Evaluating Self-Report Data Using Psychometric Methods

Ron D. Hays, Ph.D. (hays@rand.org) February 5, 2003 (3:00-6:00pm)

Four Types of Data Collection Errors

 Coverage Error
 Does each person in population have an equal chance of selection?

 Sampling Error Are only some members of the population sampled?

 Nonresponse Error Do people in the sample who respond differ from those who do not?

• Measurement Error

Are inaccurate answers given to survey questions?

What's a Good Measure?

- Same person gets same score (reliability)
- Different people get different scores (validity)
- People get scores you expect (validity)
- It is practical to use (feasibility)



How Are Good Measures Developed?

- Review literature
- Expert input (patients and clinicians)
- Define constructs you are interested in
- Draft items (item generation)
- Pretest
 - Cognitive interviews
 - Field and pilot testing
- Revise and test again
- Translate/harmonize across languages

Scales of Measurement and Their Properties

Property of Numbers

Type of Scale	Rank Order	Equal Interval	Absolute 0
Nominal	No	No	No
Ordinal	Yes	No	No
Interval	Yes	Yes	No
Ratio	Yes	Yes	Yes

Measurement Range for Health Outcome Measures



Indicators of Acceptability

- Response rate
- Administration time
- Missing data (item, scale)

Variability

- All scale levels are represented
- Distribution approximates bell-shaped "normal"



Measurement Error

observed = true + systematic + random score error error

(bias)

Flavors of Reliability

Test-retest (administrations)
Intra-rater (raters)
Internal consistency (items)

Test-retest Reliability of MMPI 317-362 r = 0.75

MMPI 317

	True	False	
True MMPI 362	169	15	184
False	21	95	116
	190	110	

I am more sensitive than most other people.

Kappa Coefficient of Agreement (Corrects for Chance)

> kappa = (observed - chance) (1 - chance)

Example of Computing KAPPA



Example of Computing KAPPA (Continued)



Guidelines for Interpreting Kappa

<u>Conclusion</u> Poor	<u>Kappa</u> <.40	<u>Conclusion</u> Poor	<u>Kappa</u> < 0.0
Fair	.4059	Slight	.0020
Good	.6074	Fair	.2140
Excellent	> .74	Moderate	.4160
		Substantial	.6180
		Almost perfect	.81 - 1.00

Landis and Koch (1977)



Ratings of Height of Houseplants

Plan	t	Baseline Height	Follow-up Height	Experimental Condition
A1				
	R1	120	121	1
	R2	118	120	
A2				
	R1	084	085	2
	R2	096	088	
B1				
	R1	107	108	2
	R2	105	104	
B2				
	R1	094	100	1
	R2	097	104	
C1				
	R1	085	088	2
	R2	091	096	

Ratings of Height of Houseplants (Cont.)

Plan	t	Baseline Height	Follow-up Height	Experimental Condition
C2				
	R1	079	086	1
	R2	078	092	
D1				
	R1	070	076	1
	R2	072	080	
D2				
	R1	054	056	2
	R2	056	060	
E1				
	R1	085	101	1
	R2	097	108	
E2				
	R1	090	084	2
	R2	092	096	

Reliability of Baseline Houseplant Ratings

Ratings of Height of Plants: 10 plants, 2 raters Baseline Results

Source	DF	SS	MS	F
Plants	9	5658	628.667	35.52
Within	10	177	17.700	
Raters	1	57.8	57.800	
Raters x Plants	9	119.2	13.244	
otal	19	5835		

Sources of Variance in Baseline Houseplant Height

Source	dfs	MS	5
Plants (N)	9	628.67	(BMS)
Within	10	17.70	(WMS)
Raters (K)	1	57.80	(JMS)
Raters × Plants	9	13.24	(EMS)

Total

19

Intraclass Correlation and Reliability

Model	Reliability	Intraclass Correlation
One-Way	MS _{BMS} - MS _{WMS} MS _{BMS}	MS _{BMS} - MS _{WMS} MS _{BMS} + (K-1)MS _{WMS}
Two-Way Fixed	MS _{BMS} - MS _{EMS}	$\frac{\text{MS}_{\text{BMS}} - \text{MS}_{\text{EMS}}}{\text{MS}_{\text{EMS}} + (\text{K-1})\text{MS}_{\text{EMS}}}$
Two-Way Random	N (MS BMS - MS EMS) NMS BMS + MS JMS - MS EMS	$\frac{MS_{BMS} - MS_{EMS}}{MS_{BMS} + (K-1)MS_{EMS}} + \frac{K(MS_{JMS} - MS_{EMS})}{K}$

