

## Lesson 2: Reliability and measurement error

Patrícia Martinková

Department of Statistical Modelling  
Institute of Computer Science, Czech Academy of Sciences

Institute for Research and Development of Education  
Faculty of Education, Charles University, Prague

NMST570, October 9, 2018

# Table of contents

1. Introduction
2. Estimation Procedures
3. More on Cronbach's alpha
4. More on IRR
5. Conclusion

# Classical test theory

In behavioral research we are typically interested in the **true score**  $T$  but have available only the **observed score**  $X$  which is contaminated by some (uncorrelated) **measurement error**  $e$ , such that  $X = T + e$ .

## Examples:

- Admission tests: we are interested in **applicant's knowledge or ability**  $T$ , but have available only the test score  $X$
- Grading of essays: We are interested in **essay's quality**  $T$  but we have available only the grader's evaluation  $X$
- Questionnaires on satisfaction: main interest is **respondent satisfaction**, but available are only his/her responses on the questionnaire.

The observed score might vary if we chose different items or different graders.

# Classical test theory

Natural questions:

- How much information about the true score is indeed contained in the measurement?
- What is the strength of the relationship between true and observed score?

## Reliability

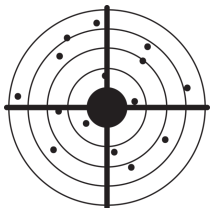
- Reliability is defined as squared correlation of the true and observed score  $\text{Rel}(X) = \rho_X = \text{cor}^2(T, X) = \rho_{T,X}^2$
- $\rho_X \in \langle 0, 1 \rangle$
- equivalently, reliability can be reexpressed as the ratio of the true score variance to total observed variance  $\rho_X = \frac{\text{var}(T)}{\text{var}(X)} = \frac{\sigma_T^2}{\sigma_X^2}$

## Implications of low reliability

- Less accurate estimates of the true score
- Wider (less precise) confidence intervals
- Need of higher number of subjects to demonstrate differences between groups (keeping the same test power)
- Attenuation of correlations, bound of criterion validity
  - Assume two traits  $T_1, T_2$  measured as  $X_1, X_2$  with uncorrelated errors  $e_1, e_2$  and reliabilities  $\text{Rel}(X_1), \text{Rel}(X_2)$
  - Observed correlation is attenuated

$$\begin{aligned}\text{cor}(X_1, X_2) &= \frac{\text{cov}(X_1, X_2)}{\sqrt{\text{var}(X_1)}\sqrt{\text{var}(X_2)}} = \frac{\text{cov}(T_1, T_2) + 0 + 0 + 0}{\sqrt{\text{var}(T_1) \frac{\text{var}(X_1)}{\text{var}(T_1)}} \sqrt{\text{var}(T_2) \frac{\text{var}(X_2)}{\text{var}(T_2)}}} \\ &= \text{cor}(T_1, T_2) \sqrt{\text{Rel}(X_1)\text{Rel}(X_2)}\end{aligned}$$

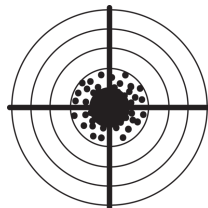
# Graphical interpretation



Low reliability thus low validity



High reliability but low validity



High reliability and high validity

- center of the target represents the value we want to measure
- shots represent independent measurements on one object
- reliability represented by variability of the shots
- validity represented by overall shots' closeness to the center

## Observations

- high reliability does not ensure high validity
- validity is bounded by reliability

# Reliability guidelines

- Conventional requirement  $\rho_x \geq .8$ , but see Lee (2012)
  - $\geq .9$  for intelligence tests
  - $\geq .7$  for personality tests
  - $\sim .6$  for essay marking
- In case of low reliability we should think of instrument revision
  - adding items
  - deleting items
  - in case of graders: training, precise instructions

## Importance of proper estimation of reliability

- Overestimation may imply adopting unreliable instrument
- Underestimation may imply (costly) revision of instrument
- Misunderstanding of reliability can imply deletion of important items and lowering validity

# Estimation procedures

The true score  $T$  is not observed, thus we can't estimate reliability from its definition ( $\rho_{T,X}^2$  nor  $\sigma_T^2/\sigma_X^2$ )

## Parallel measurements

- equally precise measurements of the same true score:
- $X_1 = T + e_1, \quad X_2 = T + e_2, \quad \text{var}(e_1) = \text{var}(e_2) = \sigma_e^2$
- the reliability of both measurements is the same  $\rho$
- if the errors are uncorrelated, then **correlation between the measurements is equal to their (common) reliability**

$$\text{cor}(X_1, X_2) = \rho_{X_1, X_2} = \frac{\text{cov}(T+e_1, T+e_2)}{\sqrt{\text{var}(T+e_1)\text{var}(T+e_2)}} = \frac{\sigma_T^2}{\sigma_X^2} = \rho$$

The methods differ in how they make use of multiple measurements.



