

Validity

Jens Søndergaard Jensen
M.Sc. in Statistics

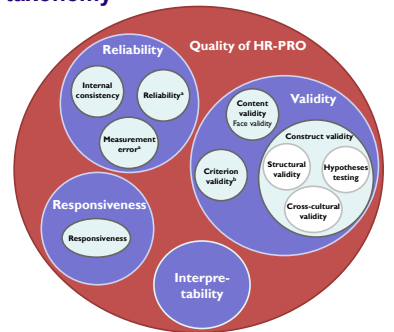
Outline

- Secrets to measuring a piece of paper
- The concept of validity
- Content and face validity
- Criterion validity
 - Concurrent validity
 - Predictive validity
- Construct validity
 - Structural validity
 - Hypotheses testing
 - Cross-cultural validity

Secrets to measuring a piece of paper

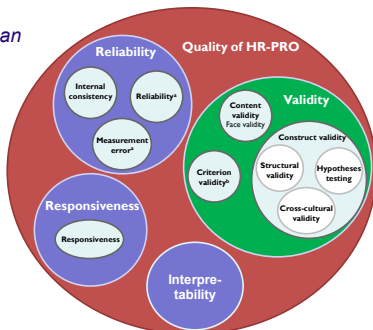
<https://www.youtube.com/watch?v=9yUZTTLpDtk>

The COSMIN taxonomy



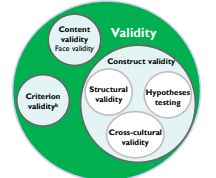
Validity

"The degree to which an instrument truly measures the construct(s) it purports to measure"



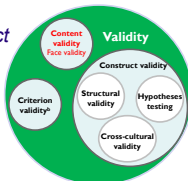
Validity

- Consists of testing hypotheses about the construct
- Important implications:
 - Knowledge about the construct
 - Complexity of the construct
 - Dependency on the situation
 - Validation of scores
 - Not measurement instruments
 - Formulation of specific hypotheses
 - Validation is a continuous process



Content validity

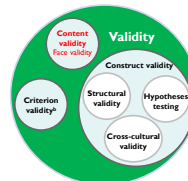
- "The degree to which the content of a measurement instrument is an adequate reflection of the construct to be measured"
 - Is the instrument relevant?
 - Is the instrument comprehensive?



The Research Clinic for Functional Disorders and Psychosomatics

Face validity

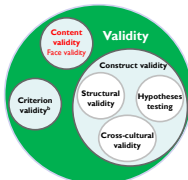
- "The degree to which a measurement instrument, indeed, looks as though it is an adequate reflection of the construct to be measured"
 - A subjective assessment
 - No standards with regard to assessment
 - Lack of face validity → strong argument for not using an instrument



The Research Clinic for Functional Disorders and Psychosomatics

Content and face validity

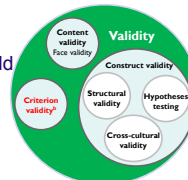
- The proces of content validation
 1. Consider information about construct and situation
 2. Consider information about content of the instrument
 3. Select an expert panel
 4. Assess whether content of the instrument corresponds with the construct (is it relevant and comprehensive)
 5. Use a strategy or framework to assess correspondence between the instrument and construct



The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity

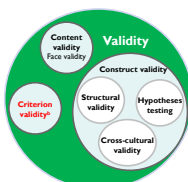
- "The degree to which the scores of an instrument are an adequate reflection of a gold standard"
- Can only be assessed when a gold standard is available
- What is a gold standard?



The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, subtypes

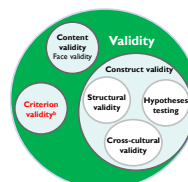
- Concurrent validity
 - Consider the score for the instrument and criterion at the same time
 - Usually assessed for instruments used for evaluative and diagnostic purposes
- Predictive validity
 - Consider whether the instrument predicts the criterion in the future
 - Used in predictive applications



The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity

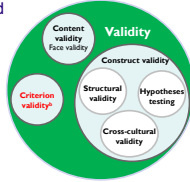
Hypothesis: The instrument under study is as good as the gold standard!



