Consequences of measurement error

Psychology 588: Covariance structure and factor models

Scaling indeterminacy of latent variables

- Scale of a latent variable is arbitrary and "determined" by a convention for convenience
- Typically set to variance of one (factor analysis convention) or to be identical to an arbitrarily chosen indicator's scale

By centering indicator variables, we set latent variables' means to zero

• Consider the following transformation:

$$\begin{split} x_{j} &= v_{j} + \lambda_{j} \xi + \delta_{j}, \quad j = 1, \dots, J, \quad \xi^{*} = a + b \xi, \quad b \neq 0 \\ &= \left(v_{j} - \frac{a}{b} \lambda_{j} \right) + \left(\frac{\lambda_{j}}{b} \right) \xi^{*} + \delta_{j} \end{split}$$

• If all *J* indicators are considered simultaneously, vector notation is more convenient:

$$\mathbf{x} = \mathbf{v} + \lambda \boldsymbol{\xi} + \boldsymbol{\delta}, \quad \boldsymbol{\xi}^* = a + b\boldsymbol{\xi}$$
$$= \left(\mathbf{v} - \frac{a}{b}\lambda\right) + \left(\frac{1}{b}\lambda\right) \boldsymbol{\xi}^* + \boldsymbol{\delta}$$

meaning that the linear transformation of ξ can be exactly compensated in the accordingly transformed $\mathbf{v}^* = \mathbf{v} - \lambda a/b$ and $\lambda^* = \lambda/b$, leaving the errors δ unchanged (i.e., same fit)

What's great about measurement errors in equation 4

- Regression weights and correlations are interpreted, implicitly assuming that the "operationally defined" variables involve no measurement error --- hardly realized for theoretical constructs (e.g., self esteem, IQ, etc.)
- Ignoring the measurement error will lead to inconsistent estimates
- We will see consequences of ignoring measurement errors

• Consider a mean-included equation for X (hours worked per week) to indicate ξ (achievement motivation):

$$X = \nu + \lambda \xi + \delta, \quad E(\xi) = \kappa, \quad E(\delta) = 0, \quad E(\xi\delta) = 0$$
$$E(X) = \mu_X = \nu + \lambda \kappa$$
$$\operatorname{var}(X) = \lambda^2 \phi + \operatorname{var}(\delta)$$

Given only one indicator per latent variable, the intercept and loading (i.e., weight) are simply scaling constants for ξ

However, if the ξ scale is set comparable to the *X* scale (i.e., $\lambda = 1$), we see that var(X) is an over-estimation of $\phi = var(\xi)$ if δ is not included in the equation

Bivariate relation and simple regression

True data structure: $x = \lambda_1 \xi + \delta$ $y = \lambda_2 \eta + \varepsilon$ η : job satisfaction y: satisfaction scale $\eta = \gamma \xi + \zeta$



$$\operatorname{cov}(\xi,\eta) = \operatorname{cov}(\xi,\gamma\xi+\zeta) = \gamma\phi$$
$$\operatorname{cov}(x,y) = \operatorname{cov}(\xi+\delta,\gamma\xi+\zeta+\varepsilon) = \gamma\phi$$

zeta

eta

У

е

• From the previous equations, $\gamma = cov(\xi, \eta)/\phi$ and by analogy with $y = \gamma^* x + \zeta^*$ if measurement errors are ignored,

$$\gamma^* = \frac{\operatorname{cov}(x, y)}{\operatorname{var}(x)} = \gamma \left(\frac{\phi}{\phi + \operatorname{var}(\delta)}\right) = \gamma \rho_{xx}$$

The parenthesized ratio (reliability) becomes 1 only with no measurement error; otherwise, γ^* is an attenuated estimate of γ and $\hat{\gamma}^* = s_{xy}/s_{xx}$ is an inconsistent estimator of γ

• If $\lambda_1 \neq \lambda_2$, the bias of regression weight has an additional factor as $\gamma^* = (\lambda_2/\lambda_1)\gamma \rho_{xx}$ --- but such scaling is unusual when there is only one indicator per latent variable

