

# Accounting for Population Uncertainty in Covariance Structure Analysis

Hao Wu

Boston College

May 21, 2013

Joint work with:

Michael W. Browne  
The Ohio State University

# Outline

Traditional Approach to Covariance Structures

The New Approach

Sampling Distributions

Simulation Studies

Summary

Model Error

Hao Wu

## Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

Simulation Studies

Summary

- ▶ Covariance matrix among observed variables are usually implied by their relations with theoretically assumed latent variables.
- ▶ For example, the factor analysis model:

$$\begin{aligned}\mathbf{x} &= \mathbf{\Lambda}\mathbf{z} + \mathbf{u} \\ \Omega(\boldsymbol{\xi}) &= \mathbf{\Lambda}\Phi\mathbf{\Lambda}' + \Psi\end{aligned}$$

where  $\Omega$ ,  $\Phi$  and  $\Psi$  are the covariance matrices for  $\mathbf{x}$ ,  $\mathbf{z}$  and  $\mathbf{u}$ , respectively, and  $\boldsymbol{\xi}$  is a vector of all parameters including the unknown loadings, factor correlations and unique variances.

# parameter estimation and test statistic

- ▶ Parameter estimates are obtained by minimizing a discrepancy function  $F$ :

$$\hat{\xi} = \arg \min F(\mathbf{S}, \Omega(\xi))$$

- ▶ We focus on the maximum Wishart likelihood (MWL) discrepancy function

$$F = -\ln |\Omega^{-1} \mathbf{S}| + \text{tr}(\Omega^{-1} \mathbf{S}) - p$$

- ▶ If the true covariance matrix  $\Sigma = \Omega(\xi_0)$  for some  $\xi_0$ ,

$$\begin{aligned} \sqrt{n}(\hat{\xi} - \xi_0) &\xrightarrow{\mathcal{L}} N(0, 2(\Delta^* \mathbf{V}^* \Delta^*)^{-1}) \\ n\hat{F} &\xrightarrow{\mathcal{L}} \chi_{df}^2 \end{aligned}$$

where  $\Delta^* = \partial \omega(\xi_0) / \partial \xi'$ ,  $\mathbf{V}^* = (\Sigma \otimes \Sigma)^{-1}$  and  $df$  is the degrees of freedom of the covariance structure.

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

Simulation Studies

Summary

when the model is not correct:  $\Sigma \neq \Omega$

- ▶ Models are only approximations to the reality.
- ▶ Covariance structures rarely hold exactly in the population.
- ▶ All models will be rejected for exact fit when sample size is large enough.
- ▶ Misspecification must be accounted.

Parameter estimates

- ▶  $\hat{\xi} \xrightarrow{P} \xi^\#$  as  $n \rightarrow \infty$ , where  $\xi^\# = \arg \min F(\Sigma, \Omega(\xi))$  is a function of  $\Sigma$ .
- ▶ Asymptotic normality with variance involving second derivatives of  $\Omega(\xi)$ .

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

Simulation Studies

Summary

# Test Statistics and RMSEA

$$n\hat{F} \xrightarrow{\mathcal{L}} \chi_{df,\delta}^2 \quad \text{as } n \rightarrow \infty$$

under the Pitman drift assumption:

- ▶  $nF^\# \rightarrow \delta$
- ▶ where  $F^\# = F(\Sigma, \Omega^\#)$  and  $\Omega^\# = \Omega(\xi^\#)$ .
- ▶ Misspecification diminishes with increasing sample size.
- ▶ Without this assumption,  $n\hat{F} \rightarrow \infty$ .

