Evaluating Multi-Item Scales



Health Services Research Design (HS 225B)

January 26, 2015, 1:00-3:00pm 71-257 CHS Listed below are a few statements about your relationships with others. How much is each statement TRUE or FALSE for you

	Definitely True	Mostly True	Don't Know	Mostly False	Definitely False
1. I am always courteous even to people who are disagreeable.	1	2	3	4	5
2. There have been occasions when I took advantage of someone.	n 1	2	3	4	5
3. I sometimes try to get even rather than forgive and forget.	er 1	2	3	4	5
4. I sometimes feel resentful when don't get my way.	I 1	2	3	4	5
5. No matter who I' m talking to, I always a good listener.	'm 1	2	3	4	5

Scoring Multi-Item Scales

- Average or sum all items in the same scale.
- Transform average or sum to
 - 0 (worse) to 100 (best) possible range
 - z-score (mean = 0, SD = 1)
 - T-score (mean = 50, SD = 10)

Linear Transformations



Y = target mean + (target SD * Zx)

Listed below are a few statements about your relationships with others. How much is each statement TRUE or FALSE for you

	Definitely True	Mostly True	Don't Know	Mostly False	Definitely False
1. I am always courteous even to people who are disagreeable.	100	75	50	25	0
2. There have been occasions when I took advantage of someone.	n 0	25	50	75	100
3. I sometimes try to get even rathe than forgive and forget.	er 0	25	50	75	100
4. I sometimes feel resentful when don't get my way.	I 0	25	50	75	100
5. No matter who I' m talking to, I always a good listener.	' m 100	75	50	25	0

Create T-score

$$z$$
-score = (score – 36)/31
T-score = (10 * z-score) + 50

$$z$$
-score = (100- 36)/31 = 2.06
T-score = 71

- 1. Same people get same scores
- 2. Different people get different scores and differ in the way you expect
- 3. Measure is interpretable
- 4. Measure works the same way for different groups (age, gender, race/ethnicity)

Aside from being practical..

1. Same people get same scores — Relicible

2. Different people get different scores and differ in the way you expect

3. Measure is interpretable

4. Measure works the same way for different groups (age, gender, race/ethnicity)

- 1. Same people get same scores
- 2. Different people get different scores and differ in the way you expect V 1 ↓
- 3. Measure is interpretable
- 4. Measure works the same way for different groups (age, gender, race/ethnicity)

- 1. Same people get same scores
- 2. Different people get different scores and differ in the way you expect
- 3. Measure is interpretable \longrightarrow \square \square
- 4. Measure works the same way for different groups (age, gender, race/ethnicity)

- 1. Same people get same scores
- 2. Different people get different scores and differ in the way you expect
- 3. Measure is interpretable
- 4. Measure works the same way for different groups (age, gender, race/ethnicity) □□□.

Reliability

- Degree to which the same score is obtained when the *target* or thing being measured (person, plant or whatever) hasn't changed.
- ✓Inter-rater (rater)

✓Need 2 or more raters of the thing being measured

✓Internal consistency (items)

✓ Need 2 or more items

✓Test-retest (administrations)

✓ Need 2 or more time points

Ratings of Performance of Six HPM 225B Lectures by Two Raters

[1 = Poor; 2 = Fair; 3 = Good; 4 = Very good; 5 = Excellent]

- 1= Marjorie Kagawa-Singer (Good, Very Good)
- 2= Tom Belin (Very Good, Excellent)
- 3= Rick Dees (Good, Good)
- 4= Ron Hays (Fair, Poor)
- 5= Jack Needleman (Excellent, Very Good)
- 6= Jane Error (Fair, Fair)

