Methodological and Statistical Issues in Research Proposals

# Measures of all things
*Reliability and validity of research measures and how these impact your research proposals*

Rich Jones

NIDUS/CEDARTREE, 7th Annual Delirium Boot Camp

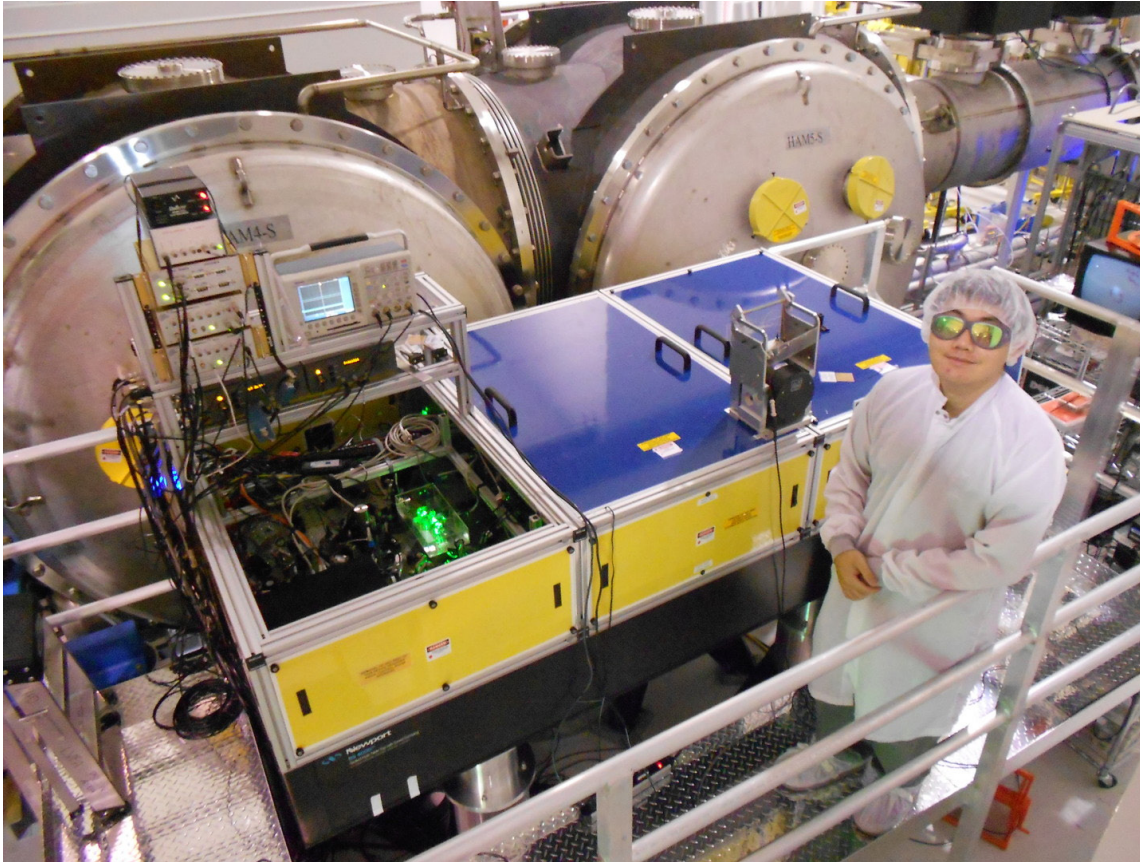October 28, 2019,  Penn Stater, State College PA

rich_jones@brown.edu   @rnjma

**Jakob Köbel** (1460 - 1533) - Geometrei. Von künstlichem Feldmessen und absehen (published first in **1535 or 1536,** reprinted in 1608.[1])

Stand at the door of a church on a Sunday and bid 16 men to stop, tall ones and small ones, as they happen to pass out when the service is finished; then make them put their left feet one behind the other, and the length thus obtained shall be a right and lawful rood to measure and survey the land with, and the 16th part of it shall be the right and lawful foot.

3

**1983**

The meter is defined as the length of the path travelled by light in a vacuum in 1/299,792,458 of a second

https://en.wikipedia.org/wiki/Metre

$c$ = speed of light in a vacuum
$c$ = 299,792,458 m/s

https://www.ligo.org/science/Publication-SqueezedVacuum/index.php

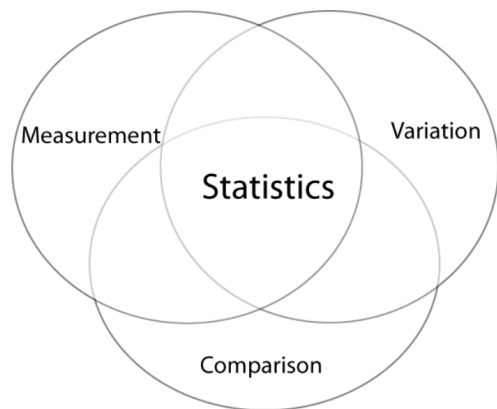Measurements are developed to address a specific practical need.

Measurements are refined as the need arises (e.g., for greater precision arises), often in the context of some new use.

Refinement of measures is linked with technological development.

