# Estimating Interrater Reliability from Planned-Missing Data
## Demonstration of the ICC4IRR application

Debby ten Hove

Department of Education and Family Studies, Vrije Universiteit Amsterdam
*D.ten.Hove@VU.nl*

August 26, 2025

**Contributors**



Dr. Letty Koopman
University of Groningen



Tasos Psychogiopoulous, MSc.
University of Amsterdam

**ICC4IRR Application**



https://tasospsy.shinyapps.io/icc4irr_app/

**Example Data**



https://github.com/icc4irr/app/blob/main/sample-data/Example_ratings_EARLI2025.csv

# Motivating example: Teaching Quality

Education researchers use **observation instruments** to evaluate teaching quality



Example: **Educational inspectors** *(raters)* evaluate **teaching skills** *(attribute)* of **teachers** *(subjects)*

# Motivating example: Teaching Quality

Such ratings are used in **practice** and in **research**

# Motivating example: Teaching Quality

Such ratings are used in **practice** and in **research**

**In practice** to make decision
about teachers or schools.

# Motivating example: Teaching Quality

Such ratings are used in **practice** and in **research**

**In practice** to make decision about teachers or schools.

**In research** to study differences across teachers or schools.
For example in intervention studies.

# Motivating example: Teaching Quality

**Important** that (the variation in) observed scores originate in differences across teachers or schools, and as little as possible in rater effects.

| Use of scores | Problems due to rater effects (noise) |
|---|---|
| Regression techniques | **(Attenuation) Bias and loss of precision** |
| Decisions about individuals | **Incorrect decisions** |

$\rightarrow$ Important to investigate the **interrater reliability** (IRR).

# Intraclass correlation coefficients

## Intraclass correlation coefficients (ICCs) for Interrater reliability (IRR)

To which degree can we differentiate between subjects, hence generalize subject scores over raters?

(Bartko, 1966; McGraw & Wong, 1996; Shrout & Fleiss, 1979)

- Applicable to quantitative data
- Rooted in Generalizability theory (Cronbach et al., 1963)
- Coefficients for absolute and relative decision making
- Available for $\geq 2$ raters

# Intraclass correlation coefficients

## Forming Inferences About Some Intraclass Correlation Coefficients
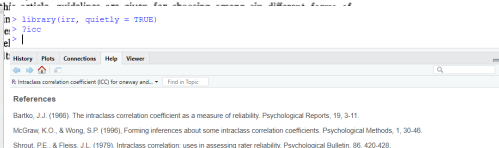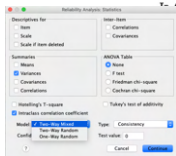
Kenneth O. McGraw
University of Mississippi

S. P. Wong
University of Memphis

Although intraclass correlation coefficients (ICCs) are commonly used in behavioral measurement, psychometrics, and behavioral genetics, procedures

...tions: Uses in Assessing
Rater Reliability

Patrick E. Shrout and Joseph L. Fleiss
Division of Biostatistics
Columbia University, School of Public Health

Reliability coefficients often take the form of intraclass correlation coefficients. In this article, guidelines are given for choosing among six different forms of the intraclass correlation for reliability studies in which $n$ targets are rated by $k$ judges. Relevant to the choice of the coefficient are the appropriate statistical model for the reliability study and the applications to be made of the reliability results. Confidence intervals for each of the forms are reviewed.

# Intraclass correlation coefficients

# ICCs for Interrater Reliability

## Two-Way Data

- Each subject is assessed by the same $k \geq 2$ raters.

| Subject | Rater | | |
|---------|-------|---|---|
| | 1 | 2 | 3 |
| 1 | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2 | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3 | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4 | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5 | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6 | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7 | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8 | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9 | $y_{91}$ | $y_{92}$ | $y_{93}$ |

# ICCs for Interrater Reliability

## Two-Way Data

- Each subject is assessed by the same $k \geq 2$ raters.

| Subject | Rater | | |
|---------|-------|-------|-------|
|         | 1     | 2     | 3     |
| 1       | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2       | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3       | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4       | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5       | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6       | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7       | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8       | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9       | $y_{91}$ | $y_{92}$ | $y_{93}$ |