The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity

- The process of criterion validation
 1. Identify a suitable criterion and method of measurement
 2. Identify an appropriate sample of the target population
 3. Define a priori the required level of agreement between instrument and criterion
 4. Obtain scores for instrument and criterion independently
 5. Determine the strength of relationship between instrument scores and criterion scores (see next slide)



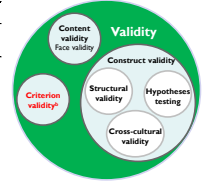
The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity

Table 6.2. Overview of statistical parameters for various levels of measurement for the gold standard and measurement instrument under study

| Level of measurement | Measurement instrument | Same units | Statistical parameter |
|----------------------|------------------------|------------|----------------------------------|
| Dichotomous | Dichotomous | Yes | Sensitivity and specificity |
| | Ordinal | NA | ROC |
| | Continuous | NA | ROC |
| Ordinal | Ordinal | Yes | Weighted kappa |
| | Continuous | NO | Spearman's r ROC/Spearman's r |
| Continuous | Continuous | Yes | Bland-Altman/ICC |
| | Continuous | No | Spearman's r or Pearson's r |

H. C. W. de Vet, et al., Measurement in Medicine p. 163



The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, dichotomous instruments

- Sensitivity**
 - Sens = $TP / (TP + FN)$
 - Probability of a positive test given that the patient is truly sick
 - The true positive rate
 - Quantifies the avoiding of false negatives
- Specificity**
 - Spec = $TN / (FP + TN)$
 - Probability of a negative test given that the patient is well
 - The true negative rate
 - Quantifies the avoiding of false positives

Table Cross-tabulation

| Gold standard | + | - | |
|---------------|----|----|---|
| + | TP | FP | TP = true positive FP = false positive |
| - | FN | TN | FN = false negative TN = true negative |

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, dichotomous instruments

- Sensitivity and specificity**
 - A perfect test has Sensitivity = Specificity = 100%
 - They are characteristics of the test (and hence they are...)
 - Independent on the prevalence of disease in the population

Table Cross-tabulation

| Gold standard | + | - | |
|---------------|----|----|---|
| + | TP | FP | TP = true positive FP = false positive |
| - | FN | TN | FN = false negative TN = true negative |

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, dichotomous instruments

- Positive and negative predicted value (PPV/NPV)**
 - $PPV = TP / (TP + FP) = P(\text{Sick} | \text{Positive test})$
 - $NPV = TN / (TN + FN) = P(\text{Not sick} | \text{Negative test})$
 - Both depend on the prevalence of disease in the population
 - Example: $N=4000$, Sens = Spec = 0.99
 - Prevalence = 50% \Rightarrow PPV = 99%
 - Prevalence = 10% \Rightarrow PPV = 92%
 - Prevalence = 5% \Rightarrow PPV = 84%
 - Prevalence = 1% \Rightarrow PPV = 52%

Table Cross-tabulation

| Gold standard | + | - | |
|---------------|----|----|---|
| + | TP | FP | TP = true positive FP = false positive |
| - | FN | TN | FN = false negative TN = true negative |

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, dichotomous instruments

- Example (Section 6.4.1)**
 - Computations are easily done by hand
 - BUT: CI's are more difficult
 - Stata can do it for us
 - Download the `diagt-` command

Table Cross-tabulation

| Gold standard | + | - | |
|---------------|----|-----|--|
| + | 30 | 114 | |
| - | 3 | 822 | |

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, dichotomous instruments

diag1 30 3 114 822

| True disease status | Test result | | Total |
|---------------------|-------------|------|-------|
| | Neg. | Pos. | |
| Normal | 822 | 114 | 936 |
| Abnormal | 3 | 33 | 33 |
| Total | 825 | 144 | 969 |

| | | Gold standard | |
|-------------|---|---------------|-----|
| | | + | - |
| Test result | + | 30 | 114 |
| | - | 3 | 822 |