# Estimation procedures

## Use of multiple administrations

Methods employ correlation coefficient btw. observed total scores

- Test-retest method (coefficient of stability)
- Alternate test forms (coefficient of equivalence)

## Use of composite measurements

Methods employ correlation coefficient btw. observed partial total scores

- Split-half coefficient
- Average split-half
- Cronbach's alpha (coefficient of internal consistency)

# Test-Retest

- Assumes independent test administrations
  - No memory
  - No improvement between administrations



- Some interval between administrations, say 6-12 weeks

# Parallel Forms

- Assumes trully paralel forms
  - Equally difficult
  - Parallel items and content
- Assumes the same conditions

# Composite measurements

- Goal is to provide multiple converging pieces of information
- E.g. educational tests, scales, questionnaires, ...

What is the relationship between reliability of composite measurement  $X = \sum_{j=1}^m X_j$  and reliability of its components?

## Spearman-Brown prophecy formula (1910)

Assume  $m$  parallel measurements  $X_1, \dots, X_m$  (independent, equally precise, with uncorrelated errors and uncorrelated with true scores). Then reliability of each  $X_i$  is the same  $\rho$  and the reliability of composite measurement  $X$  is

$$\rho_X = \frac{m \cdot \rho}{1 + (m - 1)\rho}$$

Remark: Adding parallel items increases reliability of total score.

# Generalized prophecy formula

## Spearman-Brown prophecy formula (generalized)

Assume test composed of  $m_1$  parallel measurements  $X = \sum_{j=1}^{m_1} X_j$  and its prolonged or shortened version composed of  $m_2$  parallel measurements  $X = \sum_{j=1}^{m_2} X_j$ . Then the relationship between their reliabilities is

$$\rho_{m_2} = \frac{\frac{m_2}{m_1} \cdot \rho_{m_1}}{1 + \left(\frac{m_2}{m_1} - 1\right)\rho_{m_1}}$$

Proof (hint): Notice that

$$\rho_1 = \frac{\frac{1}{m_1} \cdot \rho_{m_1}}{1 + \left(\frac{1}{m_1} - 1\right)\rho_{m_1}} = \frac{\frac{1}{m_2} \cdot \rho_{m_2}}{1 + \left(\frac{1}{m_2} - 1\right)\rho_{m_2}}$$

# Split-half coefficient

- Correlation between two subscores corrected for test length
- Test is split into two parts, two subscores  $Y_1, Y_2$  are computed
- $$\rho_{SH} = \frac{2\rho_{Y_1, Y_2}}{1 + \rho_{Y_1, Y_2}}$$
- Assumes that the two subtests are parallel
- Depends on how the split was carried out (even/odd, random, . . . )
  - even-numbered / odd-numbered
  - with intention to create two halves that are as similar as possible
  - in a random fashion
- We may also compute the mean of all possible split-half coefficients
  - average split-half
- We may also compute the worst of all possible split-half coefficients
  - Revelle's  $\beta$

# Cronbach's alpha

- Based on idea of splitting the test into individual items

$$\alpha = \frac{m}{m-1} \frac{\sum \sum_{j \neq k} \text{cov}(X_j, X_k)}{\text{var}(X)} = \frac{m}{m-1} \left( 1 - \frac{\sigma_{X_1}^2 + \dots + \sigma_{X_m}^2}{\sigma_X^2} \right)$$

- Popular estimator, provides simple and unique estimation
- Equals to composite reliability  $\sigma_T^2/\sigma_X^2$  in case of parallel (or at least  $T$ -equivalent) items and uncorrelated errors
- In general case and uncorrelated errors, alpha is lower bound to reliability  $\alpha \leq \rho_X$  (Novick & Lewis, 1967) and can be viewed as **index of internal consistency**
- In case of correlated errors, alpha can be lower or greater than reliability