Summary of Kei	iabil	ity of	f Plan	t Ratings
B	aseline	Fo	ollow-up	
	≷ _{TT}	R _{II}	R _{TT}	R _{II}
One-Way Anova C	.97 (0.95	0.97 (0.94
Two-Way Random Effects ().97 (0.95	0.97	0.94
Two-Way Fixed Effects _C).98 (0.96	0.98	0.97
Source La	bel	Baseli	<u>ne MS</u>	
Plants BA	٨S	628.6	67	
Within WMS		17.700		
Raters J/	٨S	57.80	0	
Raters X Plants EA	۸S	13.24	4	
ICC (1,1) = $\frac{BMS - WMS}{BMS + (K - 1) * V}$	WMS			
$ICC (2,1) = \frac{BMS - EMS}{BMS + (K - 1) * K}$	EMS + k	(JMS -	EMS)/n	
ICC $(3,1) = \frac{BMS - EMS}{BMS + (K - 1) * 1}$	-MS			21.1/22/40

Cronbach's Alpha

Source	df	SS	MS
Respondents (Items (JMS) Resp. x Items	BMS) 4 1 (EMS) 4	11.6 0.1 4.4	2.9 0.1 1.1
Total	9	16.1	
Alpha = <u>2.9</u> 2.	<u>- 1.1</u> = <u>1.8</u> = <u>0.6</u> 9 2.9	52	

Alpha by Number of Items and Inter-item Correlations

$$alpha_{st} = \frac{K \overline{r}}{1 + (K - 1) \overline{r}}$$

K = number of items in scale

Alpha for Different Numbers of Items and Homogeneity

Average Inter-item Correlation (\overline{r})

Number of Items (K) .0	.2	.4	.6	.8	1.0
2	.000	.333	.572	.750	.889	1.000
4	.000	.500	.727	.857	.941	1.000
6	.000	.600	.800	.900	.960	1.000
8	.000	.666	.842	.924	.970	1.000

Number of Items and Reliability for Three Versions of the Mental Health Inventory (MHI)

Measure	Number of Items	Completion time (min.)	Reliability
MHI-32	32	5-8	.98
MHI-18	18	3-5	.96
MHI-5	5	1 or less	.90

Data from McHorney et al. 1992

Spearman-Brown Prophecy Formula

alpha y =
$$\begin{pmatrix} N \cdot alpha_{x} \\ 1 + (N - 1) * alpha_{x} \end{pmatrix}$$

N = how much longer scale y is than scale x

Reliability Minimum Standards

- 0.70 or above (for group comparisons)
- 0.90 or higher (for individual assessment)

> SEM = SD (1- reliability)^{1/2}

Reliability of a Composite Score

$$Mosier = 1 - \frac{\Sigma(w_j^2)(S_j^2) - \Sigma(w_j^2)(S_j^2)(\alpha_j)}{\Sigma(w_j^2)(S_j^2) + 2\Sigma(w_j)(w_{\kappa})(S_j)(S_{\kappa})(r_{j\kappa})}$$

- w_j = weight given to component J
- w_{κ} = weight given to component K
- **S**_j = standard deviation of J
- α_j = reliability of J
- \mathbf{r}_{jK} = correlation between J and K

Hypothetical Multitrait/Multi-Item Correlation Matrix

	<u>Trait #1</u>	<u>Trait #2</u>	<u>Trait #3</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.

Multitrait/Multi-Item Correlation Matrix for Patient Satisfaction Ratings

	Technical	Interpersonal	Communication	Financial
Technical				
1	0.66*	0.63†	0.67†	0.28
2	0.55*	0.54†	0.50†	0.25
3	0.48*	0.41	0.44†	0.26
4	0.59*	0.53	0.56†	0.26
5	0.55*	0.60†	0.56†	0.16
6	0.59*	0.58†	0.57†	0.23
Interpersonal				
1	0.58	0.68*	0.63†	0.24
2	0.59†	0.58*	0.61†	0.18
3	0.62†	0.65*	0.67†	0.19
4	0.53†	0.57*	0.60†	0.32
5	0.54	0.62*	0.58†	0.18
6	0.48†	0.48*	0.46†	0.24

Note - Standard error of correlation is 0.03. Technical = satisfaction with technical quality. Interpersonal = satisfaction with the interpersonal aspects. Communication = satisfaction with communication. Financial = satisfaction with financial arrangements. *Item-scale correlations for hypothesized scales (corrected for item overlap). †Correlation within two standard errors of the correlation of the item with its hypothesized scale.



http://www.tricare.osd.mil/tricaresurveys/rel_val.html

	Item Internal Consistency Validity Correlation			
Questions	Interpersonal Relationship	Quality of Medical Care	Access to Medical Care	
Q3A	.82*	0.79	0.61	
Q3B	.85*	0.82	0.6	
Q3E	.89*	0.88	0.6	
Q3F	.84*	0.84	0.59	
Q3G	.84*	0.84	0.63	
Q3C	0.89	.87*	0.6	
Q3D	0.88	.85*	0.59	
Q3H	0.81	.88*	0.6	
Q3I	0.81	.90*	0.6	
Q3J	0.89	.91*	0.63	
Q07	0.51	0.5	.69*	
Q09	0.59	0.56	.64*	
Q10A	0.53	0.5	.79*	
Q10B	0.58	0.57	.82*	
Q10C	0.58	0.57	.80*	
Q11	0.59	0.57	.81*	
* Denotes questions include the Scale				

Lo 1 Col 1



Page 1 Sec 1 1/1 0F 1"

DEC TRE EVT OVD





What are IRT Models?

