

Ch. 4: Reliability

653

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Reliability

- History
- Classical Test Score Theory
 - Domain Sampling
 - Models of reliability
 - Sources of error
- Estimating Reliability
 - Test-Retest
 - Parallel Forms
 - Internal Consistency / Cronbach's α
- Difference Scores

692

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Constructs & Measurement

- Psychology as “soft science”
- Construct
 - exists but can't be directly measured
 - examples
- Measurement
 - “true value” - intelligence
 - measured or *observed* value (e.g. IQ test score)
 - discrepancy - “error”
- How to conceptualize *error*?

693

Psychology 402 - Fall 2018 - Dr. Michael Daeher

History 1

- 1896 - Karl Pearson - product-moment correlation (for continuous variables)
- 1904 - Charles Spearman - “*The proof and measurement of association between two things*” - *Rho* - correlation for Ordinal variables

694

Psychology 402 - Fall 2018 - Dr. Michael Daeher

History

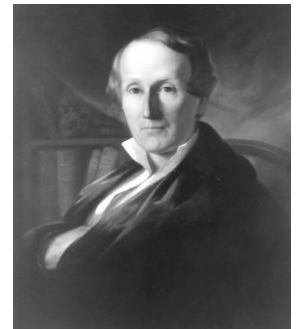
- Pearson, Spearman, Thorndike (1900-1907)
 - Basic reliability theory
- Kuder, Richardson (1937), Cronbach (1989)
 - Reliability coefficients
- Bartholomew & Knott (1990s)
 - Latent variable theory
- Drasgow et al (late 1990s)
 - Item Response Theory (IRT)

695

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Samuel George Morton

- Polygenism
 - Humans are composed of different species
- Craniometry
- Biological Determinism
- “Scientific Racism”
- d. 1851
- 50 years before Spearman's work



696

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Classical Test-Score Theory

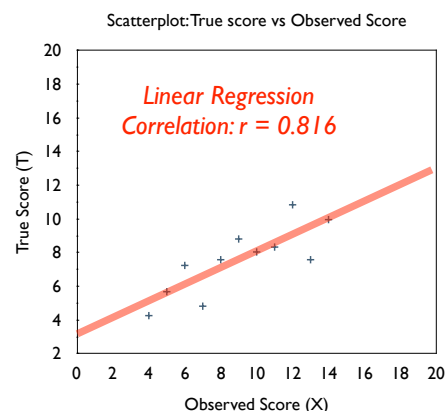
- True score (T) : the “actual” score that exists
- Observed score (X) : score as measured by a test
- Error (E) : difference between Observed and True score
- $X = T + E$
- $E = X - T$
- Assumptions: True scores have no variability. Errors are random (e.g. a normal distribution with mean of zero)
- Reliability = correlation between Observed score and True score

697

Psychology 402 - Fall 2018 - Dr. Michael Dohr

Classical Test Score Theory

X	T
10	8.04
8	7.58
13	7.58
9	8.81
11	8.33
14	9.96
6	7.24
4	4.26
12	10.84
7	4.82
5	5.68

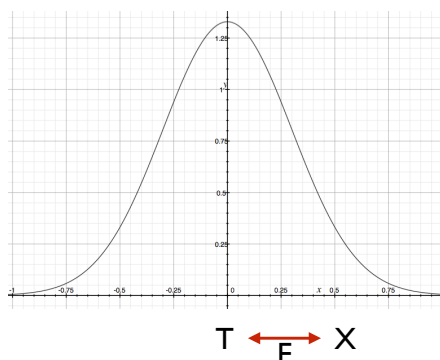


698

Psychology 402 - Fall 2018 - Dr. Michael Dohr

Classical Test-Score Theory

- T= True Score
- X = Observed
- E = Error
- $X = T + E$
- $E = X - T$

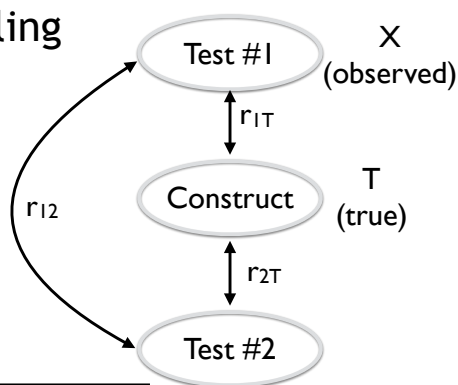


Psychology 402 - Fall 2018 - Dr. Michael Dohr

699

Domain Sampling

- How to calculate r_{1T}
- Any two tests r_{12}
- r_{1j} = average of all pairs



$$r_{1,T} = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N r_{i,j}}{N^2}}$$

Psychology 402 - Fall 2018 - Dr. Michael Dohr

701

Domain Sampling

- Problem: no way to measure True score / no possible way to measure every possible item
- Sample a limited subset of items, do this in multiple ways
- Create one or more tests
- For two given tests, correlation between the two tests will be lower than the correlation between one test and the True score
- $r_{1t} = \sqrt{r_{1j}}$