(Target = 6 presenters; assessed by 2 raters) ¹³

Reliability Formulas

Model	Reliability	Intraclass Correlation
Two-way random	$\frac{N(MS_{BMS} - MS_{EMS})}{NMS_{BMS} + MS_{JMS} - MS_{EMS}}$	$\frac{MS_{BMS} - MS_{EMS}}{MS_{BMS} + (k-1)MS_{EMS} + k(MS_{JMS} - MS_{EMS})/N}$
Two- way mixed	$\frac{MS_{BMS} - MS_{EMS}}{MS_{BMS}}$	$\frac{MS_{BMS} - MS_{EMS}}{MS_{BMS} + (k-1)MS_{EMS}}$
One- way	$\frac{MS_{BMS} - MS_{WMS}}{MS_{BMS}}$	$\frac{MS_{BMS} - MS_{WMS}}{MS_{BMS} + (k-1)MS_{WMS}}$
BMS = WMS = JMS EMS =	 Between Ratee Mean Squ Within Mean Square Item or Rater Mean Squa Ratee x Item (Rater) Mea 	are N = n of ratees k = n of items or raters re n Square

$\begin{array}{c} 01 & 13 \\ 01 & 24 \\ 02 & 14 \\ 02 & 25 \\ 03 & 13 \\ 03 & 23 \\ 04 & 12 \\ 04 & 21 \\ 05 & 15 \\ 05 & 24 \\ 06 & 12 \\ 06 & 22 \end{array}$	Two-Wo (Reliability of	iy Ranc Ratings	dom Effe of Presenta	c ts tions)
Source		df	SS	MS
Present	ers (BMS)	5	15.67	3.13
Raters	(JMS)	1	0.00	0.00
Pres. x	Raters (EMS)	5	2.00	0.40
Tota	.1	11	17.67	
2-way	$y R = \frac{6 (3.13 - 0.4)}{6 (3.13) + 0}$	<u>40)</u> .00 - 0.40	= 0.89	ICC = 0.8(

Responses of Presenters to Two Questions about Their Health

- 1= Marjorie Kagawa-Singer (Good, Very Good)
- 2= Tom Belin (Very Good, Excellent)
- 3= Rick Dees (Good, Good)
- 4= Ron Hays (Fair, Poor)
- 5= Jack Needleman (Excellent, Very Good)
- 6= Jane Error (Fair, Fair)

(Target = 6 presenters; assessed by 2 items)

01 34 02 45 03 33 04 21 05 54 06 22	lixed Ef [.]	fects (Cron	bach's Alpha)
Source	df	SS	MS
Presenters (BMS) Items (JMS) Pres. x Items (EMS)	5 1 5	$15.67 \\ 0.00 \\ 2.00$	3.13 0.00 0.40
Total	11	17.67	
Alpha = $\frac{3.13 - 0}{3.13}$.40 = 2.93 3.13	= 0.87	ICC = 0.77

Reliability Minimum Standards

- 0.70 or above (for group comparisons)
- 0.90 or higher (for individual assessment)
 - SEM = SD (1- reliability)^{1/2}
 95% CI = true score +/- 1.96 x SEM
 - if z-score = 0, then CI: -.62 to +.62 when reliability = 0.90
 Width of CI is 1.24 z-score units



Calculating KAPPA

$$\mathbf{P_{C}} = \frac{(0 \times 1) + (2 \times 1) + (2 \times 1) + (1 \times 2) + (1 \times 1)}{(6 \times 6)} = 0.19$$

$$\mathbf{P}_{obs.} = \frac{2}{6} = 0.33$$

Kappa = $\frac{0.33 - 0.19}{1 - 0.19} = 0.17$

Weighted Kappa Linear (Quadratic)

	Ρ	F	G	VG	E
Р	1	.75 (.937)	.50 (.750)	.25 (.437)	0
F	.75 (.937)	1	.75 (.937)	.50 (.750)	.25 (.437)
G	.50 (.750)	.75 (.937)	1	.75 (.937)	.50 (.750)
VG	.25 (.437)	.50 (.750)	.75 (.937)	1	.75 (.937)
E	0	.25 (.437)	.5 (.750)	.75 (.937)	1
W _I = 1 – (i/	(k – 1))				