• Correlations:

$$\rho_{\xi\eta}^{2} = \frac{\operatorname{cov}(\xi,\eta)^{2}}{\phi\operatorname{var}(\eta)} = \frac{\gamma^{2}\phi}{\operatorname{var}(\eta)}$$
$$\rho_{xy}^{2} = \frac{\operatorname{cov}(x,y)^{2}}{\operatorname{var}(x)\operatorname{var}(y)} = \frac{\gamma^{2}\phi^{2}}{\operatorname{var}(x)\operatorname{var}(y)}$$
$$= \frac{\phi}{\operatorname{var}(x)}\frac{\operatorname{var}(\eta)}{\operatorname{var}(y)}\frac{\gamma^{2}\phi}{\operatorname{var}(\eta)} = \rho_{xx}\rho_{yy}\rho_{\xi\eta}^{2}$$

which shows an attenuation of the "true" correlation due to measurement error, with the familiar correction formula:

$$\rho_{\xi\eta} = \rho_{xy} / \sqrt{\rho_{xx} \rho_{yy}}$$

• True data structure:

$$\eta = \gamma' \xi + \zeta$$
$$\mathbf{x} = \xi + \mathbf{\delta}$$

 $y = \eta + \varepsilon$



with $\Lambda_x = \mathbf{I}$ and $\lambda_y = 1$

• Ignoring measurement errors: $y = \gamma^{*'} \mathbf{x} + \zeta^{*}$

•
$$\sigma_{\xi\eta} = \operatorname{cov}(\xi,\eta) = \operatorname{cov}(\xi,\xi'\gamma+\zeta) = \Phi\gamma$$

$$\boldsymbol{\sigma}_{xy} = \operatorname{cov}(\mathbf{x}, y) = \operatorname{cov}(\boldsymbol{\xi} + \boldsymbol{\delta}, \boldsymbol{\xi}' \boldsymbol{\gamma} + \boldsymbol{\zeta} + \boldsymbol{\varepsilon}) = \boldsymbol{\Phi} \boldsymbol{\gamma}$$

• $\gamma = \Phi^{-1} \sigma_{\xi\eta}$ and by analogy with $y = \gamma^{*'} \mathbf{x} + \zeta^{*}$,

$$\boldsymbol{\gamma}^* = \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\sigma}_{xy} = \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Phi} \boldsymbol{\gamma} = \left(\boldsymbol{\Phi} + \boldsymbol{\Theta}_{\delta} \right)^{-1} \boldsymbol{\Phi} \boldsymbol{\gamma}$$

Without measurement error ($\Theta_{\delta} = 0$), $\gamma^* = \gamma$; otherwise, $\gamma^* \neq \gamma$

• Alternatively written: $\gamma^* = \Sigma_{xx}^{-1} \Sigma_{x\xi} \gamma$ since $\Sigma_{x\xi} = \Phi$ --- where $\Sigma_{xx}^{-1} \Sigma_{x\xi}$ is the OLS estimator of **B** in $\xi = \mathbf{B}\mathbf{x} + \mathbf{e}$, i.e., regression weights for prediction of ξ by \mathbf{x}

Again, without measurement error, $\Sigma_{xx}^{-1}\Sigma_{x\xi} = \mathbf{I}$

• Note: in Bollen (pp. 159-168), Γ , $\Sigma_{\xi\eta}$, Σ_{xy} are meant to be γ , $\sigma_{\xi\eta}$, σ_{xy} , respectively, for the multiple regression model

• As a very simplified case, suppose x_1 is the only fallible as:

$$x_1 = \xi_1 + \delta_1$$
$$x_i = \xi_i, \quad i = 2, \dots, q$$

with the true and estimated regression equations:

$$\eta = \gamma_1 \xi_1 + \gamma_2 \xi_2 + \dots + \gamma_q \xi_q + \zeta$$
$$\eta = \gamma_1^* x_1 + \gamma_2^* \xi_2 + \dots + \gamma_q^* \xi_q + \zeta^*$$