# The #1 neglected topic in statistics is measurement

- Andrew Gelman

https://andrewgelman.com/2015/04/28/whats-important-thing-statistics-thats-not-textbooks/
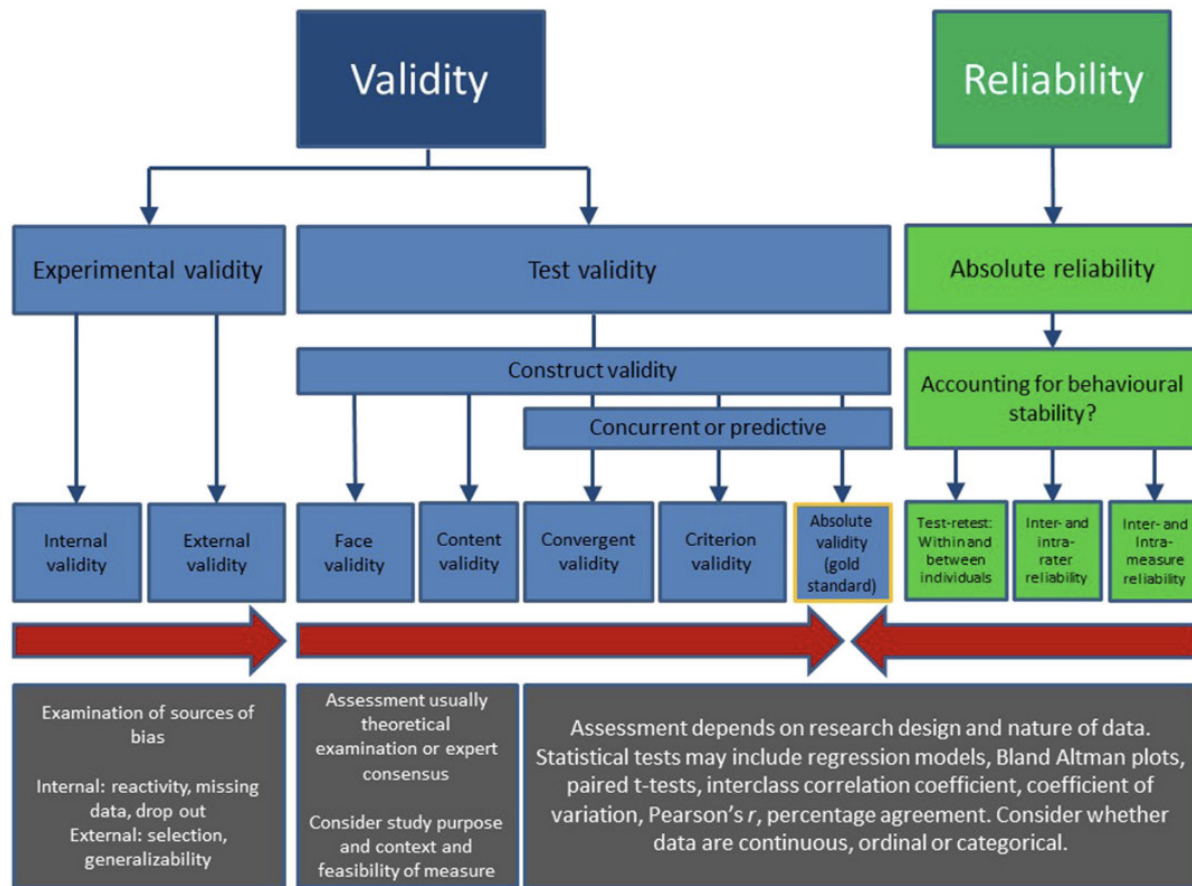
**DEBATE**

**Open Access**

# Should we reframe how we think about physical activity and sedentary behaviour measurement? Validity and reliability reconsidered

Paul Kelly[*] , Claire Fitzsimons and Graham Baker

***Terminology is used randomly, synonymously, possibly incorrectly and we all get confused***

Already we have used terms that you may have taken issue with. In many places we could have used different terms such as precision, concordance, uncertainty, or accuracy. There are also many sub-types of validity and reliability, some of which we have not yet discussed. For example, construct, comparative, absolute, relative, predictive, discriminant, representation, and translation validity; and inter-rater, intra-rater, relative, or absolute reliability.

**Fig. 3** The Edinburgh Framework v1.0 for validity and reliability in PA and SB measurement
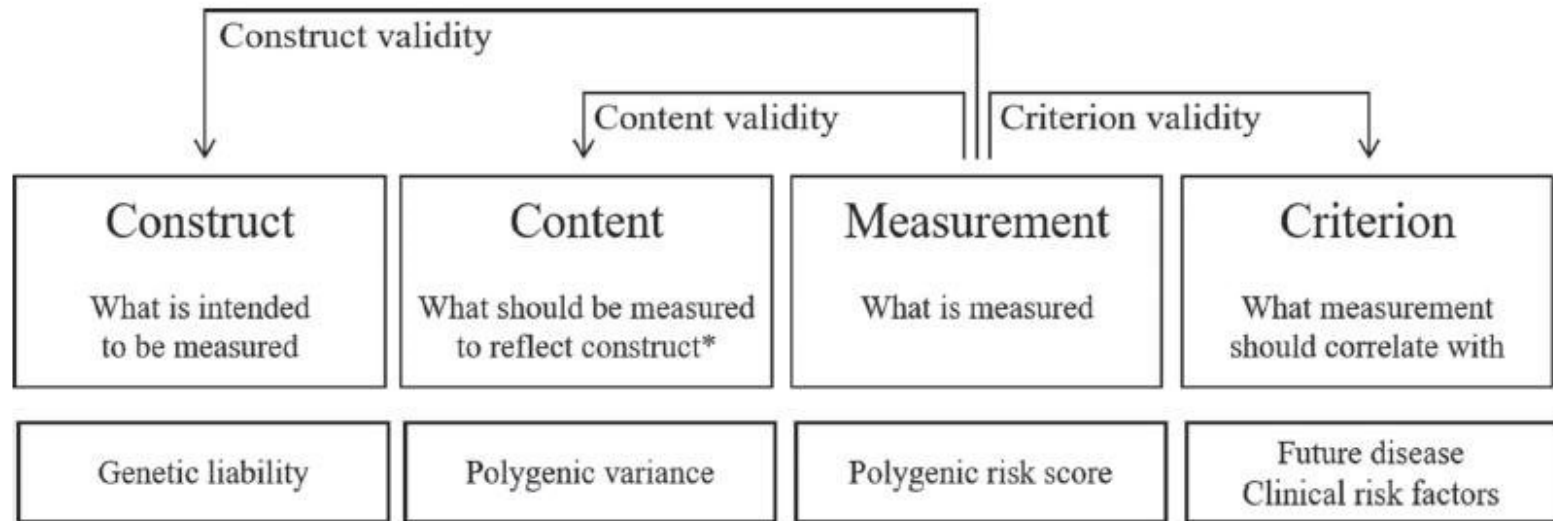
**Figure 1.** Three types of validity applied to the measurement of polygenic risk scores. Legend: * In the context of the specific application of the measurement.

