

# Reproducibility

Henrik Hein Lauridsen

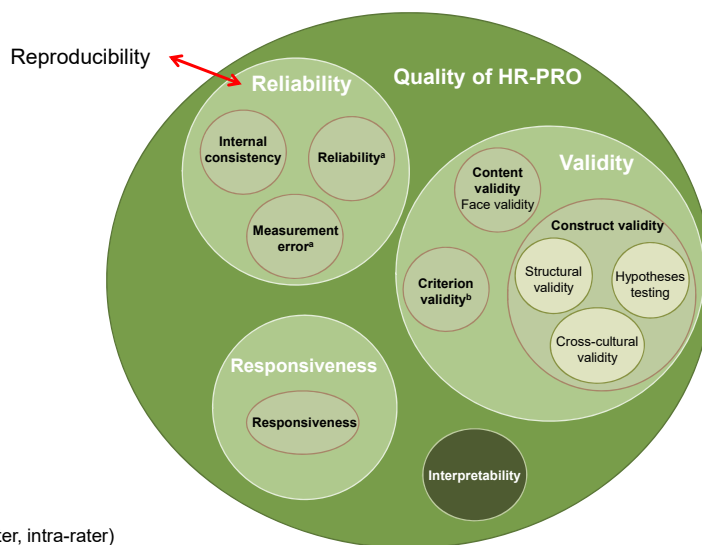
Associate Professor

Research unit for Clinical Biomechanics

University of Southern Denmark



## The COSMIN taxonomy



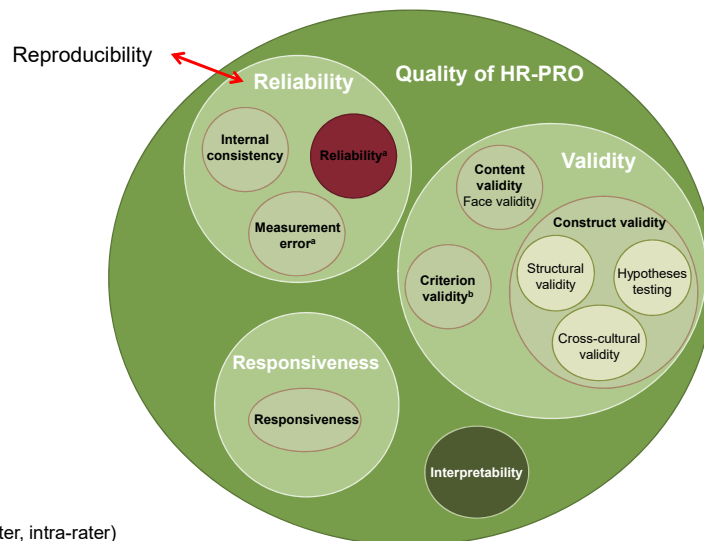
<sup>a</sup>(test-retest, inter-rater, intra-rater)

<sup>b</sup>(concurrent validity, predictive validity)

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# The COSMIN taxonomy



<sup>a</sup>(test–retest, inter-rater, intra-rater)

<sup>b</sup>(concurrent validity, predictive validity)