 $W_q = 1 - (i^2 / (k - 1)^2)$

i = number of categories ratings differ byk = n of categories

All Kappas

$$\mathbf{P_{c}} = \frac{(0 \times 1) + (2 \times 1) + (2 \times 1) + (1 \times 2) + (1 \times 1)}{(6 \times 6)} = 0.19$$

$$P_{obs.} = \frac{2}{6} = 0.33$$

Kappa =
$$\frac{0.33 - 0.19}{1 - 0.19} = 0.17$$

Linear weighted kappa = 0.52Quadratic weighted kappa = 0.77

Guidelines for Interpreting Kappa

Conclusion	Kappa	Conclusion	Kappa
Poor	< .40	Poor	< 0.0
Fair	.4059	Slight	.0020
Good	.6074	Fair	.2140
Excellent	> .74	Moderate	.4160
		Substantial	.6180
		Almost perfect	.81 - 1.00
Fleiss (1981)		Landis and Koch (1	977)

Item-scale correlation matrix

	<u>Depress</u>	<u>Anxiety</u>	<u>Anger</u>
ltem #1	0.80*	0.20	0.20
Item #2	0.80*	0.20	0.20
Item #3	0.80*	0.20	0.20
Item #4	0.20	0.80*	0.20
Item #5	0.20	0.80*	0.20
Item #6	0.20	0.80*	0.20
Item #7	0.20	0.20	0.80*
Item #8	0.20	0.20	0.80*
Item #9	0.20	0.20	0.80*



*Item-scale correlation, corrected for overlap.

Item-scale correlation matrix

	<u>Depress</u>	<u>Anxiety</u>	<u>Anger</u>
ltem #1	0.50*	0.50	0.50
Item #2	0.50*	0.50	0.50
Item #3	0.50*	0.50	0.50
Item #4	0.50	0.50*	0.50
Item #5	0.50	0.50*	0.50
Item #6	0.50	0.50*	0.50
Item #7	0.50	0.50	0.50*
Item #8	0.50	0.50	0.50*
Item #9	0.50	0.50	0.50*



*Item-scale correlation, corrected for overlap.

Confirmatory Factor Analysis

	<u>Depress</u>	<u>Anxiety</u>	<u>Anger</u>
ltem #1	0.80*	0.00	0.00
Item #2	0.80*	0.00	0.00
Item #3	0.80*	0.00	0.00
Item #4	0.00	0.80*	0.00
Item #5	0.00	0.80*	0.00
Item #6	0.00	0.80*	0.00
Item #7	0.00	0.00	0.80*
Item #8	0.00	0.00	0.80*
Item #9	0.00	0.00	0.80*

*Factor loading.



Validity

Does scale represent what it is supposed to be measuring?

- Content validity: Does measure "appear" to reflect what it is intended to (expert judges or patient judgments)?
 - Do items operationalize concept?
 - Do items cover all aspects of concept?
 - Does scale name represent item content?
- Construct validity
 - Are the associations of the measure with other variables consistent with hypotheses?

Relative Validity Example

Sensitivity of measure to important (clinical) difference

	Severity	of Kidney	UNE WAY		
	None	Mild	Severe	F-ratio	Relative Validity
Burden of Disease #1	87	90	91	2	
Burden of Disease #2	74	78	88	10	5
Burden of Disease #3	77	87	95	20	10

Scale	Age (years)		
(Better) Physical Functioning	(-)		

Scale	Age (years)		
(Better) Physical Functioning	Medium (-)		

Scale	Age (years)		
(Better) Physical Functioning	Medium (-)		

Effect size (ES) = D/SD

D = Score difference SD = SD

Small (0.20), medium (0.50), large (0.80)

Scale	Age (years)		
(Better) Physical Functioning	Medium (-) r ~0.24		

Cohen effect size rules of thumb (d = 0.20, 0.50, and 0.80):

Small r = 0.100; medium r = 0.243; large r = 0.371

```
\underline{\mathbf{r}} = \underline{\mathbf{d}} / [(\underline{\mathbf{d}}^2 + 4)^{.5}] = \underline{\mathbf{0.80}} / [(0.80^2 + 4)^{.5}] = 0.80 / [(0.64 + 4)^{.5}] = 0.80 / [(4.64)^{.5}] = 0.80 / 2.154 = \underline{0.371}
```

Scale	Age (years)	Obese yes = 1, no = 0	Kidney Disease yes = 1, no = 0	In Nursing home yes = 1, no = 0
(Better) Physical Functioning	Medium (-)	Small (-)	Large (-)	Large (-)