[95% Confidence Interval]

| | | | | |
|---------------------------|-------------------|-------|-------|-------|
| Prevalence | Pr(A) | 2.48 | 2.48 | 4.75% |
| Sensitivity | Pr(+ A) | 90.9% | 75.7% | 98.1% |
| Specificity | Pr(- N) | 87.8% | 85.6% | 89.8% |
| ROC area | (Sens. + Spec.)/2 | 89% | 84% | 94% |
| Likelihood ratio (+) | Pr(+ A)/Pr(+ N) | 7.46 | 6.09 | 9.14 |
| Likelihood ratio (-) | Pr(- A)/Pr(- N) | .104 | .0352 | .305 |
| Odds ratio | LR(+)/LR(-) | 72.1 | 23 | 226 |
| Positive predictive value | Pr(A +) | 20.8% | 14.5% | 28.4% |
| Negative predictive value | Pr(N -) | 99.6% | 98.9% | 99.9% |

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, continuous instruments

Bland-Altman method

- Plot the difference against the average
- Compare the two scores using a paired t-test
- Calculate 95% limits of agreement ($\delta \pm 1.96\text{-SD}$)
- Add horizontal lines at $y = 0$, $y = \delta$ (the mean difference) and at each of the limits of agreement
- δ is a measure of the systematic difference (error) between the instrument and the criterion
- 1.96-SD is a measure of the random error
- Hence $\delta \pm 1.96\text{-SD}$ indicates the size of the measurement error

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, continuous instruments

- Bland-Altman plot and limits of agreement
 - Assumptions
 - The differences should be independent
 - The differences should have the same distribution
 - The differences should be normally distributed
- One article you should read:
 - Bland JM, Altman DG. *Measuring agreement in method comparison studies*. Stat Methods Med Res 1999;8(2):135-60.

The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, continuous instruments

Example

- PEFR measurements on 17 persons
- Each person measured twice; one time with Mini Wright peak-flow meter and one time using a Standard Wright peak-flow meter
- Aim:
 - To assess criterion validity by comparing the test instrument (Mini Wright) to the criterion (Standard Wright)
 - If the two PEFR meters were unlikely to give readings which differed by more than, say, 10 l/min, we would have criterion validity

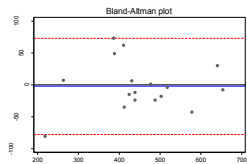
The Research Clinic for Functional Disorders and Psychosomatics