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



## Reproducibility

### Definition:

*"The degree to which the measurement is free from measurement error"*

*"Graden af målefejl ved en given måling"*

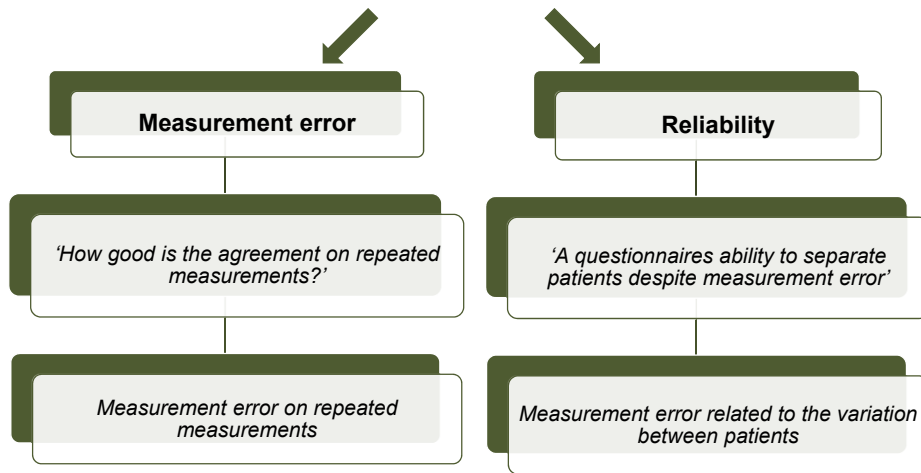
Mokkink et al. 2010

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Reproducibility

*'The degree to which the measurement is free from measurement error'*

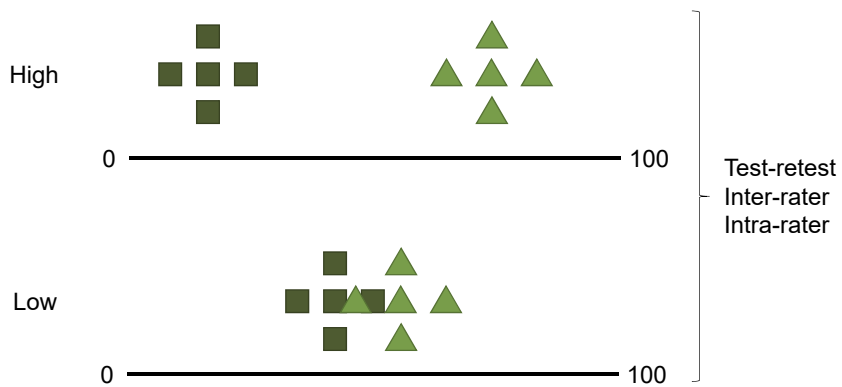


DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Reliability

5 measurements in 2 persons



DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Parameters

Scale	Parameters
Reliability	
Continuous	ICC
Ordinal	ICC or weighted kappa
Nominal	Unweighted kappa

- Scale from 0 -1
- Difficult to interpret

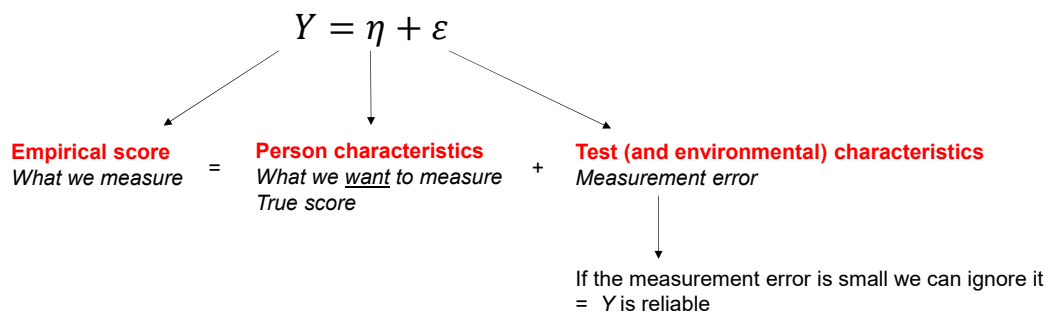
De Vet et al. "Measurement in Medicine", 2011

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Classical test theory

Assumes that the true score is the sum of the observed score and random error



DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Reliability

$$Y = \eta + \varepsilon$$

Can only be disentangled when there are repeated measures  
 $i = \text{repeated measurement (raters, occasions etc.)}$

$$Y_i = \eta + \varepsilon_i$$

**Assumption:**  
 Rewritten to variance as  $\varepsilon_i$  and  $\eta$  are uncorrelated

$$\sigma^2(Y_i) = \sigma^2(\eta) + \sigma^2(\varepsilon_i)$$

Total variance      True variance      Error variance

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Reliability

$$Y = \eta + \varepsilon$$

Can only be disentangled when there are repeated measures  
 $i = \text{repeated measurement (raters, occasions etc.)}$

$$Y_i = \eta + \varepsilon_i$$

**Assumption:**  
 Rewritten to variance as  $\varepsilon_i$  and  $\eta$  are uncorrelated

$$\sigma^2(Y_i) = \sigma^2(\eta) + \sigma^2(\varepsilon_i)$$

**Assumption:**  
 $\sigma^2(\varepsilon_i)$  is constant for every repetition  $\rightarrow$   
 implies that  $\sigma^2(Y_i)$  is also constant

$$\sigma^2(Y) = \sigma^2(\eta) + \sigma^2(\varepsilon)$$

$$\text{Reliability} = \frac{\text{True variance}}{\text{Total variance}}$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Reliability

$$Y = \eta + \varepsilon$$

Can only be disentangled when there are repeated measures  
*i* = repeated measurement (raters, occasions etc.)