# Cronbach's alpha: 2-way mixed ANOVA approach

- $X_{ij}$  responses of  $n$  students on  $m$  items
- $X_{ij} = T_i + b_j + e_{ij}$ 
  - $T_i \sim N(0, \sigma_T^2)$  random, student ability
  - $b_j$  fixed,  $\sum b_j = 0$ , describe item difficulty
  - $e_{ij} \sim N(0, \sigma_e^2)$  random error
  - total scores  $X_i = mT_i + \sum_j b_j + \sum_j e_{ij}$
- reliability:  $\rho_X = \frac{\text{var}(mT_i)}{\text{var}(X_i)} = \frac{m^2\sigma_T^2}{m^2\sigma_T^2 + m\sigma_e^2} = \frac{\sigma_T^2}{\sigma_T^2 + \frac{1}{m}\sigma_e^2}$
- Cronbach's alpha:
 
$$\alpha = \frac{m}{m-1} \frac{\sum \sum_{j \neq k} \text{cov}(X_{ij}, X_{ik})}{\text{var}(X_i)} = \frac{m}{m-1} \frac{m(m-1)\sigma_T^2}{m^2\sigma_T^2 + m\sigma_e^2} = \frac{\sigma_T^2}{\sigma_T^2 + \frac{1}{m}\sigma_e^2}$$
- estimate of Cronbach's alpha:  $\hat{\alpha} = \frac{m}{m-1} \frac{\sum \sum_{j \neq k} s_{jk}}{\sum \sum_{j,k} s_{jk}}$ , where  $s_{jk} = \frac{1}{n-1} \sum_{t=1}^n (X_{tj} - \bar{X}_{\bullet j})(X_{tk} - \bar{X}_{\bullet k})$

Martinková P, & Vlčková K. [Hodnocení reliability znalostních a psychologických testů.](#) (Estimation of Reliability of Educational and Psychological Measurements. In Czech.) *Informační bulletin České statistické společnosti*, 4, pp. 1-15, 2014.



# Cronbach's alpha: 2-way mixed ANOVA approach (2)

## Sums of squares

- $SS_T = \sum \sum (\bar{X}_{i\bullet} - \bar{X}_{\bullet\bullet})^2 \sim (m\sigma_T^2 + \sigma_e^2)\chi^2(n-1)$
- $SS_e = \sum \sum (X_{ij} - \bar{X}_{\bullet j} - \bar{X}_{i\bullet} + \bar{X}_{\bullet\bullet})^2 \sim \sigma_e^2\chi^2((n-1)(m-1))$

## Expectations of Mean sums of squares

- $EMS_T = ESS_T/(n-1) = m\sigma_T^2 + \sigma_e^2$
- $EMS_e = ESS_e/((n-1)(m-1)) = \sigma_e^2$

## Cronbach's alpha

$$\alpha = \frac{\sigma_T^2}{\sigma_T^2 + \frac{1}{m}\sigma_e^2} = \frac{EMS_T - EMS_e}{EMS_T}$$

## Cronbach's alpha estimate

$$\hat{\alpha} = \frac{m}{m-1} \frac{\sum \sum_{j \neq k} s_{jk}}{\sum \sum_{j,k} s_{jk}} = \frac{MS_T - MS_e}{MS_T} = 1 - \frac{1}{F}$$

## Cronbach's alpha: 2-way mixed ANOVA approach (3)

Estimate of Cronbach's alpha can be reexpressed as

$$\hat{\alpha} = \frac{MS_T - MS_E}{MS_T} = 1 - \frac{1}{F}$$

- $F$  statistic used to test the submodel with no subject effect ( $H_0 : \sigma_T^2 = 0$ )
- Interpretation: alpha close to 1 for  $F$  high, i.e. when we reject  $H_0$ , i.e. when admission test well discriminates between students
- Gives confidence intervals
- Estimate is not generally appropriate for more complicated designs

# Cronbach's alpha - limitations

Cronbach's alpha is a good estimator of reliability for

- parallel (or at least T-equivalent) items and and
- uncorrelated errors

Corrections needed for:

- Correlated errors
  - Example: Reading test, group of items associated with one text.
  - Corrections for correlated errors (Rae, 2006)
- Multidimensional measurement
  - Example: Math test, items measuring arithmetic skills but also reading skills etc.
  - Factor-analysis based estimation of reliability (Raykov & Maurcoulides, 2011)
- More sources of error (multilevel models, G-index)
- Other than normal distribution of item responses (what happens in case of binary items?)