702

Psychology 402 - Fall 2018 - Dr. Michael Dohr

Domain Sampling Example

- Correlation of any 2 random sample tests
- $r_{1t} = \sqrt{r_{1j}}$
- $r_{1t} = \sqrt{0.64}$
- $r_{1t} = 0.80$
 - unbiased estimate of “true” reliability

703

Psychology 402 - Fall 2018 - Dr. Michael Dohr

Models of Reliability

- Most reliability measures are Correlation coefficients
- Alternate definition: Reliability is the ratio of the variance of True scores to the variance of the Observed scores
 - $\rho^2_{XT} = \frac{\sigma^2_T}{\sigma^2_X}$
- Or, it's the "Signal to Noise" ratio
 - $\rho^2_{XT} = \frac{\sigma^2_T}{\sigma^2_T + \sigma^2_E}$
- A test with reliability of $r=0.40$ means that 40% of variation in test scores is due to variation in the "true" score, and 60% of variation is random or chance factors.

705

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Sources of Error

- "Error" is considered the difference between True score and Observed score
- Where does Error arise?
 - Measurement errors
 - Change in True score

706

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Test-Retest Reliability

- Test-Retest
 - administer same test across some time period
 - compute correlation between two administrations
 - Issue -- what is "error"?
 - actual change in true score
 - carryover or practice effects

707

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Parallel Forms Reliability

- Also called "Alternate Forms"
 - administer two versions of the test to same subjects (often on same day)
 - compute correlation between two administrations
- Pros: more rigorous method of determining reliability
- Cons: difficult to do, is not often done

708

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Internal Consistency Reliability

- Give single test, calculate internal consistency of various subsets of items
- Split halves methods exist, but have generally been supplanted by...
- Cronbach's Alpha (α)
 - estimates a lower bound for reliability
 - α of .70 to .80 is borderline
 - α of .80 is ok
 - α of .90 or higher is good

709

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Inter-Rater Reliability

- Observational data differs from self-report data.
- Even though most behavioral rating systems attempt to be precise, errors occur (e.g. was that a "hit" or a "punch"?)
- We must consider the reliability of different observers (also called "raters")
- Cohen's Kappa
 - ranges from -1 to +1
 - "poor" < .40
 - "good" .40 to .75
 - "excellent" > .75

710

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Reliability: errors & methods

	Description	Name	Statistic
Time Sampling	1 test given two times	test-retest reliability	correlation between scores at two times
Item Sampling	2 different tests given once	Alternate or Parallel forms	correlation between scores on 2 versions
Internal Consistency	One test, multiple items	Split Half or internal reliability	Cronbach's Alpha
Observer Differences	One test w/ 2+ observers	inter-observer reliability	Kappa

711

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Summary

- Reliability
 - how consistent measured scores are
- Error
 - $E = X - T$
- What kind of Error?
 - test-retest, domain sampling, internal consistency, observer-differences

720

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Standard Error of Measurement

- Desire to answer question “how close is this test result to the true result”
- If we know the Reliability (r) of the test, we can estimate the likely range of true values
- Given
 - S = std dev of measured scores
 - r = reliability coefficient of test

$$SEM = S\sqrt{1 - r}$$

726

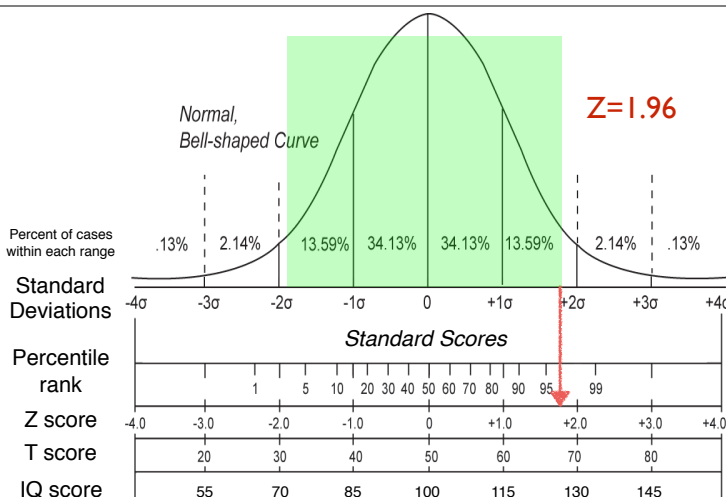
Psychology 402 - Fall 2018 - Dr. Michael Daeher

SEM example: IQ

- Example: a person scored 106 on an IQ test, that has a reliability of 0.89. What is the 95% confidence interval of the their true score
- $S = 14$
- $r = 0.89$
- $SEM = S\sqrt{1 - r}$
- $SEM = 14\sqrt{1 - 0.89}$
- $SEM = 4.64$
- Next, compute a confidence interval