$$Y_i = \eta + \varepsilon_i$$

**Assumption:**  
 Rewritten to variance as  $\varepsilon_i$  and  $\eta$  are uncorrelated

$$\sigma^2(Y_i) = \sigma^2(\eta) + \sigma^2(\varepsilon_i)$$

**Assumption:**  
 $\sigma^2(\varepsilon_i)$  is constant for every repetition  $\rightarrow$   
 implies that  $\sigma^2(Y_i)$  is constant

$$\sigma^2(Y) = \sigma^2(\eta) + \sigma^2(\varepsilon)$$

$\sigma^2(\eta) = \sigma_p^2; p = \text{patient}$

$$\text{Reliability} = \frac{\text{True variance}}{\text{Total variance}} = \frac{\sigma_p^2}{\sigma_Y^2} = \frac{\sigma_p^2}{\underbrace{\sigma_p^2 + \sigma_e^2}_{\text{COSMIN taxonomy}}}$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Intraclass Correlation Coefficients

## ICC

- 6 different types (Shrout & Fleiss)
- 10 different types (McGraw & Wong)
- Variation ratio between patients and error (random error)
- ANOVA analyses

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# An example (test-retest)

10 patients measured at 2 time points

Patients	Time points	
	t <sub>1</sub>	t <sub>2</sub>
1	9	7
2	6	5
3	3	4
4	2	1
5	4	5
6	6	4
7	8	9
8	4	5
9	2	1
10	10	9
	5.4	5.0

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Sources of variation

Interaction effect (random error) =  $\sigma_{po,e}^2$

Variance between patients =  $\sigma_p^2$

$\sigma_o^2$  = Variance between time points (systematic error)

Patients	Time points	
	t <sub>1</sub>	t <sub>2</sub>
1	9	7
2	6	5
3	3	4
4	2	1
5	4	5
6	6	4
7	8	9
8	4	5
9	2	1
10	10	9
	5.4	5.0

*Note: In the original image, a red box highlights the patient numbers 1-10, and a red arrow points from the interaction effect term to the patient numbers. A green arrow points from the interaction effect term to the time points columns, and a black arrow points from the variance between time points term to the time points columns.*

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# ICC

$$ICC = \frac{\text{signal}}{\text{signal} + \text{noise}} = \frac{\text{true variance}}{\text{true variance} + \text{error variance}}$$

$$ICC = \frac{Var_{patient}}{Var_{patient} + Var_{error}} = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{error}^2}$$

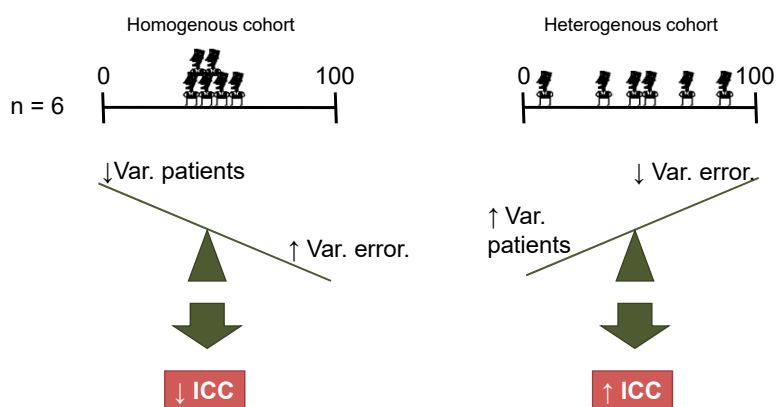
Shrout & Fleis, 1979; McGraw & Wong, 1996

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# ICC

$$ICC = \frac{Var_{patient}}{Var_{patient} + Var_{error}}$$



Shrout & Fleis, 1979; McGraw & Wong, 1996

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS





# Types of reliability

Types	Definitions
Inter-rater reliability	Variation between 2 or more raters who measure the same group of subjects
Intra-rater reliability	Variation of data measured by 1 rater across 2 or more trials
Test-retest reliability	Variation in 2 measurements taken by an instrument on the same subject under the same conditions*

\* Generally indicative of reliability in situations when raters are not involved or rater effect is neglectable, such as self-report survey instruments

# Choice of ICC formula

*Depends on:*

## I. Model

- One-way or two-way ANOVA
- Inclusion of 'Random' or 'fixed' factors

## II. Definition of relationship which is important

- Ranking ('consistency') or absolute agreement ('agreement')

## III. Type

- Single rating or means of multiple ratings

# Model

## One-way ANOVA

- When only interested in distinguishing between subjects and not in other factor(s).