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr.e}^2$$

# ICCs for Interrater Reliability

## Two-Way Data

- Each subject is assessed by the same $k \geq 2$ raters.

| Subject | Rater | | |
|---|---|---|---|
| | 1 | 2 | 3 |
| 1 | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2 | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3 | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4 | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5 | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6 | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7 | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8 | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9 | $y_{91}$ | $y_{92}$ | $y_{93}$ |

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr.e}^2$$

## Interrater Agreement

$$ICC(A, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr.e}^2}{k}}$$

# ICCs for Interrater Reliability

## Two-Way Data

- Each subject is assessed by the same $k \geq 2$ raters.

| Subject | Rater | | |
|---------|-------|-------|-------|
| | 1 | 2 | 3 |
| 1 | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2 | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3 | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4 | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5 | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6 | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7 | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8 | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9 | $y_{91}$ | $y_{92}$ | $y_{93}$ |

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr.e}^2$$

## Interrater Agreement

$$ICC(A, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr.e}^2}{k}}$$

## Interrater Consistency

$$ICC(C, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_{sr.e}^2}{k}}$$

# Practice: Planned-Missing Data

## Complete data

- ICC definitions *and* estimation methods require complete data

| Subject | Rater | | |
|---------|-------|-------|-------|
|         | 1     | 2     | 3     |
| 1       | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2       | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3       | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4       | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5       | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6       | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7       | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8       | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9       | $y_{91}$ | $y_{92}$ | $y_{93}$ |

# Practice: Planned-Missing Data

## Complete data

- ICC definitions *and* estimation methods require complete data

| Subject | Rater | | |
|---------|-------|-------|-------|
| | 1 | 2 | 3 |
| 1 | $y_{11}$ | $y_{12}$ | $y_{13}$ |
| 2 | $y_{21}$ | $y_{22}$ | $y_{23}$ |
| 3 | $y_{31}$ | $y_{32}$ | $y_{33}$ |
| 4 | $y_{41}$ | $y_{42}$ | $y_{43}$ |
| 5 | $y_{51}$ | $y_{52}$ | $y_{53}$ |
| 6 | $y_{61}$ | $y_{62}$ | $y_{63}$ |
| 7 | $y_{71}$ | $y_{72}$ | $y_{73}$ |
| 8 | $y_{81}$ | $y_{82}$ | $y_{83}$ |
| 9 | $y_{91}$ | $y_{92}$ | $y_{93}$ |

## Incomplete data

- Most educational studies use a planned-missing design. For example:

| Teacher | Educational Inspector | | |
|---------|-------|-------|-------|
| | 1 | 2 | 3 |
| 1 | $y_{11}$ | $y_{12}$ | − |
| 2 | $y_{21}$ | − | $y_{23}$ |
| 3 | − | $y_{32}$ | $y_{33}$ |
| 4 | $y_{41}$ | $y_{42}$ | − |
| 5 | $y_{51}$ | − | $y_{53}$ |
| 6 | − | $y_{62}$ | $y_{63}$ |
| 7 | $y_{71}$ | $y_{72}$ | − |
| 8 | $y_{81}$ | − | $y_{83}$ |
| 9 | − | $y_{92}$ | $y_{93}$ |

Ten Hove, Jorgensen and Van der Ark (2024). Updated Guidelines on Selecting ICCs for IRR.