# Logistic alpha

$F$  statistic in

$$\hat{\alpha} = 1 - \frac{1}{F}$$

assumes normality of items

- How does the estimate of reliability behave for binary items?
- Would a new estimate

$$\hat{\alpha}_{log} = 1 - \frac{n-1}{X^2}$$

based on statistic used in similar situation in logistic regression (difference of deviances  $X^2 = D(B) - D(A+B)$ ) give better results for case of binary data?

---

Martinková P, & Zvára K. Reliability in the Rasch Model. *Kybernetika*, 43(3), pp. 315-26, 2007. <http://www.kybernetika.cz/content/2007/3/315/paper.pdf>

## Definition of reliability in binary items

- Classical model not applicable (binary outcome can't be expressed as sum of  $T$  and independent error  $e$ )
- IRT models usually assumed
- Reliability can be defined as (Martinková & Zvára, 2007)

$$\rho_x = \frac{\text{var}(E(X|T))}{\text{var}(E(X|T)) + E(\text{var}(X|T))} = \frac{\text{var}(E(X|T))}{\text{var}(X)}$$

- Resulting integrals can be evaluated numerically, not explicitly
- Not equal to parallel-forms reliability, but differences negligible (Kim, 2012)
- S-B formula holds only approximately (Martinková, Zvara 2010)

---

Martinková P, & Zvára K. Reliability in the Rasch Model. *Kybernetika*, 43(3), pp. 315-26, 2007. <http://www.kybernetika.cz/content/2007/3/315/paper.pdf>

## Cronbach's alpha in binary items

- Cronbach's alpha is readily applicable also for binary items
- Cronbach's alpha represents generalization of so-called Kuder-Richardson formulae (*Psychometrika*, 1937):

- $\hat{\rho}_{KR-20} = \frac{p}{p-1} \left[ 1 - \frac{\sum \hat{r}_k(1-\hat{r}_k)}{\hat{\sigma}_X} \right]$ , where  $\hat{r}_k$  is easiness of  $k$ -th item

- For test with items of common difficulties

$$\hat{\rho}_{KR-21} = \frac{p}{p-1} \left[ 1 - \frac{\hat{\mu}(p-\hat{\mu}_k)}{p\hat{\sigma}_X} \right], \text{ where } \hat{\mu} \text{ is average total score}$$

# Logistic alpha: Simulation study

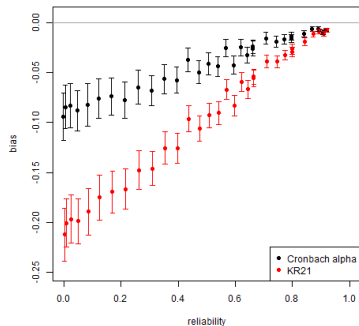
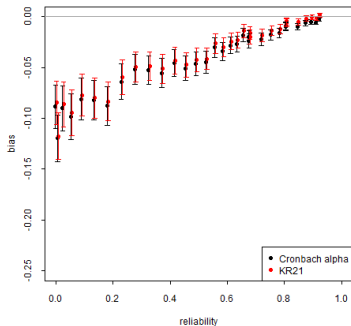
Pre-defined values:

- number of students  $n = 25, 50, 100, 500$
- number of items  $m = 10, 20, 50, 100$
- IRT parameters (difficulty, discrimination, guessing for each item)
- 55 values of  $\sigma_T$  (defines true reliability)
- number of simulates  $N = 1000$

For each combination of  $n$ ,  $m$  and  $\sigma_T$ :

- true reliability computed
  - $N$  data sets generated:
    - set of  $n$  student abilities generated  $T_i \sim N(0, \sigma_T^2)$
    - $Y_{ij}$  generated from IRT model
    - estimates computed from the data
- ⇒  $N$  estimates  $\hat{\alpha}_{CR}$ , KR-21 and  $\hat{\alpha}_{log}$
- bias and MSE of the estimates plotted out

# Simulations: Cronbach's alpha (KR-20) and KR-21



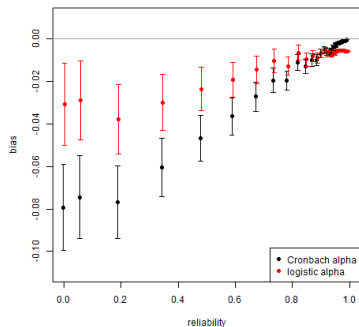
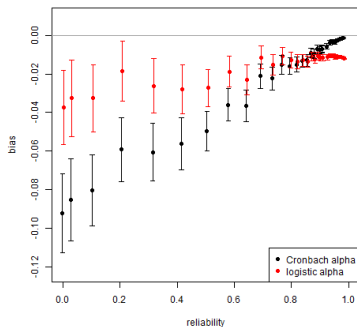
Bias and MSE of two estimators of reliability, item difficulties from  $(-0.1, 0.1)$ . Number of students  $n = 25$ , number of items  $m = 10$ , number of simulates  $N = 1000$ .