## Two way ANOVA

- When also interested in other factors, e.g. subjects AND time-points (2 factors)

# Example: one- or two-way

## One-way ANOVA

- Variation in birth weight of twins

## Two-way ANOVA

- Is birth weight of *first-born* baby higher than *second-born* in twins



# Model

## Random factors

- One wants to generalise to all other representors of that factor (observers or moments)

## Fixed factors

- One is only interested in the representors of the factor under study (observers or moments)

# Example: random or fixed factors

## Random factors

- How reliable are the judgements of X-rays by an 'average' radiologist?

## Fixed factors

- How reliable are the judgements of the two radiologists who rate the X-rays in this study?

# Definition of relationship which is important

## Consistency

- Only random error is included – ranking is important

## Agreement

- Systematic errors are also included – absolute agreement is important

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Type

## Single rating

- Comparing 1 rating from each of e.g. 4 raters

## Multiple ratings

- Comparing the mean rating of e.g. 4 raters

If the mean value is taken of multiple ratings, the error variance becomes smaller as random errors diminish by averaging

Data on individual ratings more common

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# ICC forms

McGraw and Wong (1996)				Shrout and Fleiss (1979)
<i>Model</i>		<i>Relationship</i>	<i>Type</i>	
One-way	Random effects	Agreement	Single rater/measurement	ICC (1,1)

Model

Type

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# ICC forms

McGraw and Wong (1996)				Shrout and Fleiss (1979)
Model		Relationship	Type	
One-way	Random effects	Agreement	Single rating/measurement	ICC (1,1)
One-way	Random effects	Agreement	Multiple ratings/measurements	ICC (1,K)
Two-way	Random effects	Agreement	Single rating/measurement	ICC (2,1)
Two-way	Random effects	Agreement	Multiple ratings/measurements	ICC (2,K)
Two-way	Random effects	Consistency	Single rating/measurement	-
Two-way	Random effects	Consistency	Multiple ratings/measurements	-
Two-way	Mixed effects	Agreement	Single rating/measurement	-
Two-way	Mixed effects	Agreement	Multiple ratings/measurements	-
Two-way	Mixed effects	Consistency	Single rating/measurement	ICC (3,1)
Two-way	Mixed effects	Consistency	Multiple ratings/measurements	ICC (3,K)

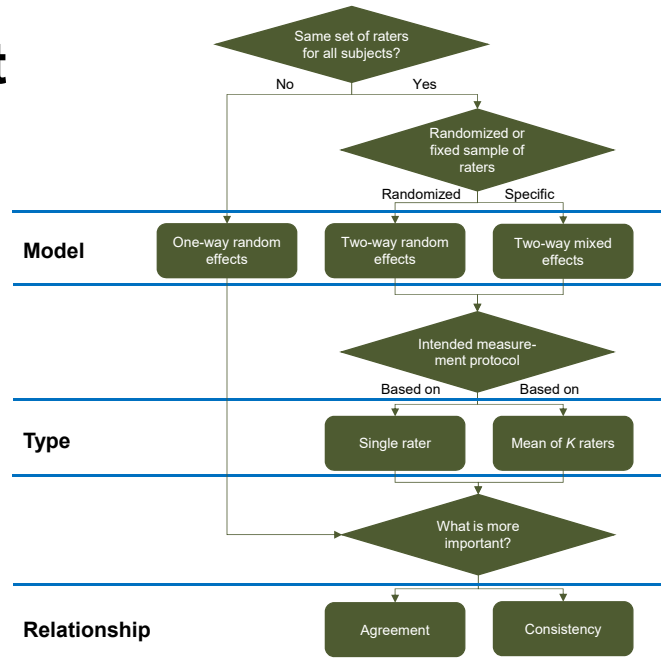
Choice depends on study design: inter-rater, intra-rater or test-retest