• In this special case, the regression weight matrix has a simple multiplicative form of bias (hint: use $\tilde{\Phi} = \Phi + \Theta_s$):

$$\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{x\xi} = \left(\boldsymbol{\Phi} + \boldsymbol{\Theta}_{\delta}\right)^{-1}\boldsymbol{\Phi} = \mathbf{I} + \left[\mathbf{c}, \mathbf{0}_{q \times (q-1)}\right] = \left[\frac{1 + c_1 |\mathbf{0}'|}{\mathbf{c}_2 |\mathbf{I}|}\right]$$

• Consequently, resulting bias factors are:

$$\gamma_1^* = b_{\xi_1 x_1 \cdot \xi_2 \cdots \xi_q} \gamma_1$$

$$\gamma_i^* = \gamma_i + b_{\xi_1 \xi_i \cdot x_1 \{\xi_{\sim i}\}} \gamma_1, \quad i = 2, \dots, q$$

- > Bias-factor for x_1 is less than 1 in absolute value (1 without measurement error), and so γ_1^* is biased toward 0 --- the bias factor $b_{\xi_1 x_1 \cdot \xi_2 \cdots \xi_q}$ indicates regression weight b_1 in $\xi_1 = b_0 + b_1 x_1 + b_2 \xi_2 + \ldots + b_q \xi_q$
- > Consequences for x_i , i = 2, ..., q are additive, depending on relationships between ξ_1 and ξ_i holding all other IVs constant, and γ_1

- So far all reasoning is based on rather unrealistic assumptions:
 - Only single indicator per latent variable, and so its loading becomes simply scaling constant
 - > Only one fallible IV
- Without such assumptions (e.g., all IVs fallible), consequences of measurement error become too complicated and hard to simplify algebraically --- no particular simple form of $\Sigma_{xx}^{-1}\Sigma_{x\xi}$
- One clear conclusion: all estimates are inconsistent --systematically different from what they meant to be

• Consequence in standardization:

standardized
$$\gamma_i^* = \gamma_i^* \sqrt{\frac{\phi_{ii} + \operatorname{var}(\delta_i)}{\operatorname{var}(\eta) + \operatorname{var}(\varepsilon)}}$$

• Consequence in SMC is similar to the bivariate case:

 $\operatorname{plim}(R^2) \ge \operatorname{plim}(R^{*2})$

- What should we do with essentially omnipresent measurement error?
 - Use SEM which allows for measurement errors in the model --- though we are limited in certain models regarding the model identification (e.g., Table 5.1, p. 164)

• Consequence in regression weights further complicated:

$$\boldsymbol{\gamma}^* = \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{x\xi} \boldsymbol{\gamma} + \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\sigma}_{\delta\varepsilon}$$

For simple regression:
$$\gamma^* = \gamma \rho_{xx} + \frac{\operatorname{cov}(\varepsilon, \delta)}{\operatorname{var}(x)}$$

Now, γ^* is not necessarily < γ

• If correlated measurement errors are only within IVs (i.e., $\sigma_{\delta \varepsilon} = 0$, $\Sigma_{xx} = \Phi + \Theta_{\sigma}$ where Θ_{σ} is not diagonal), $\gamma^* = \Sigma_{xx}^{-1} \Sigma_{x\xi} \gamma$ still holds (but the bias factor will have a more complicated form, also involving off-diagonal entries of Θ_{σ})

- In path models with sequential causal paths, consequences of measurement errors very hard to simply generalize --- see the union sentiment (Fig. 5.2, p. 169) and SES (Fig. 5.4, p. 173) examples
- If reliabilities are known, the corresponding error variances can be constrained; if unknown, the error variances may be modeled as free parameters provided that they are identifiable
- To keep in mind: we need more than one indicator per latent variable for identifiability and statistical testing --- leading to measurement models with multiple indicators or CFA