Janssens A. 2019 Human Molecular Genetics, https://doi.org/10.1093/hmg/ddz205

12

McArdle, J., & Prescott, C. (1992). Age-based construct validation using structural equation modeling. *Experimental Aging Research, 18*(3), 87-116.

QUANTITATIVE TOPICS
IN RESEARCH ON AGING
J.J. McArdle and S.A. Cohen, Eds.

# Age-Based Construct Validation Using Structural Equation Modeling

J.J. MCARDLE
*The University of Virginia*

CAROL A. PRESCOTT
*Medical College of Virginia*

In this paper we describe some mathematical and statistical models based on *structural equation modeling* (SEM) using computer programs like LISREL. We focus on SEM methodology for the *simultaneous* examination of the internal validity of psychological constructs and the external validity represented by age relations. To illustrate these ideas we use a latent variable path model to examine the organization of intellectual abilities measured by the WAIS-R in the standardization sample. We also examine different ways in which age can be used to structure this organization. This is primarily a methodological paper, but we try to integrate conceptual principles of modeling with some substantive issues of research on the psychology of aging.

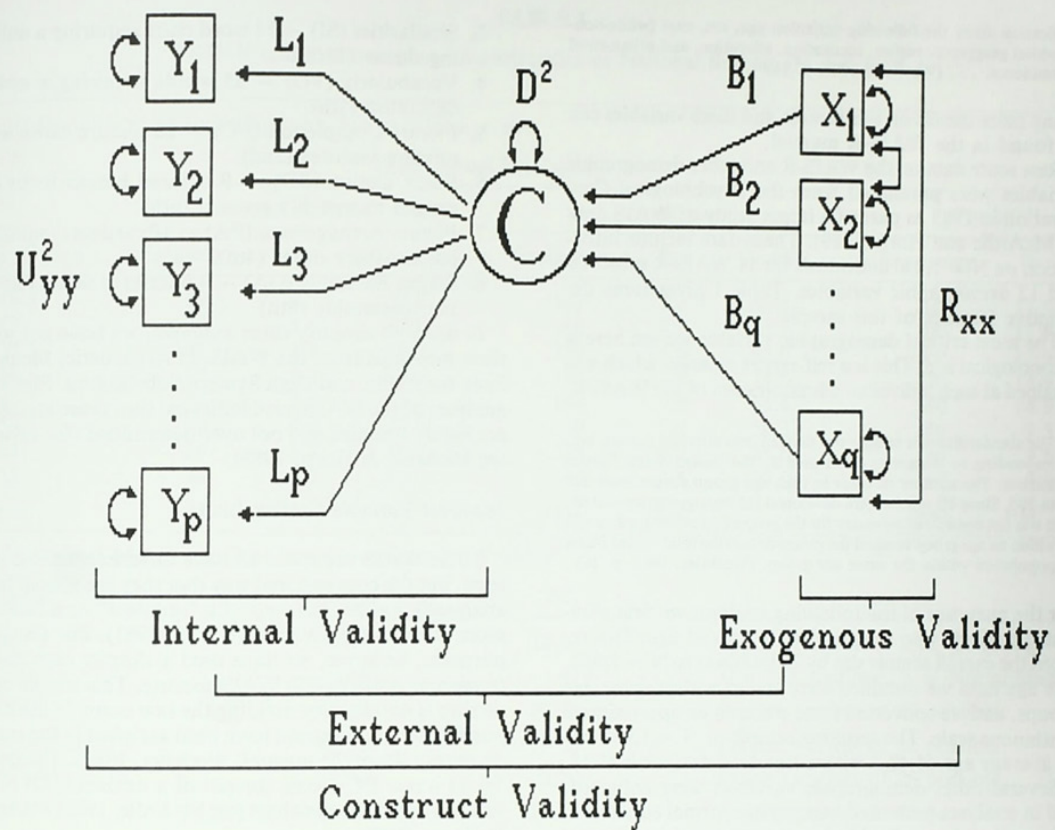**FIGURE 1.** A latent variable path diagram of a nomological network.

14

| **Reliability** | <u>How well</u> do we measure the thing? |
|---|---|
| **Validity** | <u>How well</u> do our measurements the thing we measure ***map*** on to the construct we want to measure? |

**Rich Jones** @rnjma · Sep 9

"A map is not the territory it represents, but, if correct, it has a similar structure to the territory, which accounts for its usefulness."
— Alfred Korzybski, Science and Sanity
en.wikipedia.org/wiki/Map%E2%80... @RsrveResilience

💬          🔁 1          ♡ 4          ⬆          ⬇ ᵢₗᵢ

**Outline**

Reliability ...

1. Concept
2. Basics
3. Implications
4. Paradox
5. Optimizing



https://www.frontiersin.org/articles/10.3389/fnhum.2011.00002/full

# Concept: Reliability

# RELIABILITY



Reliable
Not Valid

Low Reliability
Low Validity

Not Reliable
Not Valid

Reliable
Valid

https://www.scienceforsport.com/reliability/

# Reliability (statistics)

From Wikipedia, the free encyclopedia

*For other uses, see Reliability.*

**Reliability** in statistics and psychometrics is the overall consistency of a measure.[1] A measure is said to have a high reliability if it produces similar results under consistent conditions. "It is the characteristic of a set of test scores that relates to the amount of random error from the measurement process that might be embedded in the scores. Scores that are highly reliable are accurate, reproducible, and consistent from one testing occasion to another. That is, if the testing process were repeated with a group of test takers, essentially the same results would be obtained. Various kinds of reliability coefficients, with values ranging between 0.00 (much error) and 1.00 (no error), are usually used to indicate the amount of error in the scores."[2] For example, measurements of people's height and weight are often extremely reliable.[3][4]

https://en.wikipedia.org/wiki/Reliability_(statistics)