Koo et al., 2016

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS

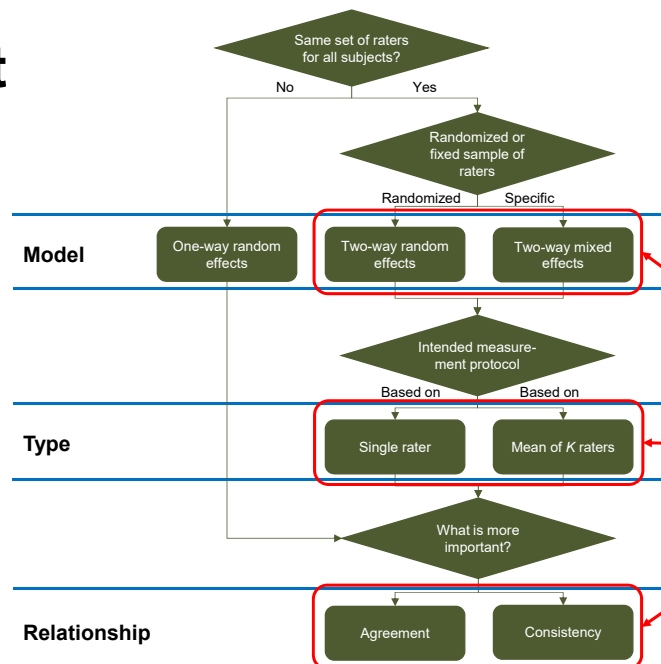


# Flowchart



# Flowchart

Inter-rater reliability



**EXAMPLE**  
Using a randomized or mixed set of raters – several choices:

- Single or mean ratings
- Ranking or absolute agreement

Common



## ICC (2,1)<sub>consistency</sub>

$$ICC(2,1)_{consistency} = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{po,e}^2}$$

↑  
Random error (residual)

- Two-way random-effects ANOVA
- Not including systematic error
- Used in analyses where the systematic error is known to be small

Shrout & Fleis 1979; McGraw & Wong 1996, Weir 2005

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



## ICC (2,1)<sub>agreement</sub>

$$ICC(2,1)_{agreement} = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_o^2 + \sigma_{po,e}^2}$$

↑  
Systematic error between t1 and t2

↑  
Random error (residual)

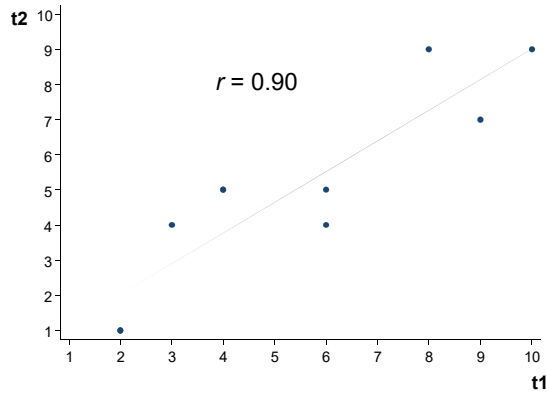
- Two-way random-effects ANOVA
- Includes systematic error
- Often used in test-retest designs

Shrout & Fleis 1979; McGraw & Wong 1996, Weir 2005

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# What about Pearson's $r$ ?



Patienter	Measurements	
	$t_1$	$t_2$
1	9	7
2	6	5
3	3	4
4	2	1
5	4	5
6	6	4
7	8	9
8	4	5
9	2	1
10	10	9
	5.4	5.0

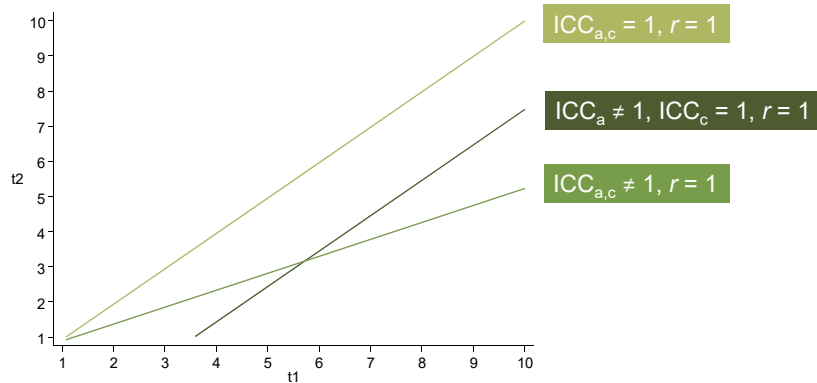
Does not consider:

- Systematic error
- Variance differences between  $t_1$  and  $t_2$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Pearson vs. ICC



DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS





# Kappa

*"A measure of reliability between two measurements which is adjusted for chance agreement."*

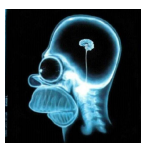
Nominal scale: kappa

Ordinal scale: weighted kappa

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



## Nominal scale - example



		Radiologist B		Total
		+	-	
Radiologist A	+	50	15	65
	-	15	20	35
Total		65	35	100

Observed agreement ( $p_o$ ):  $(50 + 20)/100 = 0.70$

Expected agreement ( $p_e$ ) = 0.54

$$\kappa = (p_o - p_e) / (1 - p_e)$$

$$\kappa = (0.70 - 0.54) / (1 - 0.54) = 0.35$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Ordinal scale - weighted kappa