Bias and MSE of two estimators of reliability, item difficulties from  $(-3, 3)$ . Number of students  $n = 25$ , number of items  $m = 10$ , number of simulates  $N = 1000$ .

- $\hat{\alpha}_{KR-21}$  is not appropriate in case of different item difficulties



# Simulations: Cronbach's and logistic alpha



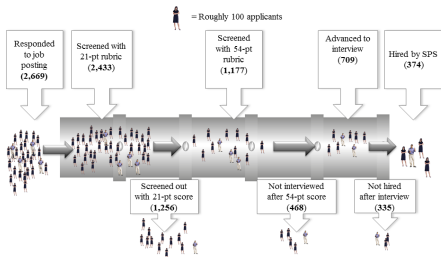
Bias and MSE of two estimators of reliability, number of students  $n = 25$ , number of items  $m = 50$ , number of simulates  $N = 1000$ .

Bias and MSE of two estimators of reliability, number of students  $n = 25$ , number of items  $m = 100$ , number of simulates  $N = 1000$ .

- $\hat{\alpha}_{log}$  has promising properties especially for high number of items

# More on inter-rater reliability

## Motivation: Teacher Selection Process



Applicants to classroom job openings in Spokane Public Schools during years (2008/09 - 2012/13)

Martinková P, Goldhaber D, & Erosheva E. (2018). Disparities in ratings of internal and external applicants: A case for model-based inter-rater reliability. *PLOS ONE*, 13(10): e0203002. <https://doi.org/10.1371/journal.pone.0203002>

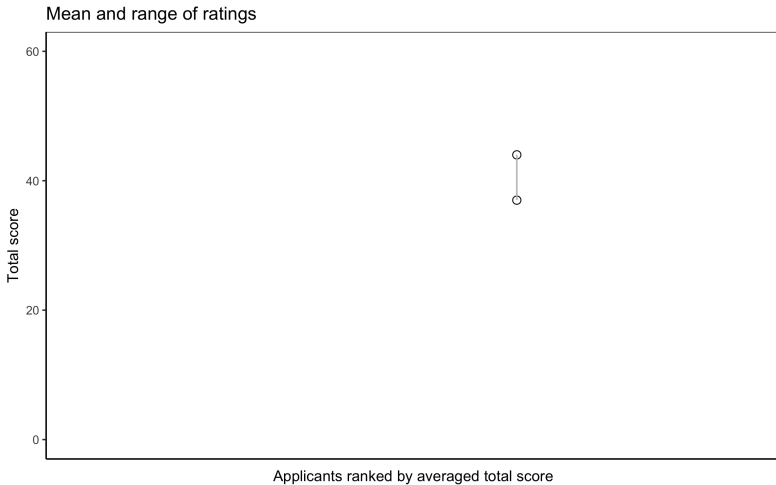
# Motivation: Ratings as Source of Error