Outcomes the construct should influence

Test

**Utility**

# Basics

| test(1) | test(2) | reliability |
|---|---|---|
| Form A time 1 | Form A time 2 | Retest |
| Form A rater 1 | Form A rater 2 | Inter-rater |
| Form A | Form B | Parallel forms |
| Form A half 1 | Form A half 2 | Split half |
| Form A half $k$ | Form A half $\bar{k}$ | Internal consistency |

## Classical test theory notion of reliability

Observed *test score* is a function of a true score and (random) *error*

The *true score* is the score that would be obtained on an arbitrarily large number of repeated assessments under identical conditions

A test score reliability is the fraction of the variance in the test score that is attributable to the true score

Since we don't observe test scores under an arbitrary large number of repeated observations, we have to do some tricks to get an estimate of reliability

**Classical test theory notion of reliability**

$$x = t + e$$

$$\sigma_x^2 = \sigma_t^2 + \sigma_e^2$$

$$\rho_{xx'} = \frac{\sigma_t^2}{\sigma_x^2} = 1 - \frac{\sigma_e^2}{\sigma_x^2}$$



error

test

true_score

CrossMark

# Comparing test-retest reliability of dynamic functional connectivity methods

Ann S. Choe [a,b], Mary Beth Nebel [c,d], Anita D. Barber [e], Jessica R. Cohen [f], Yuting Xu [g], James J. Pekar [a,b], Brian Caffo [g], Martin A. Lindquist [g,*]
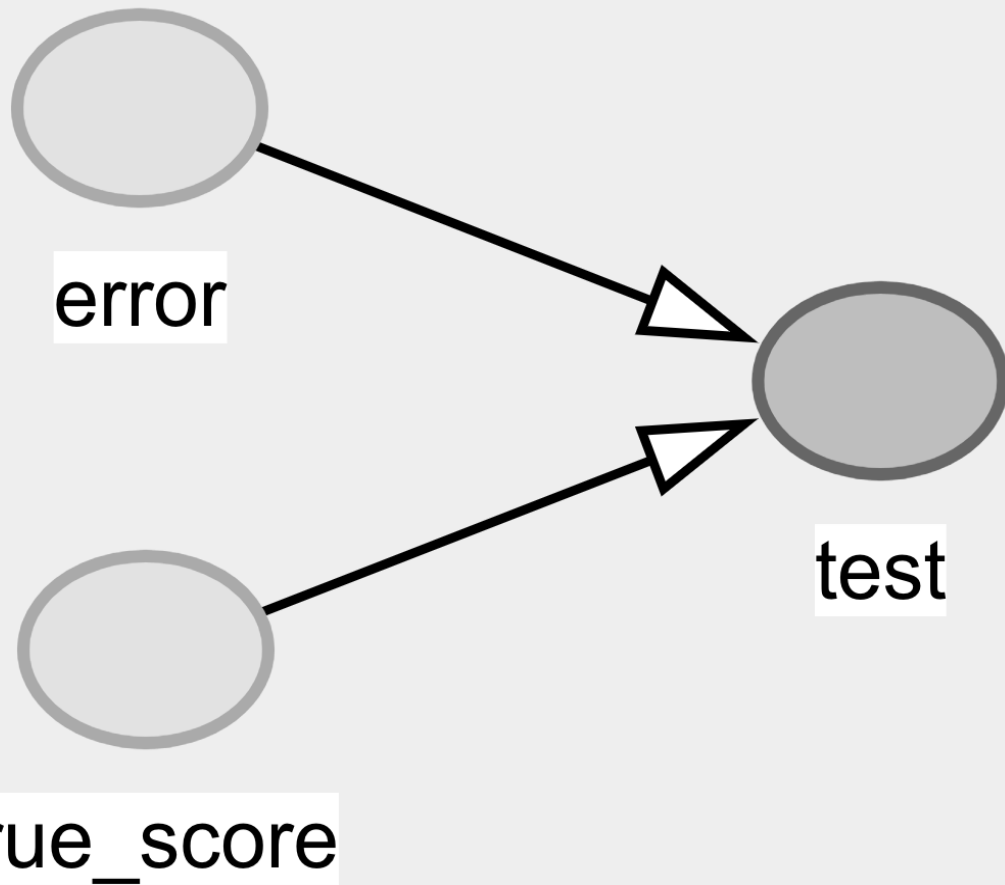
However, across all estimation methods, reliability of the brain state-derived measures was low.

*3.1.1.1. The reliability of dynamic correlation means was highly consistent* ... e mean ... oduced ... by the overlapping confidence intervals presented in the left panel of Fig. 1A, the I2C2 of dynamic correlation means was similar across all estimation methods (95% confidence intervals (CIs) for SW, TSW and DCC methods were $[0.51, 0.65]$, $[0.50, 0.64]$, and $[0.51, 0.62]$ respectively). For comparison, the 95% CI for the static correlation was $[0.52, 0.66]$.

# Implications

# Reliability standards

| Reliability | Fleiss (1981) | Landis & Koch (1970) | Nunnally & Bernstein (1994) | NIH PROMIS | Reliability |
|---|---|---|---|---|---|
| 0.0 | Poor | Poor | Inadequate for group differences research | | 0.0 |
| 0.1 | | Slight | | | 0.1 |
| 0.2 | | Fair | | | 0.2 |
| 0.3 | | | | | 0.3 |
| 0.4 | Fair to good | Moderate | | | 0.4 |
| 0.5 | | | | | 0.5 |
| 0.6 | | Substantial | | | 0.6 |
| 0.7 | Excellent | | | | 0.7 |
| 0.8 | | Almost perfect | Adequate for group-level inference | | 0.8 |
| 0.9 | | | Suitable for individual level inference | Target stopping rule for CATs | 0.9 |
| 1.0 | | | | | 1.0 |

Inter-observer agreement / Inter-observer agreement / Reliability / Scale information (Item response theory)

*If important decisions are made with respect to specific test scores, a reliability of 0.90 is the bare minimum, and a reliability of 0.95 should be considered the desirable standard.*

# Study design: sample size

# Sample size

Lehr's equation - number needed per group (n) to detect a standardized effect size (d) with type-I error level of 5% and type-II error level of 20% (16)