RMSEA:  $\varepsilon = \sqrt{F^\#/df}$

- ▶ CI can be obtained from the non-central  $\chi^2$  sampling distribution of  $n\hat{F}$ .
- ▶ Test of close fit:  $\varepsilon \leq 0.05$  vs.  $\varepsilon > 0.05$

# problems of the traditional approach

- ▶ Misspecification is not modeled or explained.
  - ▶ The **same** procedure is used as estimating a correctly specified model.
  - ▶ Misspecification is acknowledged in a post hoc manner.
- ▶ Misspecification is assumed a fixed effect in the population.
  - ▶ Fixed effect: does not change over replications;
  - ▶ Random effect: changes over replications.
  - ▶ Does misspecification change when the measurement is replicated?
- ▶ The Pitman drift assumption is impractical:
  - ▶ misspecification is a property of the population and should not be related to sample size.

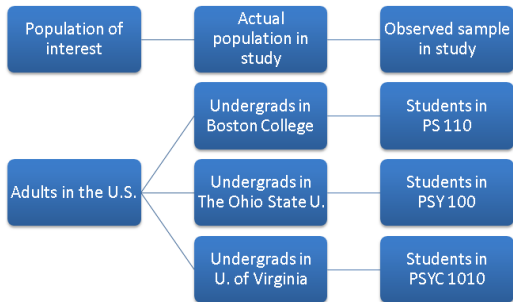
# How does misspecification arise?

Example: In a study, investigators have students in PS 110 fill out a questionnaire and find that the desired structure of the questionnaire is rejected.

- ▶ The study is more likely targeting an ideal population
  - ▶ such as “all adults in the U.S.”;
  - ▶ a more general/standard measurement environment
- ▶ The sample is not representative of the ideal population.
  - ▶ all observations from people of very similar age group
  - ▶ possibly measured under the same incidental effect of unknown sources
- ▶ Misspecification may come from the difference between the ideal and operational populations.



- ▶ If this study were to be replicated, would it be replicated in the same university?
  - ▶ Very unlikely.
  - ▶ Incidental effects may also change: different time of year (e.g. semester) etc.
  - ▶ The operational population has uncertainty, so misspecification is a random effect (model error).



# A new approach

In the new approach to modeling misspecification, we assume

- ▶ Misspecification arises due to stochastic model error.
- ▶ Different replications of the same measurement involve different population covariance matrix  $\Sigma$ .
- ▶ These unstructured  $\Sigma$  is a random effect centered on a structure  $\Omega(\xi_0)$ .

- ▶ Sampling error is modeled in the traditional way:

$$\mathbf{S}|\Sigma \sim W_p(n^{-1}\Sigma, n)$$

- ▶ Model error is also modeled **randomly** with a conjugate distribution:

$$\Sigma|m, \Omega \sim IW_p(m\Omega, m)$$

- ▶ Parameters of interest:
  - ▶  $\xi$  as in covariance structure  $\Omega(\xi)$
  - ▶  $\nu = 1/m$ : misspecification parameter (dispersion parameter)

$\mathbf{\Omega}(\boldsymbol{\theta})$ 

$$\boldsymbol{\Sigma} \mid m, \mathbf{\Omega} \sim \text{IW}_p(m\mathbf{\Omega}, m)$$

 $\boldsymbol{\Sigma}$ 

$$\mathbf{S} \mid \boldsymbol{\Sigma} \sim \text{W}_p(n^{-1}\boldsymbol{\Sigma}, n)$$

 $\mathbf{S}$ 

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

Simulation Studies

Summary

- ▶ The marginal distribution is a matrix-variate Beta distribution:

$$\mathbf{S} | m, \Omega \sim B_p^{\text{II}}\left(\frac{m}{n}\Omega, \frac{n}{2}, \frac{m}{2}\right)$$

- ▶ Parameters can be estimated by maximizing an adjusted marginal likelihood - Maximum Beta Likelihood (MBL).
- ▶ The marginal likelihood function is adjusted to correct for a downwards bias in estimating  $\nu = 1/m$ .