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr}^2$$

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr}^2$$

### Interrater Agreement

$$ICC(A, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr}^2}{k}}$$

Account for unbalanced number of raters:

$$ICC(A, \hat{k}) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr}^2}{\hat{k}}}$$

# ICCs for planned-missing data

## Variance Decomposition

$$\sigma_y^2 = \sigma_s^2 + \sigma_r^2 + \sigma_{sr}^2$$

### Interrater Agreement

$$ICC(A, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr}^2}{k}}$$

Account for unbalanced number of raters:

$$ICC(A, \hat{k}) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_r^2 + \sigma_{sr}^2}{\hat{k}}}$$

### Interrater Consistency

$$ICC(C, k) = \frac{\sigma_s^2}{\sigma_s^2 + \frac{\sigma_{sr}^2}{k}}$$

Account for partial non-overlapping raters:

$$ICC(Q, \hat{k}) = \frac{\sigma_s^2}{\sigma_s^2 + q * \sigma_r^2 + \frac{\sigma_{sr}^2}{\hat{k}}}$$

$$\hat{k} = \left( \frac{k_1^{-1} + k_2^{-1} + \ldots + k_I^{-1}}{N_s} \right)^{-1}; q = \frac{1}{\hat{k}} - \frac{\sum_i \sum_{i'} \frac{k_{i,i'}}{k_i k_{i'}}}{N_s(N_s - 1)}$$

# Overview ICC Selection

See: Ten Hove, Jorgensen and Van der Ark (2024). Updated Guidelines on Selecting an ICC for IRR.

# Estimating ICCs from Planned-Missing Data

Compared three estimation methods for ICCs:

- **MCMC**: Markov chain Monte Carlo Estimation of hierarchical models (LoPilato et al., 2015; Ten Hove et al., 2020, 2021)
- **MLE-R**: Maximum likelihood estimation of **R**andom effects models (Marcoulides, 1990; cf. Jiang, 2018; Ten Hove et al., 2021)
- **MLE-CF**: Maximum likelihood estimation of **C**ommon-**F**actor models (Jorgensen, 2021; Marcoulides, 1996; Vispoel et al., 2018a, 2019))

For various design factors (e.g., $K, N, \hat{k}, \sigma^2$)

Based on (*among other things*):

- Computational accuracy
- 95% (B)CI coverage rates

**SCAN ME**

# Estimating ICCs from Planned-Missing Data

## AND what if..

A different observational design were used?

## Examples of Planned Missing Designs

| (1) **Random** Rater Assignment | | | | | (2) **Anker** Rater | | | | | (3) **Blocks** of Raters | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subject | Rater | | | | Subject | Rater | | | | Subject | Rater | | | |
| | 1 | 2 | 3 | 4 | | 1 | 2 | 3 | 4 | | 1 | 2 | 3 | 4 |
| 1 | x | | x | | 1 | x | x | | | 1 | x | x | | |
| 2 | | x | x | | 2 | x | | x | | 2 | x | x | | |
| 3 | | | x | x | 3 | x | | | x | 3 | x | x | | |
| 4 | x | x | | | 4 | x | x | | | 4 | | | x | x |
| 5 | x | | | x | 5 | x | | x | | 5 | | | x | x |
| 6 | | x | x | | 6 | x | | | x | 6 | | | x | x |

- **Simulation 1: MLE of random effects models** useful for estimating ICCs for IRR from typical observational studies
  - *Also* most **User-friendly** because it <u>converges</u> in most conditions and <u>only takes seconds</u>.
- **Simulation 2:** Type of (planned missing) design does not matter much with respect to (SE) bias and coverage of ICCs

Conclusion: **MLE of random effects models** very useful for (interrater)reliability studies.

# Current Work: ICC4IRR application

- **Under Development**: Shiny app to estimate ICCs for IRR

  Psychogiopoulos, Koopman & Ten Hove (2025)

- **In progress**: Tutorial and guidance in planning rater studies

SCAN ME

# ICC4IRR application



**ICC4IRR**  ESTIMATE IRR  FLOWCHART  COMPUTE DESIGN FACTORS  ESTIMATE DESIGN FACTORS  **ABOUT**

## ABOUT THIS APP

ICC4IRR is a shiny application to estimate interrater reliability (IRR) from quantitative planned incomplete data, resulting from observation studies in which raters (partly) vary across subjects.