If d = .5 , n = 64
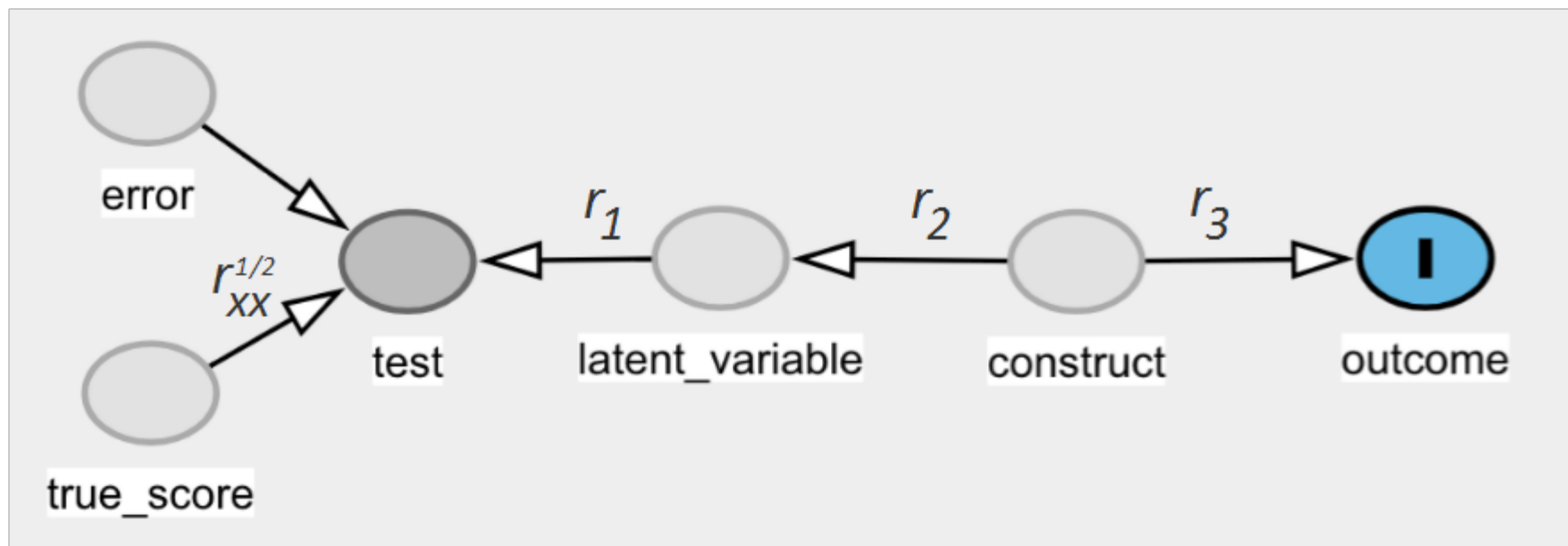
$$n = \frac{16}{d^2}$$

# Imagine you are planning a study

You have a treatment that can produce a 0.5 SD difference in the means *of the true score* across treated and non-treated participants (d = 0.5)
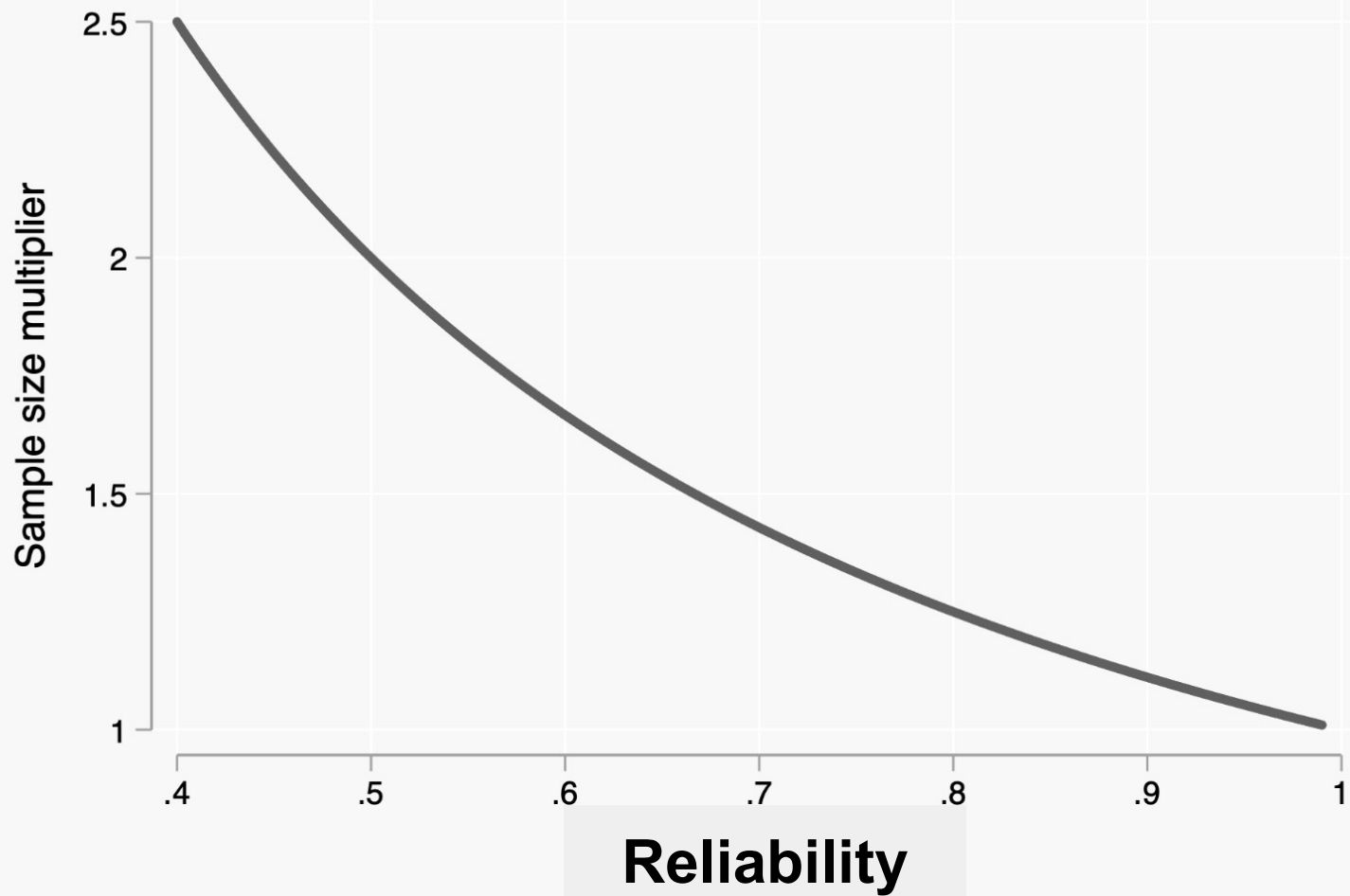
You have a measure of the outcome that is perfectly reliable (*REL* = 1.0)

