Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

# Generalized Measurement Invariance Tests for Factor Analysis

Ed Merkle[1]    Achim Zeileis[2]

[1]University of Missouri

[2]Universität Innsbruck

---

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

---

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Measurement Invariance

- Measurement invariance: Sets of tests/items consistently assigning scores across diverse groups of individuals.

- Notable violations of measurement invariance:
  - SAT for different ethnic groups (Atkinson, 2001)
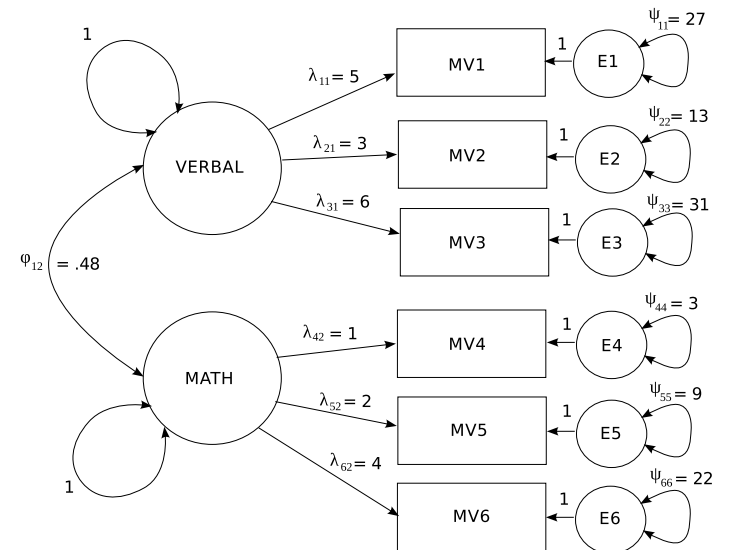  - Intelligence tests & the Flynn effect (Wicherts et al., 2004)

---

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Example (Age $\leq 16$)

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Example (Age > 16)



Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Hypotheses

- Hypothesis of "full" measurement invariance:

$$H_0 : \boldsymbol{\theta}_i = \boldsymbol{\theta}_0, i = 1, \ldots, n$$
$$H_1 : \text{Not all the } \boldsymbol{\theta}_i = \boldsymbol{\theta}_0$$

where $\boldsymbol{\theta}_i = (\lambda_{i,1,1}, \ldots, \psi_{i,1,1}, \ldots, \varphi_{i,1,2})^\top$ is the full $p$-dimensional parameter vector for individual $i$.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Hypotheses

- $H_0$ from the previous slide is difficult to fully assess due to all the ways by which individuals may differ.

- We typically place people into groups based on a meaningful auxiliary variable, then study measurement invariance across those groups (via Likelihood Ratio tests, Lagrange multiplier tests, Wald tests).

- If we did not know the groups in advance, we could conduct a LR or LM test for each possible grouping, then take the maximum. Requires different critical values! (Can be obtained from proposed tests.)
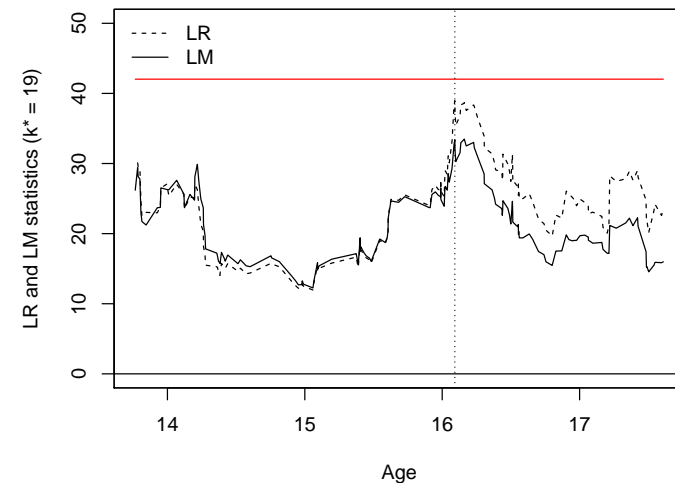
Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Lack of Grouping