E.g. Classification of malignant melanoma

	absent	mild	moderate	severe
absent				
mild				
moderate				
severe				

Some 'errors' are weighted more than others

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Interpretation

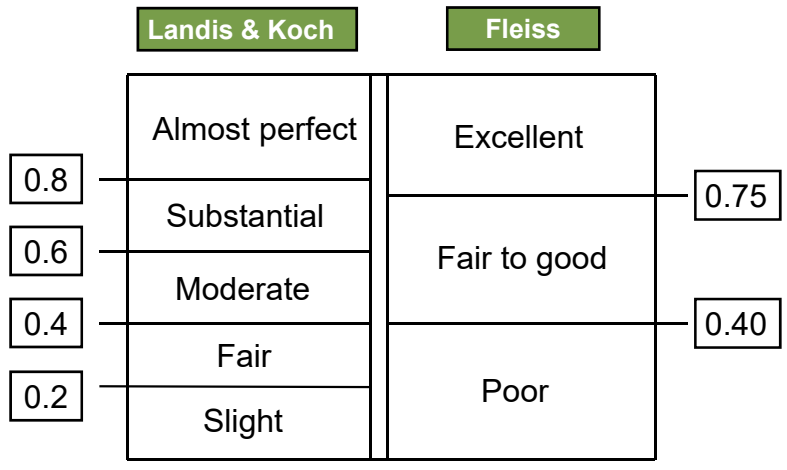
Kappa depends on:

- 1. The number of response options**
  - More response options → lower kappa
- 2. The score prevalence**
  - Skewed distributions → lower kappa
- 3. Presence of systematic differences**
  - Slightly higher kappa

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Interpretation

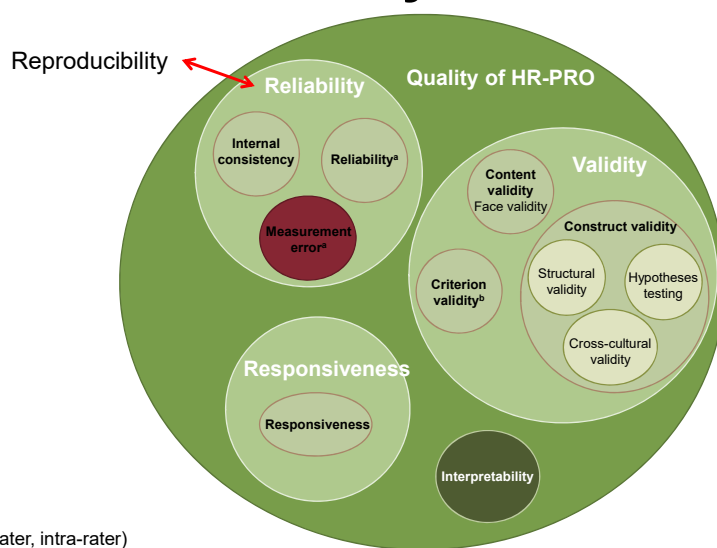


Landis and Koch, the measurement of observer agreement for categorical data, Biometrics 1977; 33: 159-174  
 Fleiss J.L. Statistical methods for rates and proportions, 2nd ed. New York, John Wiley & Sons, 1981.

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# The COSMIN taxonomy



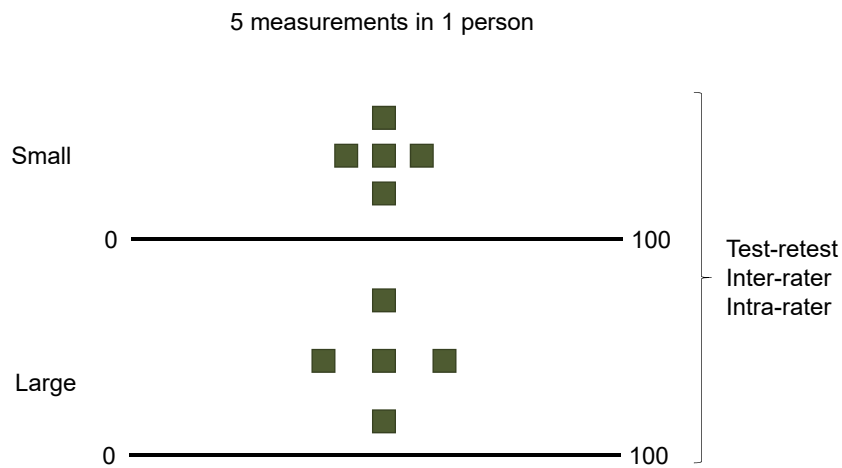
<sup>a</sup>(test-retest, inter-rater, intra-rater)