How many people do you need to randomize per group to have 80% power with a two sided type-I error level of 5%

If
  d = .5 , r = .7
then
  n = 92

$$n = \frac{16}{d^2} \times \frac{1}{REL}$$

**Katie Corker** @katiecorker
Sep 8

I don't know who needs to hear this, but RELIABILITY IS A PRECONDITION FOR VALIDITY.
#noonesaidwedontcareaboutvalidity #reliabilityfirst

32    172

# Clinical work

# Imagine you are a clinician: a geriatrician (1)

1.  You want to <u>identify</u> older adults who would be unsafe drivers because of cognitive impairment, using a test of mental status
2.  You decide that adults in the <u>worst quintile</u> of *true cognitive status* would be unsafe drivers
3.  You want to make the <u>right decision 9 times out of 10</u>

<u>How reliable</u> does your test have to be?

# Imagine you are a clinician: a geriatrician (2)

Answer:

If you consider "true positive" and "true negatives" as correct decisions (i.e., you would like a "hit rate" of 90%), you need a reliability of .8

*The PPV when reliability = .8 is about 75%*

# Imagine you are a clinician: a geriatrician (3)

Answer:

If you only consider "true positive" as correct decisions (i.e., you would like a positive predictive value of 90%), you need a reliability of .97 (!)

*The PPV when reliability  = .8 is about 75%*

- Standards for reliability depend on lots of things

- In research, highly reliable measures can improve power

- In clinical settings, highly reliable measures are the foundation of good clinical practice

MAJOR MESSAGE

**Riyan Portuguez** @riyanportuguez · Sep 3, 2018

Points in **Reliability**:

📚The correlation coefficient for clinical setting is .90 or higher.
📚 The correlation coefficient for research .70 or higher.
📚The goal in **reliability** is to reduce the degree of **measurement** error.

#BLEPP2018

💬 4          🔁 71                ♡ 209              ⬆️⬇️

# Paradox

# THE ATTENUATION PARADOX IN TEST THEORY[1]

JANE LOEVINGER
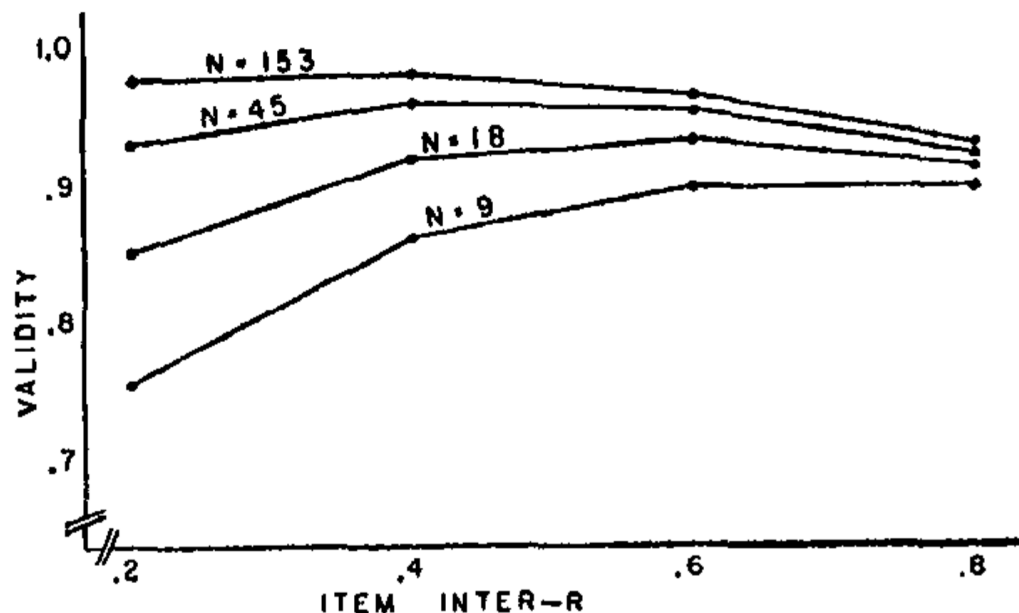
*Washington University*[2]