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

# Proposed Tests

- In contrast to existing tests of measurement invariance, the proposed tests offer the abilities to:
  - Test for measurement invariance when groups are ill-defined (e.g., when the grouping variable is continuous).
  - Test for measurement invariance in any subset of model parameters.
  - Interpret the nature of measurement invariance violations.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

# Proposed Tests

- The proposed family of tests rely on first derivatives of the model's log-likelihood function.
- We consider individual terms (*scores*) of the gradient. These scores tell us how well a particular parameter describes a particular individual.

$$\sum_{i=1}^{n} s(\hat{\boldsymbol{\theta}}; \mathbf{x}_i) = \mathbf{0}, \text{ where}$$

$$s(\hat{\boldsymbol{\theta}}; \mathbf{x}_i) = \frac{\partial}{\partial \boldsymbol{\theta}} \log \mathsf{L}(\mathbf{x}_i, \boldsymbol{\theta})\big|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}}$$

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

# Proposed Tests

- Under measurement invariance, parameter estimates should roughly describe everyone equally well. So people's scores should fluctuate around zero.

- If measurement invariance is violated, the scores should stray from zero.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

# Aggregating Scores

- We need a way to aggregate scores across people so that we can draw some general conclusions.
  - Order individuals by an auxiliary variable.

  - Define $t \in (1/n, n)$. The *empirical cumulative score process* is defined by:

$$\mathbf{B}(\hat{\boldsymbol{\theta}}; t) = \frac{1}{\sqrt{n}} \sum_{i=1}^{\lfloor nt \rfloor} s(\hat{\boldsymbol{\theta}}; \mathbf{x}_i).$$

  where $\lfloor nt \rfloor$ is the integer part of $nt$.

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Tests

- Under the hypothesis of measurement invariance, a functional central limit theorem holds:

$$\mathbf{I}(\widehat{\boldsymbol{\theta}})^{-1/2}\mathbf{B}(\widehat{\boldsymbol{\theta}};\cdot) \xrightarrow{d} \mathbf{B}^0(\cdot),$$

where $\mathbf{I}(\widehat{\boldsymbol{\theta}})$ is the observed information matrix and $\mathbf{B}^0(\cdot)$ is a $p$-dimensional Brownian bridge.

- Testing procedure: Compute an aggregated statistic of the empirical score process and compare with corresponding quantile of aggregated Brownian motion.

- Test statistics: Special cases include double maximum (DM), Cramér-von Mises (CvM), maximum of LM statistics.

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Simulation

- Simulation: What is the power of the proposed tests?

  - Two-factor model, with three indicators each.
  - Measurement invariance violation in three factor loading parameters, with magnitude from 0–4 standard errors.
  - Sample size in $\{100, 200, 500\}$.
  - Model parameters tested in $\{3, 19\}$.
  - Three test statistics.

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Simulation

Measurement Invariance

Ed Merkle, Achim Zeileis

Background

Proposed Tests

Illustration

Conclusions

## Example

- Example: Studying stereotype threat via factor analysis (Wicherts et al., 2005)
  - Stereotype threat: Knowledge of stereotypes about one's social group might cause one to fulfill the stereotypes.
  - Wicherts et al. study: 295 students were administered three intelligence tests. Stereotypes were primed for half of the students.
  - Groups defined by: Ethnicity (majority/minority) and whether or not stereotypes were primed.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background
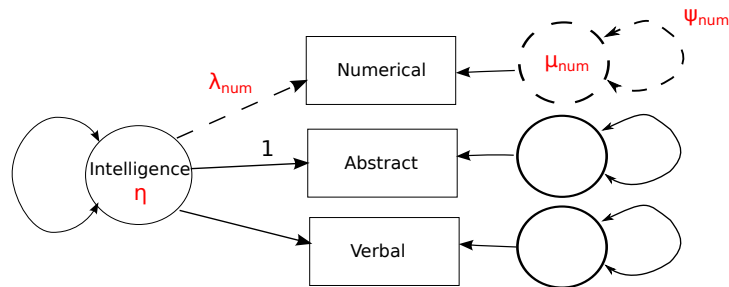
Proposed
Tests