Criterion validity, continuous instruments

Example

```

1. see https://www.youtube.com/watch?v=K18L1G1L1G1
2. opening definitions
3. open definitions
4. measure agreement
5. open definitions
6. limits of agreement (precision interval)
7. measure agreement (precision interval)
8. measure agreement (precision interval)
9. measure agreement (precision interval)
10. measure agreement (precision interval)
11. measure agreement (precision interval)
12. measure agreement (precision interval)
13. measure agreement (precision interval)
14. measure agreement (precision interval)
15. measure agreement (precision interval)
16. measure agreement (precision interval)
17. measure agreement (precision interval)
18. measure agreement (precision interval)
19. measure agreement (precision interval)
20. measure agreement (precision interval)
21. measure agreement (precision interval)
22. measure agreement (precision interval)
23. measure agreement (precision interval)
24. measure agreement (precision interval)
25. measure agreement (precision interval)
26. measure agreement (precision interval)
27. measure agreement (precision interval)
28. measure agreement (precision interval)
29. measure agreement (precision interval)
30. measure agreement (precision interval)
31. measure agreement (precision interval)
32. measure agreement (precision interval)
33. measure agreement (precision interval)
34. measure agreement (precision interval)
35. measure agreement (precision interval)
36. measure agreement (precision interval)
37. measure agreement (precision interval)
38. measure agreement (precision interval)
39. measure agreement (precision interval)
40. measure agreement (precision interval)
41. measure agreement (precision interval)
42. measure agreement (precision interval)
43. measure agreement (precision interval)
44. measure agreement (precision interval)
45. measure agreement (precision interval)
46. measure agreement (precision interval)
47. measure agreement (precision interval)
48. measure agreement (precision interval)
49. measure agreement (precision interval)
50. measure agreement (precision interval)
51. measure agreement (precision interval)
52. measure agreement (precision interval)
53. measure agreement (precision interval)
54. measure agreement (precision interval)
55. measure agreement (precision interval)
56. measure agreement (precision interval)
57. measure agreement (precision interval)
58. measure agreement (precision interval)
59. measure agreement (precision interval)
60. measure agreement (precision interval)
61. measure agreement (precision interval)
62. measure agreement (precision interval)
63. measure agreement (precision interval)
64. measure agreement (precision interval)
65. measure agreement (precision interval)
66. measure agreement (precision interval)
67. measure agreement (precision interval)
68. measure agreement (precision interval)
69. measure agreement (precision interval)
70. measure agreement (precision interval)
71. measure agreement (precision interval)
72. measure agreement (precision interval)
73. measure agreement (precision interval)
74. measure agreement (precision interval)
75. measure agreement (precision interval)
76. measure agreement (precision interval)
77. measure agreement (precision interval)
78. measure agreement (precision interval)
79. measure agreement (precision interval)
80. measure agreement (precision interval)
81. measure agreement (precision interval)
82. measure agreement (precision interval)
83. measure agreement (precision interval)
84. measure agreement (precision interval)
85. measure agreement (precision interval)
86. measure agreement (precision interval)
87. measure agreement (precision interval)
88. measure agreement (precision interval)
89. measure agreement (precision interval)
90. measure agreement (precision interval)
91. measure agreement (precision interval)
92. measure agreement (precision interval)
93. measure agreement (precision interval)
94. measure agreement (precision interval)
95. measure agreement (precision interval)
96. measure agreement (precision interval)
97. measure agreement (precision interval)
98. measure agreement (precision interval)
99. measure agreement (precision interval)
100. measure agreement (precision interval)
    
```

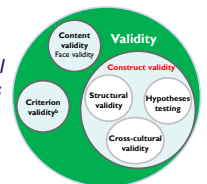


- Conclusion
 - $\delta = -2.1$, 95% CI (-22.0 ; 17.8)
 - Limits of agreement (-78.1 ; 73.9)
 - Criterion validity not OK

The Research Clinic for Functional Disorders and Psychosomatics

Construct validity

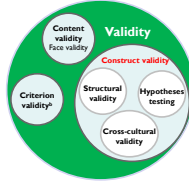
"The degree to which the scores of an instrument are consistent with hypotheses, e.g. with regard to internal relationships, relationships with scores of other instruments or differences between relevant groups"



The Research Clinic for Functional Disorders and Psychosomatics

Construct validity

- Assumes the instrument validly measures the construct
- Three aspects
 - Structural validity
 - Hypotheses testing
 - Cross-cultural validity

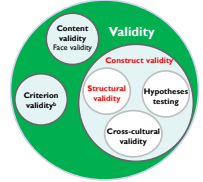


The Research Clinic for Functional Disorders and Psychosomatics

Structural validity

"The degree to which the scores of a measurement instrument are an adequate reflection of the dimensionality of the construct to be measured"

- Assessed by
- Confirmatory Factor Analysis (CFA)
 - Item Response Theory (IRT) (Indirectly)



The Research Clinic for Functional Disorders and Psychosomatics

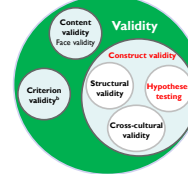
Factor analysis

- Exploratory Factor Analysis (EFA)
 - Used when:
 - You have no idea of the number of factors
 - You have no idea of which items belong to each factor
- Confirmatory Factor Analysis (CFA)
 - A statistical model
 - Use fit-parameters to assess whether data fit a hypothesized factor structure
 - Possible to compare fit of several models
 - Used when:
 - You know the number factors (there should be in theory)
 - You know which items belong to each factor (theoretically)