AUTHORS:
Tasos Psychogyiopoulos, Letty Koopman and Debby ten Hove

CONTACT:
d.ten.hove@vu.nl

CITE AS:
Psychogyiopoulos, A., Koopman L. & Ten Hove, D. (2025). *ICC4IRR: A shiny application to estimate interrater reliability using intraclass correlation coefficients* . https://tasospsy.shinyapps.io/icc4irr_app/

Example citation: We investigated the interrater consistency [or agreement] using intraclass correlation coefficients (ICCs) [that accounted for partially non-overlapping raters across subjects] using the R/shiny application ICC4IRR (Psychogyiopoulos, Koopman & Ten Hove, 2025).

## ADDITIONAL REFERENCES

- Ten Hove, D., Jorgensen, T. D., & van der Ark, L. A. (2024). Updated guidelines on selecting an intraclass correlation coefficient for interrater reliability, with applications to incomplete observational designs. *Psychological Methods, 29* (5), 967–979. https://doi.org/10.1037/met0000516
- Ten Hove, D., Jorgensen, T. D., & Van der Ark, L. A. (2025). How to estimate intraclass correlation coefficients for interrater reliability from planned incomplete data. *Multivariate Behavioral Research.* https://doi.org/10.1080/00273171.2025.2507745
- R Project
- Shiny

# ICC4IRR application

## ICC4IRR application

# ICC4IRR application

# ICC4IRR application

# ICC4IRR application

# ICC4IRR application

# ICC4IRR application

## ABOUT THIS APP

ICC4IRR is a shiny application to estimate interrater reliability (IRR) from quantitative planned incomplete data, resulting from observation studies in which raters (partly) vary across subjects.

AUTHORS:

Tasos Psychogyiopoulos, Letty Koopman and Debby ten Hove

CONTACT:

d.ten.hove@vu.nl

CITE AS:

Psychogyiopoulos, A., Koopman L. & Ten Hove, D. (2025). *ICC4IRR: A shiny application to estimate interrater reliability using intraclass correlation coefficients* . https://tasospsy.shinyapps.io/icc4irr_app/

Example citation: We investigated the interrater consistency [or agreement] using intraclass correlation coefficients (ICCs) [that accounted for partially non-overlapping raters across subjects] using the R/shiny application ICC4IRR (Psychogyiopoulos, Koopman & Ten Hove, 2025).

## ADDITIONAL REFERENCES

- Ten Hove, D., Jorgensen, T. D., & van der Ark, L. A. (2024). Updated guidelines on selecting an intraclass correlation coefficient for interrater reliability, with applications to incomplete observational designs. *Psychological Methods, 29* (5), 967–979. https://doi.org/10.1037/met0000516
- Ten Hove, D., Jorgensen, T. D., & Van der Ark, L. A. (2025). How to estimate intraclass correlation coefficients for interrater reliability from planned incomplete data. *Multivariate Behavioral Research.* https://doi.org/10.1080/00273171.2025.2507745
- R Project
- Shiny

# Thanks for your attention!

## Questions or suggestions?

D.ten.Hove@VU.nl

All software code on **GITHUB**:



ICC4IRR Shiny app:

# Key References

- McGraw, K. O., & Wong, S. P. (1996). Forming inferences about some intraclass correlation coefficients. *Psychological methods*. https://psycnet.apa.org/doi/10.1037/1082-989X.1.1.30

- Psychogiopoulos, T., Koopman, L., & Ten Hove, D. (2025). *ICC4IRR: A shiny application to estimate interrater reliability using intraclass correlation coefficients*. https://tasospsy.shinyapps.io/icc4irr_app

- Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: uses in assessing rater reliability. *Psychological bulletin*. https://psycnet.apa.org/doi/10.1037/0033-2909.86.2.420

- Ten Hove, D., Jorgensen, T. D., & Van der Ark, L. A. (2024). Updated guidelines on selecting an intraclass correlation coefficient for interrater reliability, with applications to incomplete observational designs. *Psychological Methods*. https://doi.org/10.1037/met0000516

- Ten Hove, D., Jorgensen, T. D., & Van der Ark, L. A. (2025). How to Estimate Intraclass Correlation Coefficients for Interrater Reliability from Planned Incomplete Data. *Multivariate Behavioral Research*. https://doi.org/10.1080/00273171.2025.2507745