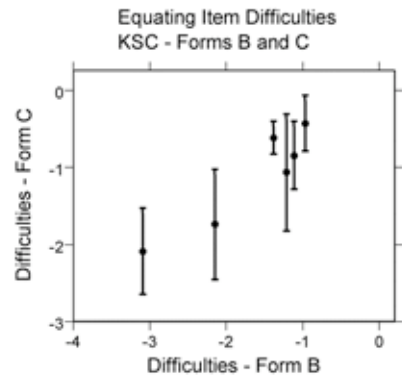
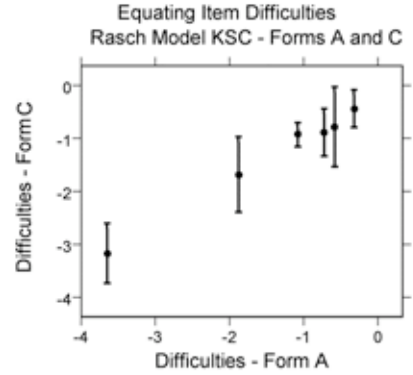
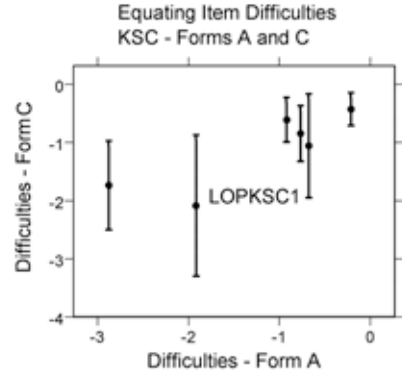
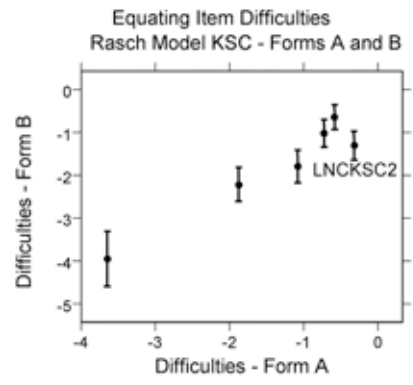
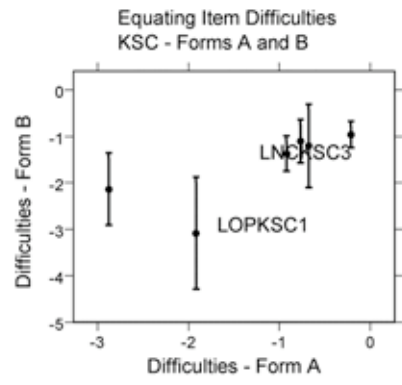
First Approach: Fitting a Rasch Model

- Reliability and information similar to general IRT Model.
- Greatly improves equating of forms

<http://www.soe.umich.edu/lmt/>

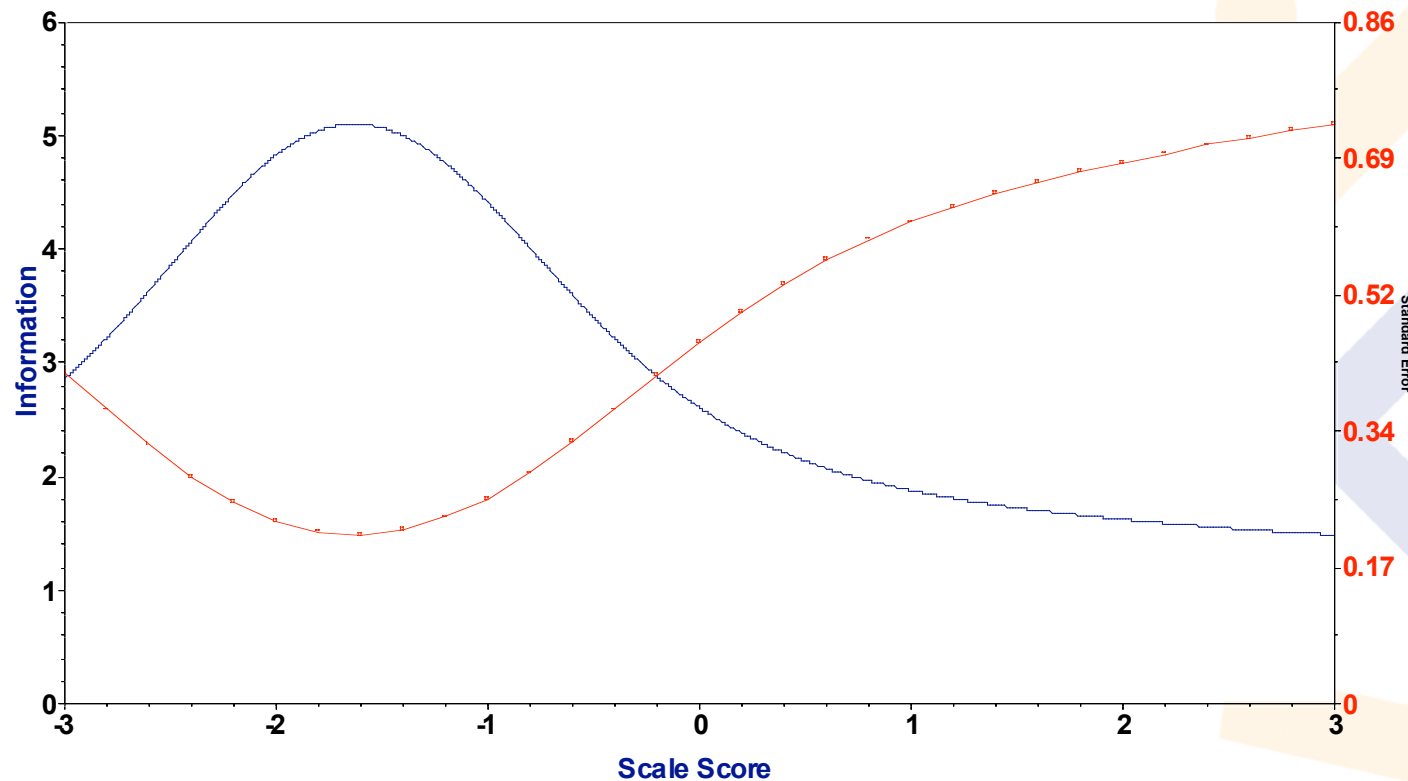


54



Second Approach: Fit a General Model to Items Close to the F_2 axis

Test Information and Measurement Error - KSCAI Selected Items



Reliability = 0.75

<http://www.soe.umich.edu/lmt/>



Conclusions

- Unidimensional and multidimensional IRT provide important tools for evaluating the structure of CKT-M
- Argument based approach provides structure for applying evidence to interpretive argument and reformulating that argument
- Implications of results for constructing measures in a multidimensional item universe

Discussants

- Michael Kane, National Bar Association
- Mark Reckase, Michigan State University
- Anna Sfard, University of Haifa and Michigan State University

<http://www.soe.umich.edu/lmt/>