The Research Clinic for Functional Disorders and Psychosomatics

Hypothesis testing

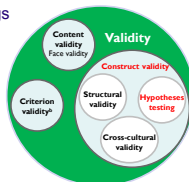
- Describe the construct in high detail (use the conceptual model)
- Formulate hypotheses about:
 - Positive correlations with similar constructs (**Convergent validity**)
 - No/slight correlation with unrelated constructs. (**Discriminant validity**)
 - Differences between subgroups of patients (**Discriminative validity**)
 - Focus on effect sizes and CI's; not statistical significance



The Research Clinic for Functional Disorders and Psychosomatics

Hypothesis testing

- Describe the instruments to which your instrument is compared.
- Describe characteristics of subgroups to be discriminated.
- Gather empirical data
- Test your hypotheses – are results consistent with hypotheses?
- Discuss your findings



The Research Clinic for Functional Disorders and Psychosomatics

Convergent/discriminative validity

- Example
 - Validation of several questionnaires used to assess functional health status in children with Acute Otitis Media (AOM)
 - 383 children with AOM
 - 7 questionnaires assessing functional health status
 - SF-36 (functioning and emotional behavior)
 - FSQ generic (age appropriate functioning and emotional behavior)
 - FSQ specific (general impact of illness on -----||-----)
 - OM-6 (physical functioning)
 - NRS child (HRQoL of child)
 - NRS caregiver (HRQoL of caregiver)
 - FFQ (Family functioning)

The Research Clinic for Functional Disorders and Psychosomatics

Convergent/discriminant validity

- Hypotheses
 - Weak correlation ($r = 0.1 - 0.3$) between FSQ generic and NRS caregiver
 - Moderate to strong correlation ($r > 0.4$) between SF-36 and NRS caregiver
 - Moderate to strong correlation ($r > 0.4$) between OM-6, FSQ specific, NRS child, NRS caregiver and FFQ

Table 6.5: Construct validity: correlations between questionnaires

| | RAND | FSQ Generic | FSQ Specific | OM-6 | NRS child | NRS caregiver | FFQ |
|---------------|------|-------------|--------------|------|-----------|---------------|------|
| RAND | 1.00 | 0.52 | 0.49 | 0.34 | 0.33 | 0.43 | 0.49 |
| FSQ Generic | | 1.00 | 0.80 | 0.37 | 0.25 | 0.43 | 0.24 |
| FSQ Specific | | | 1.00 | 0.49 | 0.26 | 0.63 | 0.24 |
| OM-6 | | | | 1.00 | 0.23 | 0.74 | 0.28 |
| NRS child | | | | | 1.00 | 0.22 | 0.47 |
| FFQ | | | | | | 1.00 | 0.39 |
| NRS caregiver | | | | | | | 1.00 |

The Research Clinic for Functional Disorders and Psychosomatics

Convergent/discriminant validity

- Conclusions
 - NRS child does not perform as hypothesized
 - NRS caregiver shows lower correlation, with FSQ specific and OM-6, than hypothesized

Table 6.6: Construct validity: correlations between questionnaires

| | RAND | FSQ Generic | FSQ Specific | OM-6 | NRS child | NRS caregiver | FFQ |
|---------------|------|-------------|--------------|------|-----------|---------------|------|
| RAND | 1.00 | 0.52 | 0.49 | 0.34 | 0.33 | 0.43 | 0.49 |
| FSQ Generic | | 1.00 | 0.80 | 0.37 | 0.25 | 0.43 | 0.24 |
| FSQ Specific | | | 1.00 | 0.49 | 0.26 | 0.63 | 0.24 |
| OM-6 | | | | 1.00 | 0.23 | 0.74 | 0.28 |
| NRS child | | | | | 1.00 | 0.22 | 0.47 |
| FFQ | | | | | | 1.00 | 0.39 |
| NRS caregiver | | | | | | | 1.00 |

The Research Clinic for Functional Disorders and Psychosomatics

Discriminative validity

- Hypotheses:
 - Children with ≥ 4 episodes of AOM per year ($n=242$) have lower scores on all instruments than children with only two and three episodes per year ($n=141$)
 - NO magnitude was specified