<sup>b</sup>(concurrent validity, predictive validity)

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Measurement error



DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Measurement error

Scale	Parameters
Continuous	Standard error of the measurement (SEM) Limits of agreement (LOA)
Ordinal	% agreement
Nominal	% agreement

- Same units as the measurement instrument
- Easy to interpret
- A characteristic of the questionnaire

De Vet et al. "Measurement in Medicine", 2011

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Standard error of the measurement (SEM)

## Definition

"Is the standard deviation of errors of measurement that are associated with 'test' scores"

## Purpose

- Quantify the extend to which a 'test' provides accurate scores
  - ↓ SEM → ↑ score accuracy and vice versa

## Use

- SEM used to calculate confidence intervals around obtained scores
- E.g.
  - 68% CI = Score ± (1\*SEM)
  - 95% CI = Score ± (1.96\*SEM)
  - 99% CI = Score ± (2.58\*SEM)

## Example

- Depression score = 50 (score range [0-100])
- How confident are we that the person's true depression score is 50?

De Vet et al. 2006

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# SEM

$$SEM = SD(\sqrt{1 - reliability})$$

## Example

- Depression score with SD = 15 (cross-sectional design)
- What is SEM of a single depression score = 50 under three different conditions:

- Reliability = 0.70:  $SEM = 15(\sqrt{1 - 0.7}) = 8.22$       95% CI =  $50 \pm (1.96 \times 8.22) = [66.11; 33.89]$
- Reliability = 0.80:  $SEM = 15(\sqrt{1 - 0.8}) = 6.71$       95% CI =  $50 \pm (1.96 \times 6.71) = [63.15; 36.85]$
- Reliability = 0.90:  $SEM = 15(\sqrt{1 - 0.9}) = 4.74$       95% CI =  $50 \pm (1.96 \times 4.74) = [59.29; 40.71]$

This is the reason single scores have to be more reliable → usually > 0.9  
In group scores, the measurement error ↓ by a factor  $\sqrt{n}$  → usually > 0.7

de Vet et al. Measurement in Medicine, p 142.

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# SEM

## Longitudinal design

- SEM used to calculate the measurement error of an instrument over time
- Established in a test-retest design on stable patients
- Can be extracted from the ICC when all variance components are available

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# SEM<sub>consistency</sub>

$$\text{Reliability} = \text{ICC}_{\text{consistency}} = \frac{\sigma_{\text{patient}}^2}{\sigma_{\text{patient}}^2 + \sigma_{\text{error}}^2}$$

$$\text{Standard Error of Measurement (SEM}_{\text{consistency}}) = \sqrt{\sigma_{\text{error}}^2}$$

$$\text{SEM}_{\text{consistency}} = \text{SD} \times \sqrt{(1 - \text{ICC}_{\text{consistency}})}$$

$$(\text{SD} = \text{SD}_{\text{pooled}} = (\text{SD}_{\text{time1}} + \text{SD}_{\text{time2}})/2)$$

$$\text{SEM}_{\text{consistency}} = \text{SD}_{\text{change}}/\sqrt{2}$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



## SEM<sub>agreement</sub>

$$\text{Reliability} = ICC_{\text{agreement}} = \frac{\sigma_{\text{patient}}^2}{\sigma_{\text{patient}}^2 + \sigma_{\text{measurements}}^2 + \sigma_{\text{error}}^2}$$

$$\text{Standard Error of Measurement (SEM}_{\text{agreement}}) = \sqrt{\sigma_{\text{measurements}}^2 + \sigma_{\text{error}}^2}$$

$$SEM_{\text{agreement}} = SD \times \sqrt{(1 - ICC_{\text{agreement}})} \quad (SD = SD_{\text{pooled}} = (SD_{\text{time1}} + SD_{\text{time2}})/2)$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