45

THE ATTENUATIO

JA
W



FIG. 1. ATTENUATION PARADOX AS A FUNCTION OF NUMBER OF ITEMS FOR TESTS COMPOSED OF MEDIAN EQUIVALENT ITEMS. DATA FROM BROGDEN'S (1) TABLE 2.
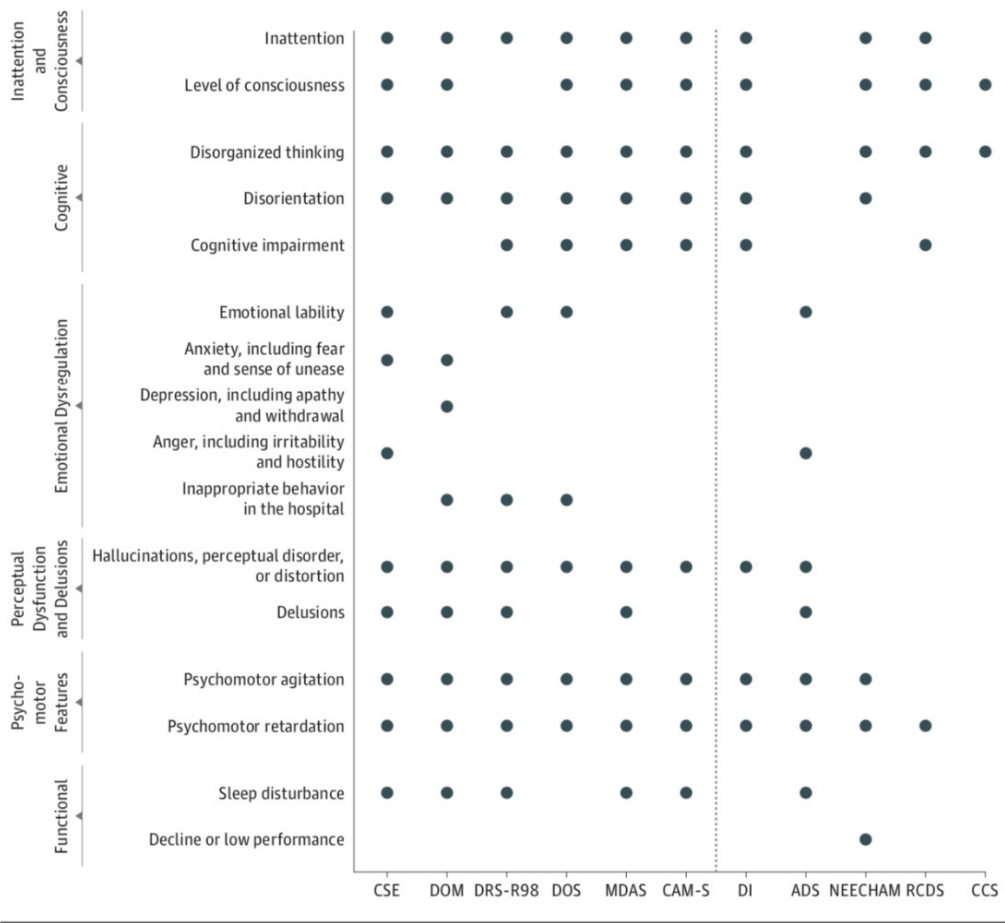
The real message behind the Attenuation Paradox is: the CTT notion of reliability is limited, and, wrong. Should use item response theory (IRT) instead.

# Optimizing

# How to optimize reliability of measures

1. Use rigorously designed and <u>validated</u> measures (c.f., COSMIN)
   - Be critical and skeptical in your review of literature. Look out for
     - "Bloated specific" measures that sacrifice fidelity for bandwidth
     - Forms pruned with "alpha if item deleted"
   - Pilot test instruments in samples from your target population
2. Resist the urge to use *fixed* short forms
3. Satisfy the urge to use short forms with Computerized Adaptive Tests
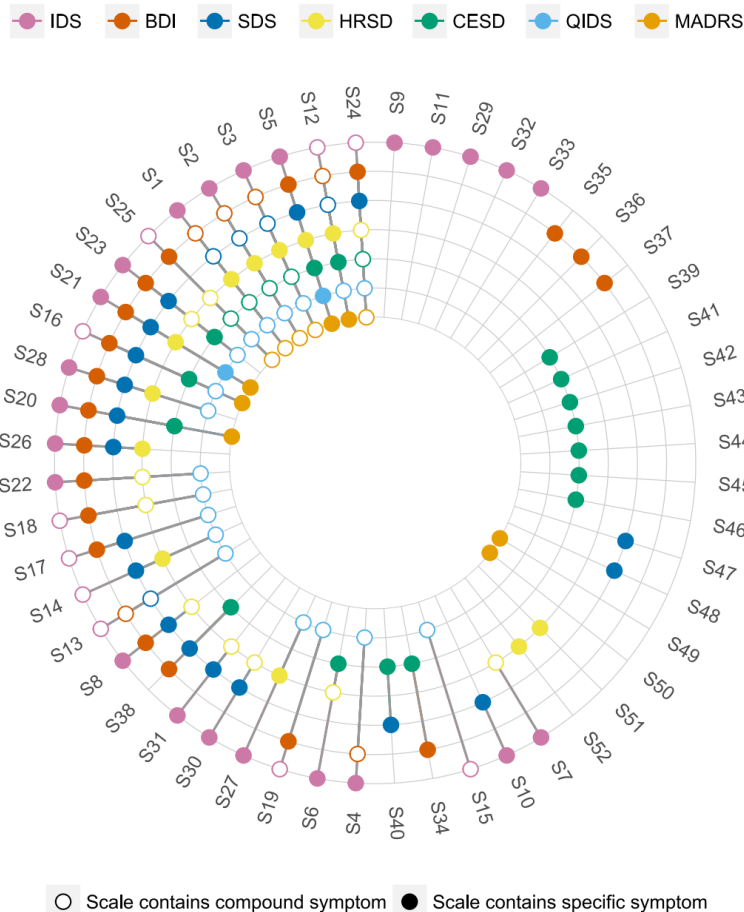4. Staff training and continual quality assessment in conducting research

Figure 2. Domain Coverage of 11 Multi-item Delirium Severity Instruments

**Assessment of Instruments for Measurement of Delirium Severity: A Systematic Review**

Jones RN et al. JAMA Intern Med. 2019:179(2);231-239

Black dot indicates representation of a domain in the instrument; either partial or full coverage of a domain met criteria for inclusion by the expert panel. ADS indicates Agitation Distress Scale; CAM-S, Confusion Assessment Method–Severity Score; CCS, Communication Capacity Scale; CSE, Confusion State Examination; DI, Delirium Index; DOM, Delirium-O-Meter; DOS, Delirium Observation Screening; DRS-R98, Delirium Rating Scale-Revised-98; MDAS, Memorial Delirium Assessment Scale; NEECHAM, Neelon and Champagne Confusion Scale; RCDS, Recoverable Cognitive Dysfunction Scale.