# Relationship to RMSEA

When  $\nu$  is small,

$$\hat{\nu}^{IW} \stackrel{\mathbf{a}}{=} \varepsilon^2$$

- ▶ Our measure of misspecification is asymptotically equivalent to (the square of) RMSEA.
- ▶ Note  $\hat{\nu}^{IW}$  is an estimate from the population model of inverted Wishart distribution.

# Replication framework and consistency

- ▶ Both model error ( $\Sigma|\Omega, m$ ) and sampling error ( $\mathbf{S}|\Sigma, n$ ) are sources of randomness.
- ▶  $\mathbf{S}$  has a Beta distribution.
- ▶  $\hat{\xi} = \xi_0 + o_p(1)$  and  $\hat{v} = v_0 + o_p(v_0 + 1/n)$  as **both**  $v_0 \rightarrow 0$  and  $n \rightarrow \infty$ .

# A Weak Pitman's Drift

The sampling distributions are derived under the assumption of  $n \rightarrow \infty$  and  $v_0 = 1/m_0 \rightarrow 0$ . This assumption is similar to the Pitman's drift assumption

- ▶ the sample size is assumed large, and
- ▶ misspecification is assumed small.

However, there is one major difference:

- ▶ We do **not** restrict the bivariate limit to a particular rate, while Pitman drift assumes  $nF^\# \rightarrow \delta$ , which would be equivalent to  $m_0 = O(n)$  in our system.

This assumption is weaker and is more plausible in practice than Pitman drift assumption.



Under some regularity conditions, the MBLE  $\hat{\xi}$  satisfy

$$\frac{\hat{\xi} - \xi_0}{\sqrt{v_0 + \epsilon}} \xrightarrow{\mathcal{L}} N(0, 2(\Delta^{*\prime} \mathbf{V}^* \Delta^*)^{-1})$$

or

$$\hat{\xi} \stackrel{d}{\approx} N(\xi_0, 2(v_0 + \epsilon)(\Delta^{*\prime} \mathbf{V}^* \Delta^*)^{-1})$$

where  $\epsilon = 1/n$  and  $\mathbf{V}^* = (\boldsymbol{\Omega}_0 \otimes \boldsymbol{\Omega}_0)^{-1}$ . Note the additive effects of the sampling and model error on the dispersion of the parameter estimate.

As  $n \rightarrow \infty$  and  $v_0 \rightarrow 0$ , the MBLE  $\hat{v}$  has sampling distribution

$$\frac{\hat{v}_0 + \epsilon}{v_0 + \epsilon} \xrightarrow{\mathcal{L}} \chi_{df}^2 / df$$

where  $\hat{v}_0$  is defined as

$$\hat{v}_0 = \begin{cases} \hat{v} & \hat{v} > 0 \\ (2df)^{-1} \text{tr} \left\{ (\hat{\Omega} - \mathbf{S}) \hat{\Omega}^{-1} \right\}^2 - \frac{1}{n} & \hat{v} = 0 \end{cases}$$

Confidence bounds or intervals can be obtained by inverting this sampling distribution.

# A t distribution

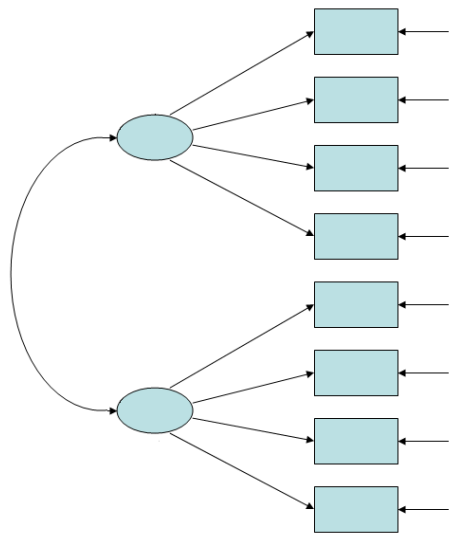
$\hat{\xi}$  and  $\hat{v}_0$  are asymptotically independent. As a result,

$$\frac{\hat{\xi}_i - \xi_{i0}}{\sqrt{2(\hat{v}_0 + \epsilon)[\hat{\Delta}'\hat{V}\hat{\Delta}]^{ii}}} \xrightarrow{\mathcal{L}} t_{df}$$

and a CI on  $\xi_{i0}$  is given by

$$\hat{\xi}_i \pm t_{df, 1-\alpha/2} \sqrt{2(\hat{v}_0 + \epsilon)[\hat{\Delta}'\hat{V}\hat{\Delta}]^{ii}}.$$