Mathematical equations that relate observed survey responses to a persons location on an unobservable latent trait (i.e., intelligence, patient satisfaction).

Latent Trait and Item Responses



34 1/23/18

IRT Model Assumptions

Unidimensionality

-One construct measured by items in scale.

Local Independence

-Items uncorrelated when latent trait(s) have been controlled for.

Types of IRT Models

- Unidimensional and multidimensional
- Dichotomous and polytomous
- Parameterization
 - One parameter: difficulty (location)
 - Two Parameter: difficulty and slope (discrimination)

Item difficulty

Transform proportion of people endorsing the item (p) to correspond to (1-p)th percentile from z distribution

- $Z = \ln (1-p)/p)/1.7 = (\ln (1-p) \ln (p))/1.7$
 - = (ln (.228) ln (.772))/1.7
 - = (-1.47840965 + .258770729)/1.7
 - = -1.21963892/1.7

Item Discrimination

Item-scale correlation, corrected for item overlap

$$-Z' = \frac{1}{2} [\ln (1 + r) - \ln (1 - r)]$$

- if
$$r = 0.80$$
, $z = 1.10$

- if
$$r = 0.95$$
, $z = 1.83$

(0.5 -> 2 is typical range)

1-Parameter Logistic Model for (Dichotomous Outcomes)

$$P_i(\Theta) = \frac{e^{(\Theta - b_i)}}{1 + e^{(\Theta - b_i)}}$$

 $P_i(\Theta)$ Probability that a randomly selected respondent with ability Θ (trait level) answers item i correctly.

b_i Item i difficulty.

Item Characteristic Curves (1-Parameter Model)



----- Item 1 (Difficulty = -1) ------ Item 3 (Difficulty = 1)

2-Parameter Logistic Model (Dichotomous Outcomes)

$$P_i(\Theta) = \frac{e^{Da_i(\Theta - b_i)}}{1 + e^{Da_i(\Theta - b_i)}}$$

 $P_i(\Theta)$ Probability that a randomly selected respondent with ability Θ (trait level) answers item i correctly. b_i Item i difficulty.

- a_i Item i slope.
- D Scaling constant.

Item Characteristic Curves (2-Parameter Model)



Item Responses and Trait Levels



Information Conditional on Trait Level

 Item information proportional to inverse of standard error:

$$SE(\Theta) = \frac{1}{\sqrt{I(\Theta)}}$$

 Scale/Test information is the sum over item information:

$$I(\Theta) = \sum_{i=1}^{n} I_i(\Theta)$$

Item Information (2-parameter model)



Linking Item Content to Trait Estimates



Dichotomous Items Showing DIF (2-Parameter Model)



Forms of Validity

- Content
- Criterion
- Construct Validity

Construct Validity

- Does measure relate to other measures in ways consistent with hypotheses?
- Responsiveness to change

Relative Validity Analyses

- Form of "known groups" validity
- Relative sensitivity of measure to important clinical difference
- · One-way between group ANOVA

Relative Validity Example

Severity of Heart Disease

	None	Mild	Severe	F-ratio	Relative Validity
Scale #1	87	90	91	2	
Scale #2	74	78	88	10	5
Scale #3	77	87	95	20	10

Responsiveness to Change and Minimally Important Difference

 HRQOL measures should be responsive to interventions that changes HRQOL

- Evaluating responsiveness requires assessment of HRQOL
 - pre-post intervention of known efficacy
 - at two times in tandem with gold standard

Two Essential Elements

- External indicator of change (Anchors)

 mean change in HRQOL scores among people who have a "minimal" change in HRQOL.
- Amount of HRQOL change

External Indicator of Change (A)

Overall has there been any change in your asthma since the beginning of the study?

Much improved; Moderately improved; Minimally improved

No change

Much worse; Moderately worse; Minimally worse

External Indicator of Change (B)

Rate your overall condition. This rating should encompass factors such as social activities, performance at work or school, seizures, alertness, and functional capacity; that is, your overall quality of life.

7 response categories; ranging from <u>no impairment</u> to <u>extremely severe impairment</u>

External Indicator of Change (C)

 "changed" group = seizure free (100% reduction in seizure frequency)

 "unchanged" group = < 50% change in seizure frequency

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"

Effect Size Benchmarks

- Small: 0.20->0.49
- Moderate: 0.50->0.79
- Large: 0.80 or above



Treatment Impact on PCS



Treatment Impact on MCS