Legend: IDS, BDI, SDS, HRSD, CESD, QIDS, MADRS

1 Early insomnia
2 Middle insomnia
3 Late insomnia
4 Hypersomnia
5 Sad mood
6 Anxious
7 Panic
8 Irritable
9 Mood reactivity
10 Diurnal variation
11 Grief
12 Appetite decrease
13 Appetite increase
14 Weight decrease
15 Weight increase
16 Concentration
17 Indecisiveness
18 Guilt
19 Worthlessness
20 Pessimism
21 Suicidal ideation
22 Interest loss
23 Pleasure loss
24 Fatigue
25 Energy loss
26 Libido
27 Retardation
28 Agitation
29 Somatic complaints
30 Sympathetic arousal
31 Gastrointestinal
32 Interpersonal sensitivity
33 Leaden paralysis
34 Past failure
35 Punishment
36 Self-dislike
37 Self-criticalness
38 Crying
39 Lonely
40 Effort
41 Talked less
42 People are unfriendly
43 People disliked me
44 Feeling bothered
45 Feeling good
46 Feeling happy
47 Feeling needed
48 Life is full
49 Inner tension
50 Inability to feel
51 Hypochondriasis
52 Loss of insight

The 52 symptoms of major depression: Lack of content overlap among seven common depression scales.

Fried EI. J Affect Disord. 2017.

https://www.ncbi.nlm.nih.gov/m/pubmed/27792962/

○ Scale contains compound symptom ● Scale contains specific symptom

**Fig. 1.** Co-occurrence of 52 depression symptoms across 7 depression rating scales. Colored circles for a symptom indicate that a scale directly assesses that symptom, while empty circles indicate that a scale only measures a symptom indirectly. For instance, the IDS assesses item 4 hypersomnia directly; the BDI measures item 4 indirectly via a general question on sleep problems; and the SDS does not capture item 4 at all. Note that the 9 QIDS items analyzed correspond exactly to the DSM-5 criterion symptoms for MDD. Please see the online version for colors; in the black and white version, the circles respresent (from outer to inner circle): IDS, BDI, SDS, HRSD, CESD, QIDS, and MADRS.
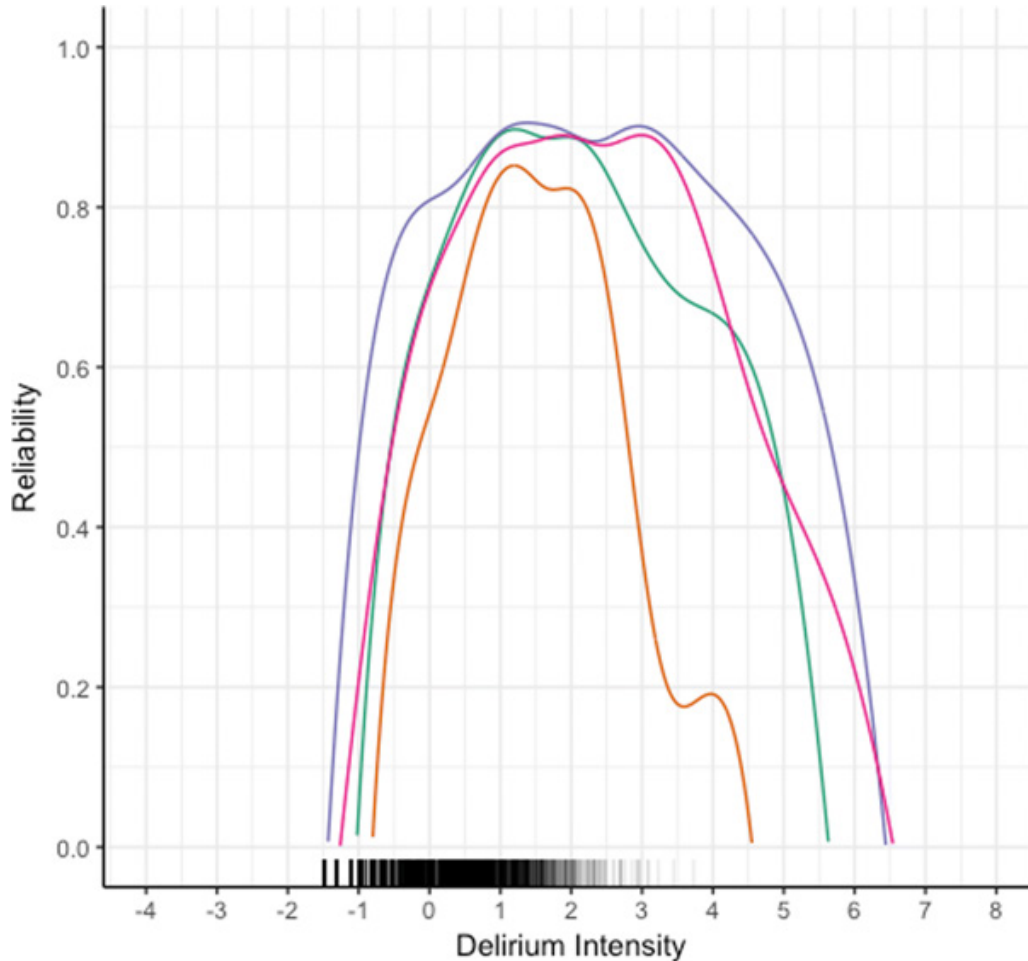
# Harmonization of delirium severity instruments: a comparison of the DRS-R-98, MDAS, and CAM-S using item response theory

Alden L. Gross[1,2]*, Doug Tommet[3], Madeline D'Aquila[4], Eva Schmitt[4], Edward R. Marcantonio[4,5], Benjamin Helfand[3,6], Sharon K. Inouye[4,5†], Richard N. Jones[3†] and for the BASIL Study Group

**Delirium Instrument**
- DRS-R-98
- MDAS
- CAM-S LF
- CAM-S SF

# Discussion and questions

Test (A)

Reference standard

Latent variable

$k_a$

$g_r$

$b$

Construct

$k_b$

$g_o$

Test (B)

Outcome