Table 6.6: Discriminative validity scores of children with 2-3 versus 4 or more episodes of AOM in the previous year

| | 2-3 AOM episodes | ≥ 4 AOM episodes | p-value |
|------------------|------------------|-----------------------|---------|
| Generic | 21.1 | 19.6 | 0.004 |
| RAND SF-36 | 76.5 | 72.2 | 0.002 |
| FSQ Generic | 83.9 | 78.4 | 0.001 |
| FSQ Specific | | | |
| Disease specific | | | |
| OM-6 | 18.9 | 17.0 | <0.001 |
| NRS child | 5.2 | 5.4 | 0.48 |
| FFQ | 84.9 | 78.5 | <0.001 |
| NRS caregiver | 6.6 | 6.2 | 0.22 |

P-values calculated by Mann-Whitney test

The Research Clinic for Functional Disorders and Psychosomatics

Discriminative validity

- Conclusions:
 - Statistical significant difference between the two groups on all instruments except the two NRS questionnaires
- Overall conclusions:
 - The NRS questionnaires does not perform as well as expected
 - Results from convergent/discriminative validity support each other
 - Only the NRS questionnaires showed poor convergent validity and low to moderate discriminative validity

The Research Clinic for Functional Disorders and Psychosomatics

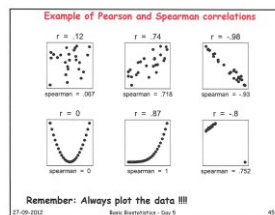
A warning about correlations

$$r = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$$

- Pearson's r
 - Measures strength of a linear association between two continuous normally distributed variables
- Spearman's ρ
 - Measures strength of a monotone association between two variables
 - The non-parametric version of Pearson's r. Uses the rank of observations instead of the actual observations
- Statistical significance (ie. $p \neq 0$) is rarely important
- Remember warning from Basic Biostatistics

The Research Clinic for Functional Disorders and Psychosomatics

Correlations



Correlations some comments

The Pearson correlation is only a valid measure of association if:

- We have independent observations, i.e. the pairs (x_i, y_i) are independent.
- Both the x 's and the y 's have a normal distribution.
- There is a linear relationship between x and y .

Note: these assumptions are stronger than the ones behind the simple linear regression.

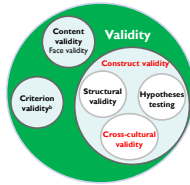
The test of no association based on Spearman rank correlation is valid if 1. and 3b. The is a monotone relationship between x and y .

27-09-2022 Basic Biostatistics - Day 5 47

The Research Clinic for Functional Disorders and Psychosomatics

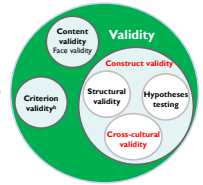
Cross-cultural validity

"The degree to which the performance of the items on a translated or culturally adapted PRO instrument are an adequate reflection of the performance of items in the original version of the instrument".



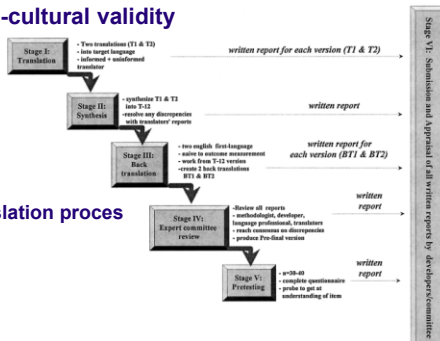
Cross-cultural validity

- Often assessed after translation of a questionnaire
- Checked by
 - Assessing its construct validity
 - Assessment of measurement invariance, differential item functioning (DIF) (CFA, logistic regression or IRT)



Cross-cultural validity

The translation process



Validation in context

- Minimal n required
 - Convergent/discriminant validity (correlations)
 - 50 patients (preferably over 100 patients)
 - Discriminative validity (t-test, rank test)
 - 50 patients per subgroup
 - Confirmatory Factor Analysis
 - 4-10 cases per item (at least 100 patients)
 - Item Response Theory
 - >200 to create stable models
 - Too many is also a problem
- Missing values
 - Should, as always, be reported and investigated
 - More than 15% missing → problems with generalizability