# Model Being Used



Model Error

Hao Wu

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

**Simulation Studies**

Summary

A factor analysis model for **correlation** structure

$$\mathbf{\Omega} = \mathbf{D}(\mathbf{\Lambda}\mathbf{\Phi}\mathbf{\Lambda}' + \mathbf{D}_{\psi})\mathbf{D}$$

- ▶ Factor loadings

$$\mathbf{\Lambda} = \begin{pmatrix} 0.5 & 0.5 & 0.6 & 0.6 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.7 & 0.7 & 0.8 & 0.8 \end{pmatrix}'$$

- ▶ Factor correlation  $\rho = 0.5$
- ▶ Unique variances

$$\psi = [0.75, 0.75, 0.64, 0.64, 0.51, 0.51, 0.36, 0.36]'$$

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

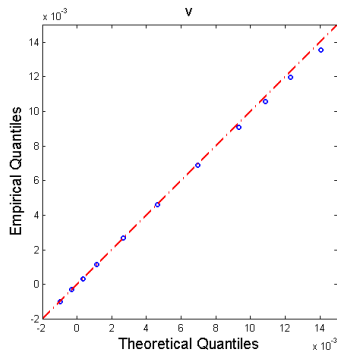
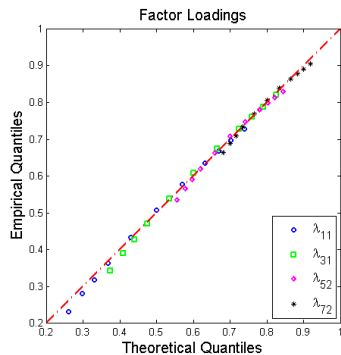
Simulation Studies

Summary

$n$  and  $m$  chosen to be 200, 500, 1000 and  $\infty$ , Monto Carlo sample size 50,000.

- ▶ Study I: 15 combinations of  $(n, m_0)$ , exam of the MBL estimates.
- ▶ Study II:  $n = m = 1000$ , comparison of MBL and MWL (traditional approach).

Selected results from  $n = m = 200$  condition:



# Study II

The missing rates (%) of the MBL and MWL 95% CIs with  $n = m = 1000$ .

	MBL	MWL
$\lambda_{11}$	5.49	32.09
$\lambda_{31}$	5.30	32.05
$\lambda_{52}$	5.34	31.79
$\lambda_{72}$	5.39	31.83
$\rho$	5.59	32.22
$\psi_1$	5.63	32.24
$\psi_3$	5.43	32.15
$\psi_5$	5.31	31.86
$\psi_7$	5.37	31.87



# Summary

A new model for misspecified covariance structures

- ▶ It assumes misspecification arises from a random model error due to population uncertainty, and
- ▶ uses a weak version of the Pitman drift assumption in asymptotics.

We found

- ▶  $\hat{v}^{IW} \approx \varepsilon^2$ .
- ▶ the asymptotic variance of  $\hat{\xi}$  has two additive parts due to sampling and model error respectively.
- ▶ the test statistic has a central  $\chi^2$  distribution

Simulation results show

- ▶ The asymptotic distributions work well.
- ▶ Assuming a random model error, MWL CIs have poor coverage as it fails to account for the extra source of variability.

Outline

Traditional  
Approach to  
Covariance  
Structures

The New Approach

Sampling  
Distributions

Simulation Studies

Summary

Thank you!

Questions?