Illustration

Conclusions

# Model

- To study the data, Wicherts et al. employed a series of four-group, one-factor models.
  - General finding: Minorities with stereotype primes have different measurement parameters than other groups.
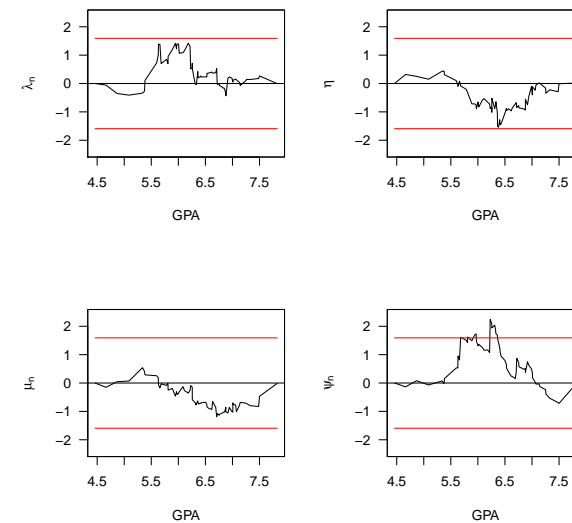  - Current example: Is measurement further impacted by academic performance (as measured by student GPA)?

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

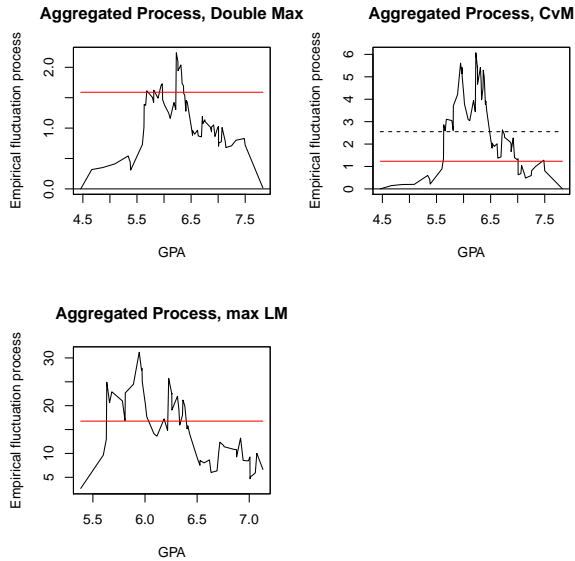Illustration

Conclusions

# Model

- We utilize a model employed by Wicherts et al., where four model parameters are specific to the "minority, stereotype prime" group.
  - Test for measurement invariance in these parameters wrt the student GPA variable (either all four together or individually).
  - Violations of measurement invariance imply that stereotype threat is more problematic for students of low or high GPA.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

# Model



# Results for Single Parameters

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Aggregated Results

**Aggregated Process, Double Max**

**Aggregated Process, CvM**

**Aggregated Process, max LM**

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Conclusions

- Measurement invariance tests utilizing stochastic processes have important advantages over existing tests:
  - Isolating specific parameters that violate measurement invariance, allowing the researcher to define specific types of measurement invariance "post hoc" instead of "a priori".
  - Isolating groups of individuals whose parameter values differ.
  - Studying the impact of continuous variables on model estimates, without "ruining" the rest of the model.
- Power is reasonable, with specific tests being better in specific circumstances.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Software

- To carry out the tests, we utilize
  - `lavaan` for model estimation.
  - `estfun()` for score extraction, which is currently a combination of our own code and `lavaan` code.
  - `strucchange` for carrying out the proposed tests with the scores.
    - Required input: Fitted model, function for score extraction, and information matrix (optional).
    - `gefp()` constructs the process.
    - `sctest()` and `plot()` calculate and visualize test statistics.

Measurement
Invariance

Ed Merkle,
Achim Zeileis

Background

Proposed
Tests

Illustration

Conclusions

## Current Work

- Continued test implementation via `strucchange` and `lavaan` (and possibly `OpenMx`).

- Detailed examination of test properties.

- Extension to related psychometric issues.

- Working paper:
  `http://econpapers.repec.org/RePEc:inn:wpaper:2011-09`

- Questions?