727

Psychology 402 - Fall 2018 - Dr. Michael Daeher



728

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Confidence Interval

- “How likely is a true score to fall within a range”
- Z = z-score associated with % range
- Confidence interval = $Z * SEM$
- Example:
 - 95% confidence interval : $Z = 1.96$
 - $SEM = 4.64$
 - $1.96 * 4.64 = 9.1$
 - 95% CI = ± 9.1 points
 - Range = $X \pm CI$
 - $106 \pm 9.1 = \text{range from } 96.9 \dots 115.1$

729

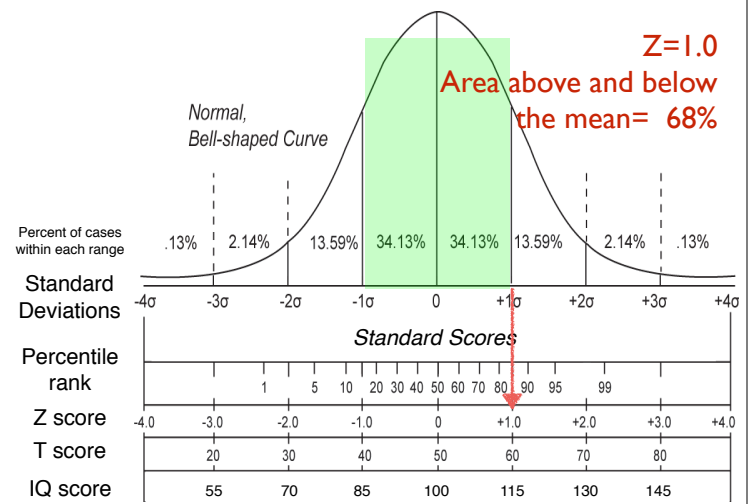
Psychology 402 - Fall 2018 - Dr. Michael Daeher

Common Z scores & Confidence Intervals

Z Score	Area above mean	Area above + below Mean	Proportion as %
0.00	0.000		0%
0.13	0.051		
0.67	0.249		
1.00	0.341	0.682	68%
1.64	0.449		
1.96	0.475		95%
2.57	0.495		

733

Psychology 402 - Fall 2018 - Dr. Michael Diehr



735

Psychology 402 - Fall 2018 - Dr. Michael Diehr

Reliability of Difference Scores

- Common need is to compute the difference between two scores or two tests, with known reliability
- Unfortunately, taking the difference dramatically reduces the reliability
- E.g. for two tests with reliability .90 and .70 that are correlated to each other by .70, a difference score has a reliability of .33

737

Psychology 402 - Fall 2018 - Dr. Michael Diehr

How reliable?

- $r = .70$ or $.80$ or higher is often considered “good enough” for much research
- $r > .90$ is very good, may not be worth time / effort to get higher

738

Psychology 402 - Fall 2018 - Dr. Michael Diehr

Increasing Reliability

- Increase N (number of questions, items or tests) ...
- Focus on common characteristic...
 - tests are more reliable if all items measure a single characteristic
- Other methods (covered later)
 - Use **Factor Analysis** to determine sub-characteristics of a single test
 - Use **Item Analysis** (“discriminability analysis”) to find items that best measure a single characteristic

739

Psychology 402 - Fall 2018 - Dr. Michael Diehr

740

Psychology 402 - Fall 2018 - Dr. Michael Diehr

Increase N

- N = number of questions or items or tests
- Formulas exist to determine how much to increase N by to reach a certain level of reliability
- $N_d = r_d (1 - r_o) / r_o (1 - r_d)$
 N_d = new N (times old N)
 r_d = desired level of reliability
 r_o = observed level of reliability
- Example: 20-item CES-D has reliability of .87. We need .95. $N_d = 2.82$, so new N is $2.82 \times 20 = 56$

741

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Increase N - Examples

- $N_d = r_d (1 - r_o) / r_o (1 - r_d)$
- Example:
 - 20-item CES-D has reliability of .87. We need .95. $N_d = 2.82$, so new N is $2.82 \times 20 = 56$
 - Your 40-item test has reliability of .50. You want .90. $N_d = 9.0$, so new N is $9 \times 40 = 360$!

742

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Focus Test

- Reliability increases the more the test focuses on a single concept or characteristic
- Trying to capture multiple concepts in a single test reduces reliability
- Methods:
 - Ad-hoc / informal -- face validity of items and remove those that don't fit
 - Statistical:
 - Factor Analysis
 - Discriminability Analysis.

743

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Chapter 4 Summary

- Measurement Error occurs in all fields -- Psychology focuses on it
- Reliability : several kinds, estimate it by deciding *where* the measurement error comes from
- Unreliable tests can't be Valid.
- Improving Reliability: more items, focusing test, factor analysis
- Reliability is useful: Used to calculate SEM and Confidence Intervals

745

Psychology 402 - Fall 2018 - Dr. Michael Daeher

Reliability Summary

- Reliability: consistency of measurement
- Source of error → how to measure reliability
- Reliability coefficients are correlations
- Reliability is not Validity
- Reliability is a foundation upon which Validity can be built

746

Psychology 402 - Fall 2018 - Dr. Michael Daeher