## 54-Pt Screening Rubric:

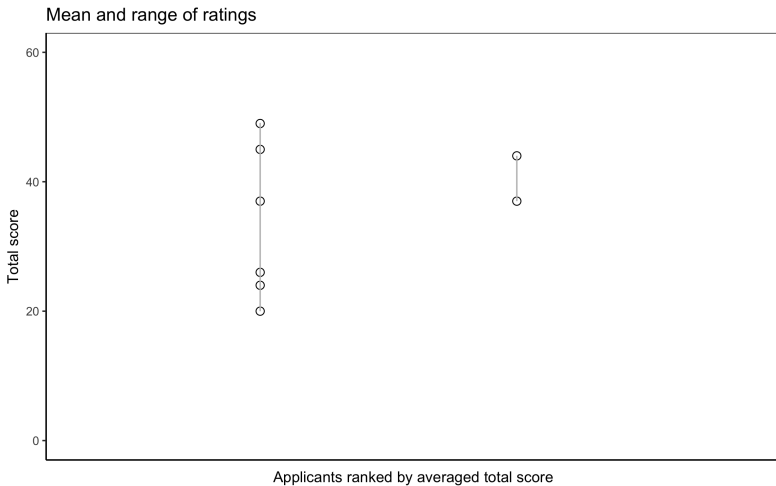
CERTIFICATED APPLICANT - PRINCIPAL / SUPERVISOR SCREENING	
DATE:	SCREENER:
Job # / Position Title:	
APPLICANT NAME:	
RATING: (3-6)	
SCREENING CRITERIA	<p>3 - 6 Strong evidence to support this as an area of strength</p> <p>3 - 4 Satisfactory evidence to support this as an area of strength</p> <p>2 - 3 Some evidence to support this as an area of strength</p>
CERTIFICATE AND EDUCATION	How completion of course of study, certificate level (baccalaureate or postgraduate) education
Washington State Certificate	Yes/No
Required Endorsement	Yes/No
RATING (1 - 6)	4
TRAINING	Level of specialty, depth and level of credentials additional training relating to the position
RATING (1 - 6)	4
EXPERIENCE	How applicant's prior experience supports the position of success - not just duration of years - if ongoing conditions could be cited briefly
RATING (1 - 6)	4
CLASSROOM MANAGEMENT	How the applicant "reflects or models" strategies - This field is not filled and checked if the applicant and checked if the applicant handles large, small or ethnically/racially diverse, at-risk groups, develops students and procedures to promote learning, establishes clear parameters, and responds appropriately.
RATING (1 - 6)	4
FLEXIBILITY	How the applicant demonstrates clearly teaching strategies, flexible teaching or ability or willingness adjust - ability to give and accept constructive feedback, successfully teaches a variety of learners, effectively uses various teaching styles
RATING (1 - 6)	4
INSTRUCTIONAL SKILLS	How the applicant reflects or models of skills in the areas - lesson, objectives, resources, teacher's impact, creating purpose, active, active, students and activities, when relevantly organize strategies appropriate to age, background and individual learning of students
RATING (1 - 6)	4
INTERPERSONAL SKILLS	Overviews and maintains effective working relationships with diverse staff, students, parents/guardians, and community
RATING (1 - 6)	4
CULTURAL COMPETENCY	Look for specific references to successful strategies for building and maintaining a relationship with each student and their family. This may not be explicitly mentioned, but the following strategies offer some evidence of cultural competency: specific instructional strategies providing a safe and student access to a rigorous curriculum; culturally-specific language about students and families; a belief that all children can achieve at high levels; explicit and implicit inclusion practices; specific instructional strategies for engaging culturally responsive students who are also rigorous and appropriate standards about their work with diverse populations. These references should be specific and address basic level
RATING (1 - 6)	4
PREFERRED QUALIFICATIONS AS INDICATED ON POSTING	
RATING (1 - 6)	4
LETTERS OF RECOMMENDATION	Look for specific letters of recommendation from most the area where appropriate. Four letters should reflect the quality and nature of the recommendation as well as the nature of the letter. (Example: Are the letters from parent or former supervisor?)
RATING (1 - 6)	4
TOTAL SCREENING SCORE	40