## Bland & Altman LOA

Graphical presentation of mean and difference between two scores

- Limits of agreement (LOA) = measurement error

$$LOA = \text{mean}_{\text{change}} \pm 1.96 \cdot SD_{\text{change}}$$

Systematic error

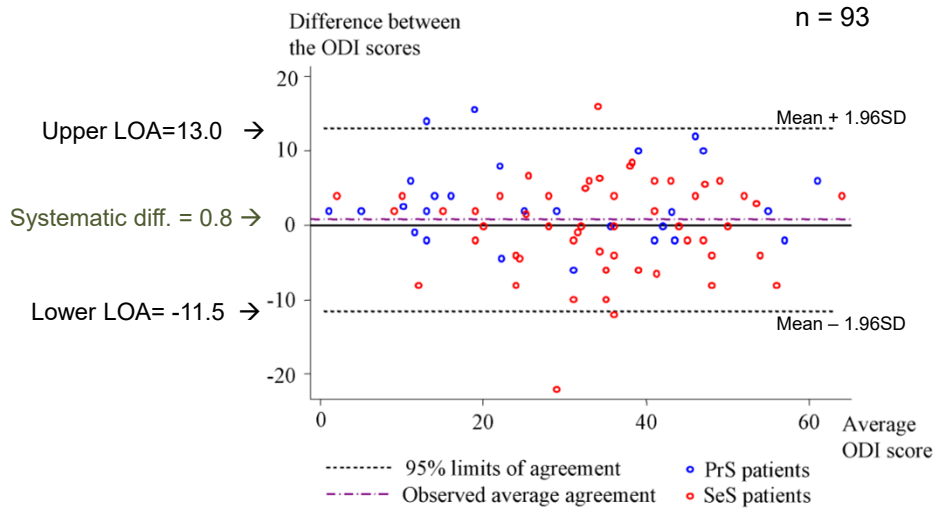
Random error

95% of the differences lies between the LOA

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Bland & Altman plot (agreement)



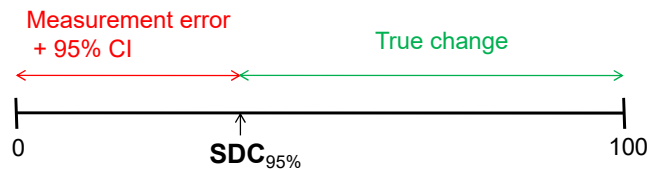
DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Smallest detectable change (SDC)

Is the smallest change in score that you **CAN** detect with the instrument, above measurement error (in individual patients)

$$SDC_{agreement} = 1.96 \cdot SEM_{agreement} \times \sqrt{2}$$



NB: Also called Minimal Detectable Change, Real change, Smallest Real Change, Significant change

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS





# SEM, LOA and SDC

## SEM

$$SEM_{consistency} = SD_{change} / \sqrt{2}$$

## Limits of agreement (LOA)

$$LOA = mean_{difference} \pm 1.96 \times SD_{change}$$

$$LOA = mean_{difference} \pm 1.96 \times SEM_{consistency} \times \sqrt{2}$$

## Smallest detectable change (SDC<sub>95%</sub>)

$$SDC_{consistency} = 1.96 \times SD_{change}$$

$$SDC_{consistency} = 1.96 \times SEM_{consistency} \times \sqrt{2}$$

$$SDC_{consistency} = 1.96 \times SD \times \sqrt{(1 - ICC_{consistency})} \times \sqrt{2} = 1.96 \times SD \sqrt{2 \times (1 - ICC_{consistency})}$$

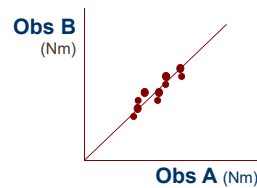
$$SDC_{agreement} = 1.96 \times SEM_{agreement} \times \sqrt{2}$$

$$SDC_{agreement} = 1.96 \times SD \sqrt{(1 - ICC_{agreement})} \times \sqrt{2} = 1.96 \times SD \sqrt{2 \times (1 - ICC_{agreement})}$$

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS

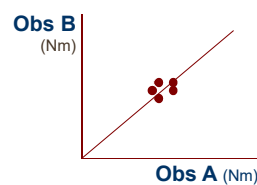


## Example: What is the reliability (Rel) and measurement error (SEM)?



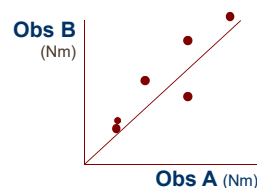
Rel = 0.986

SEM = 1.48



Rel = 0.5

SEM = 1.48



Rel = 0.993

SEM = 4.39

Stratford et al, 1989

DEPARTMENT OF SPORTS SCIENCE AND CLINICAL BIOMECHANICS



# Summary

## Reproducibility

- **Reliability**
  - Measurement error related to the variation between patients
  - Discriminative instruments
  - Parameters: ICC, Kappa
- **Measurement error**
  - Measurement error on repeated measurements
  - Evaluative instruments
  - Parameters: SEM, LOA (SDC)