Cohen effect size rules of thumb (d = 0.20, 0.50, and 0.80):

Small r = 0.100; medium r = 0.243; large r = 0.371

$$\underline{\mathbf{r}} = \underline{\mathbf{d}} / [(\underline{\mathbf{d}}^2 + 4)^{.5}] = \underline{\mathbf{0.80}} / [(0.80^2 + 4)^{.5}] = 0.80 / [(0.64 + 4)^{.5}] = 0.80 / [(4.64)^{.5}] = 0.80 / 2.154 = \underline{0.371}$$

Scale	Age (years)	Obese yes = 1, no = 0	Kidney Disease yes = 1, no = 0	In Nursing home yes = 1, no = 0
(Better) Physical Functioning	Medium (-)	Small (-)	Large (-)	Large (-)
(More) Depressive Symptoms	(?)	Small (+)	(?)	Small (+)

Cohen effect size rules of thumb (d = 0.20, 0.50, and 0.80):

Small r = 0.100; medium r = 0.243; large r = 0.371

```
\underline{\mathbf{r}} = \underline{\mathbf{d}} / [(\underline{\mathbf{d}}^2 + 4)^{.5}] = \underline{\mathbf{0.80}} / [(0.80^2 + 4)^{.5}] = 0.80 / [(0.64 + 4)^{.5}] = 0.80 / [(4.64)^{.5}] = 0.80 / 2.154 = \underline{0.371}
```

(r's of 0.10, 0.30 and 0.50 are often cited as small, medium, and large.) $_{35}$

Questions?



Responsiveness to Change

- Measures should be responsive to interventions that change the underlying construct
- Need external indicators of change (Anchors)
 Clinical measure
 - "improved" group = 100% reduction in seizure frequency
 - "unchanged" group = <50% change in seizure frequency
 - Retrospective self- or provider-report of change
 - Much better, A little better, Same, A little worse, Much worse
- Anchor correlated with change on target measure at 0.371 or higher

Responsiveness Index

- Effect size (ES) = D/SD
 - D = raw score change in "changed" (improved) groupSD = baseline SD
- Small: 0.20->0.49
- Medium: 0.50->0.79
- Large: 0.80 or above

Responsiveness Indices

(1) Effect size (ES) = D/SD

(2) Standardized Response Mean (SRM) = D/SD†(3) Guyatt responsiveness statistic (RS) = D/SD‡

D = raw score change in "changed" group;
SD = baseline SD;
SD† = SD of D;
SD‡ = SD of D among "unchanged"

Amount of Expected Change Varies

SF-36 physical function score mean = 87 (SD = 20) \checkmark Assume I have a score of 100 at baseline

Hit by Bike causes me to be

- limited a lot in vigorous activities
- limited a lot in climbing several flights of stairs
- limited a little in moderate activities

SF-36 physical functioning score drops to 75 (-1.25 SD)

Hit by Rock causes me to be

- limited a little in vigorous activities

SF-36 physical functioning score drops to 95 (- 0.25 SD)

Partition Change on Anchor

A lot better
A little better
No change
A little worse

>A lot worse

Use Multiple Anchors

- 693 RA clinical trial participants evaluated at baseline and 6weeks post-treatment.
- Five anchors:
 - 1. Self-report (global) by patient
 - 2. Self-report (global) by physician
 - 3. Self-report of pain
 - 4. Joint swelling (clinical)
 - 5. Joint tenderness (clinical)

Kosinski, M. et al. (2000). Determining minimally important changes in generic and diseasespecific health-related quality of life questionnaires in clinical trials of rheumatoid arthritis. Arthritis and Rheumatism, 43, 1478-1487.

Patient and Physician Global Reports

How are you (is the patient) doing, considering all the ways that RA affects you (him/her)?