DO NOT REUSE SCREENING FORM 153

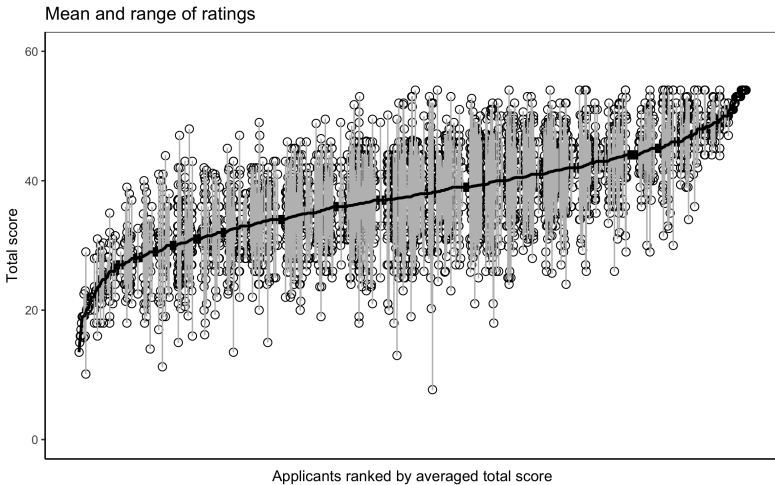
# Ratings of a single applicant



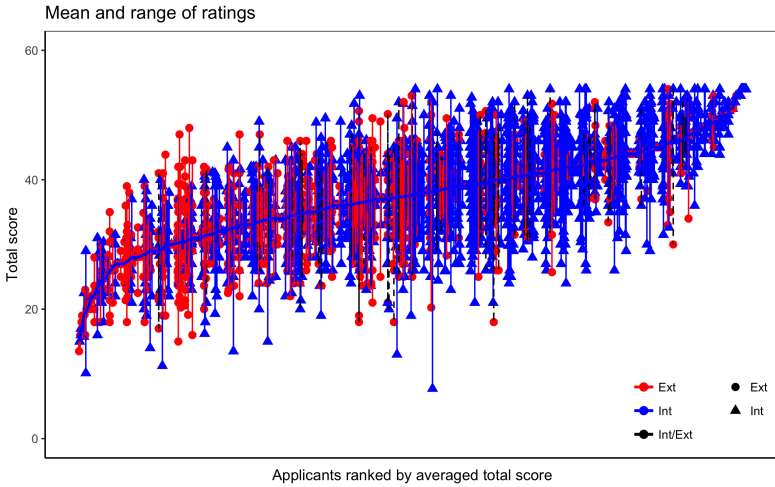
# Ratings of two applicants



# Ratings of all applicants



# Ratings of all applicants by Internal/External Status



# Inter-Rater Reliability

$$Y_{ij} = \mu + A_i + B_j + e_{ij}$$

- applicant true quality  $A_i \sim N(0, \sigma_A^2)$ ,
- rater leniency  $B_j \sim N(0, \sigma_B^2)$ ,
- error  $e_{ij} \sim N(0, \sigma_e^2)$

## Inter-Rater Reliability:

$$R = \text{cor}(Y_{ij}, Y_{ij'}) = \text{ICC} = \frac{\sigma_A^2}{\sigma_Y^2} = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2 + \sigma_e^2}$$

- $R \in [0, 1]$ , low values mean a lot of measurement error
- Aggregates (average of  $J$  raters) have higher IRR:

$$R_n = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2/J + \sigma_e^2/J}$$



# Across- and Within-School IRR (Model 1)

$$Y_{ijk} = \mu + A_i + B_j + S_k + AS_{ik} + e_{ijk}$$

- School leniency  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$
- IRR across schools:

$$R_{across} = \text{cor}(Y_{ijk}, Y_{ij'k'}) = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2 + \sigma_S^2 + \sigma_{AS}^2 + \sigma_e^2}$$

Within-school IRR:

$$R_{within} = \text{cor}(Y_{ijk}, Y_{ij'k'}) = \frac{\sigma_A^2 + \sigma_S^2 + \sigma_{AS}^2}{\sigma_A^2 + \sigma_B^2 + \sigma_S^2 + \sigma_{AS}^2 + \sigma_e^2}$$

# IRR for Internal vs. External Applicants (Model 3)

- Q: Does IRR differ in ratings of internal vs. external applicants?
- **Model 3:** Variance components may vary by group
  - e.g. Rater variance may higher when rating external applicants

$$\begin{aligned} Y_{ijk} = & \mu + \omega_i \beta_0 + (1 - \omega_i) A_{0i} + \omega_i A_{1i} \\ & + (1 - \omega_i) B_{0j} + \omega_i B_{1j} \\ & + (1 - \omega_i) S_{0k} + \omega_i S_{1k} \\ & + A S_{ik} + e_{ijk} \end{aligned}$$

- $\omega_i = 1$  for internal and 0 for external applicants
- $A_{0i} \sim N(0, \sigma_{A0}^2)$  and  $A_{1i} \sim N(0, \sigma_{A1}^2)$
- $B_{0j} \sim N(0, \sigma_{B0}^2)$  and  $B_{1j} \sim N(0, \sigma_{B1}^2)$
- $S_{0k} \sim N(0, \sigma_{S0}^2)$  and  $S_{1k} \sim N(0, \sigma_{S1}^2)$

# IRR for Internal vs. External Applicants (Model 3)

## Within-school IRR:

- For internal applicant :

$$R_1 = \text{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A1}^2 + \sigma_{S1}^2 + \sigma_{AS}^2}{\sigma_{A1}^2 + \sigma_{B1}^2 + \sigma_{S1}^2 + \sigma_{AS}^2 + \sigma_e^2}$$

- For external applicant:

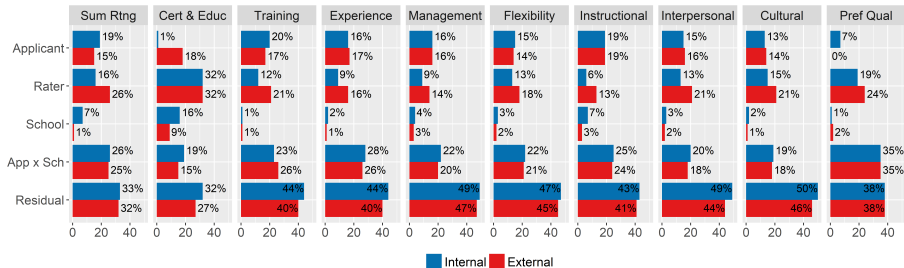
$$R_0 = \text{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A0}^2 + \sigma_{S0}^2 + \sigma_{AS}^2}{\sigma_{A0}^2 + \sigma_{B0}^2 + \sigma_{S0}^2 + \sigma_{AS}^2 + \sigma_e^2}$$

# IRR estimation and inference

## More flexible estimation using linear random-effect models

- Estimation using restricted maximum likelihood
  - `lmer()` in `lme4` in R
- Model selection using AIC, BIC, likelihood ratio tests
  - Model 3 wins for total score as well as for all subcomponents
- Bootstrapped confidence intervals or MCMC
  - using `bootMer()` in `lme4`, or `brm()` in `inbrms`

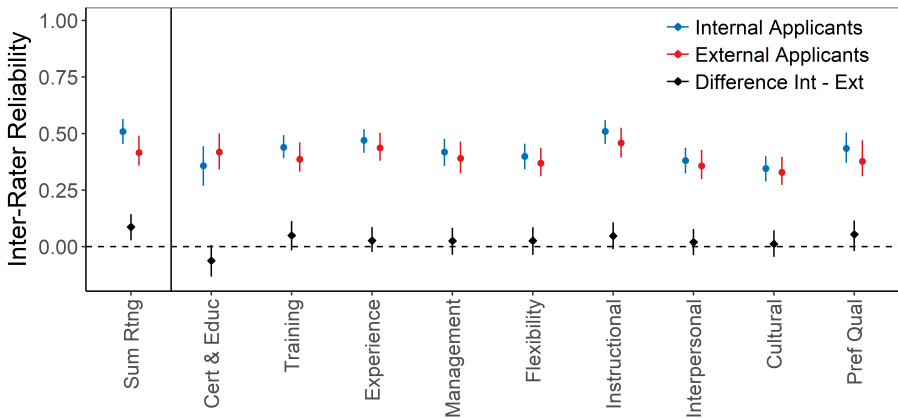
# Results: Variance decomposition (Model 3)



- High applicant-school variability
- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants

# IRR for Internal and External Applicants (Model 3)

- Significant difference in IRR between Internal and External applicants



# Conclusion

In this presentation, we have

- explained motivation behind reliability
- presented mostly used approaches for reliability estimation
  - test-retest
  - parallel forms
  - split-half coefficient
  - Cronbach's alpha
- presented research on alternative to Cronbach's alpha
- discussed use of model-based reliability estimates (for IRR)

Thank you for your attention!

[www.cs.cas.cz/martinkova](http://www.cs.cas.cz/martinkova)



- Webb NM, Shavelson RJ, & Haertel EH. Reliability coefficients and generalizability theory. In *Handbook of Statistics, Vol. 26 – Psychometrics*. pp. 81–124, Elsevier, 2007.
- Martinková P, & Vlčková K. Hodnocení reliability znalostních a psychologických testů. (Estimation of Reliability of Educational and Psychological Measurements. In Czech.) *Informační bulletin České statistické společnosti*, 4, pp. 1-15, 2014.  
[http://www.statspol.cz/cs/wp-content/uploads/IB\\_4\\_2014.pdf](http://www.statspol.cz/cs/wp-content/uploads/IB_4_2014.pdf)
- Martinková P, & Zvára K. Reliability in the Rasch Model. *Kybernetika*, 43(3), pp. 315-26, 2007.  
<http://www.kybernetika.cz/content/2007/3/315/paper.pdf>
- Martinková P, Goldhaber D, & Erosheva E. (2018). Disparities in ratings of internal and external applicants: A case for model-based inter-rater reliability. *PLOS ONE*, 13(10): e0203002.  
<https://doi.org/10.1371/journal.pone.0203002>

# Vocabulary

- Latent variable
- Reliability, measurement error
- Test-retest reliability
- Split-half
- Cronbach's alpha
- Kuder-Richardson formula
- Inter-rater reliability