- Very good (asymptomatic and no limitation of normal activities)
- Good (mild symptoms and no limitation of normal activities)
- Fair (moderate symptoms and limitation of normal activities)
- Poor (severe symptoms and inability to carry out most normal activities)
- Very poor (very severe symptoms that are intolerable and inability to carry out normal activities

--> Improvement of 1 level **over time**

Global Pain, Joint Swelling and Tenderness

- 0 = no pain, 10 = severe pain
- Number of swollen and tender joints

-> 1-20% improvement over time

Effect Sizes for SF-36 Physical Function Change Linked to Minimal Change in Anchors

Scale	Self-R	ClinR	Pain	Swell	Tender	Mean
Physical Function	<u>.35</u>	.33	.34	<u>.26</u>	.32	.32
						45

Effect Sizes for SF-36 Changes Linked to Minimal Change in Anchors

Scale	Self-R			
PF	.35			
Role-P	.56			
Pain	<u>.83</u>			
GH	<u>.20</u>			
EWB	.39			
Role-E	.41			
SF	.43			
EF	.50			
PCS	.49			
MCS	.42			

Effect Sizes (mean = 0.34) for SF-36 Changes Linked to Minimal Change in Anchors

Scale	Self-R	ClinR	Pain	Swell	Tender	Mean
PF	<u>.35</u>	.33	.34	<u>.26</u>	.32	.32
Role-P	<u>.56</u>	.52	<u>.29</u>	.35	.36	.42
Pain	<u>.83</u>	.70	.47	.69	<u>.42</u>	.62
GH	<u>.20</u>	.12	.09	.12	<u>.04</u>	.12
EWB	<u>.39</u>	.26	.25	.18	<u>.05</u>	.23
Role-E	<u>.41</u>	.28	<u>.18</u>	.38	.26	.30
SF	<u>.43</u>	.34	<u>.28</u>	.29	.38	.34
EF	<u>.50</u>	.47	<u>.22</u>	.22	.35	.35
PCS	<u>.49</u>	.48	<u>.34</u>	.29	.36	.39
MCS	.42	.27	<u>.19</u>	.27	.20	.27

Item Response Theory (IRT)

IRT models the relationship between a person's response Y_i to the question (i) and his or her level of the latent construct θ being measured by positing

$$\Pr(Y_i \ge k) = \frac{1}{1 + \exp(-a_i\theta + b_{ik})}$$

b_{ik} estimates how difficult it is to get a score of k or more on item (i).

 a_i is an estimate of the discriminatory power of the item.

Item Responses and Trait Levels



www.nihpromis.org

Computer Adaptive Testing (CAT)







Reliability Target for Use of Measures with Individuals

- Reliability ranges from 0-1
 - 0.90 or above is goal
 SEM = SD (1- reliability)^{1/2}
 - ➢ 95% CI = true score +/- 1.96 x SEM
 - > if true z-score = 0, then CI: -.62 to +.62
 > Width of CI is 1.24 z-score units
- Reliability = 0.90 when <u>SE = 3.2</u>
 - T-scores (mean = 50, SD = 10) T = 50 + (z * 10)
 - Reliability = $1 (SE/10)^2$

I was grouchy [1st question]

- Never	[39]
- Rarely	[48]
- Sometimes	[56]
- Often	[64]
- Always	[72]

Estimated Anger = 56.1 SE = 5.7 (rel. = 0.68)

I felt like I was ready to explode

[2nd question]

- Never
- Rarely
- Sometimes
- Often
- Always

Estimated Anger = 51.9 SE = 4.8 (rel. = 0.77)

- I felt angry [3rd question]
 - Never
 - Rarely
 - Sometimes
 - Often
 - Always

Estimated Anger = 50.5 SE = 3.9 (rel. = 0.85)

I felt angrier than I thought I should [4th question]

- Never
- Rarely
- Sometimes
- Often
- Always

Estimated Anger = 48.8 SE = 3.6 (rel. = 0.87)

- I felt annoyed [5th question]
 - Never
 - Rarely
 - Sometimes
 - Often
 - Always

Estimated Anger = 50.1 SE = 3.2 (rel. = 0.90)

I made myself angry about something just by thinking about it. [6th question]

- Never
- Rarely
- Sometimes
- Often
- Always

Estimated Anger = 50.2SE = 2.8 (rel = 0.92)

PROMIS Physical Functioning vs. "Legacy" Measures



Thank you. Powerpoint file is freely available at: <u>http://gim.med.ucla.edu/FacultyPages/Hays/</u>

Contact information: <u>drhays@ucla.edu</u> 310-794-2294

