

# Robust Confidence Intervals for Effects Sizes in Multiple Linear Regression

Paul Dudgeon

Melbourne School of Psychological Sciences  
The University of Melbourne. Vic. 3010  
AUSTRALIA  
dudgeon@unimelb.edu.au

Modern Modeling Methods (M<sup>3</sup>) Conference  
University of Connecticut  
20-21 May 2014

- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models
- 3 Standard Errors in Linear Regression
- 4 Standard Errors in SEM Framework
- 5 Simulation Study
- 6 Brief Summary Results of Simulation Study
- 7 Thanks

# Effect Sizes for IVs in Multiple Linear Regression

- Unstandardised regression coefficients ( $B_j$ ) in linear regression are natural effect sizes.
- Additional effect sizes measures:
  - 1 Standardized regression coefficient ( $B_j^*$ ).
  - 2 Semipartial correlation ( $sr_j$ ).
  - 3 Improvement in R-squared ( $\Delta R_j^2$ ).
- These three effect size measures are the primary focus of this talk.
  - Develop confidence intervals robust to assumption violations.
  - Structural equation modelling as a unified framework for deriving and calculating standard errors.

## Previous Research

- Limited previous research in this area.
- Algina and colleagues (2001, 2007 2008, 2010) investigated various methods for deriving confidence intervals for the improvement in R-square.
- Aloe & Becker (2012) derived a standard error for the semipartial correlation for use in meta-analysis.
- Yuan & Chan (2011) developed standard errors for standardized regression coefficients.
- With one exception, these approaches all rest on standard regression assumptions holding.

# The Linear Regression Model

- The linear regression model for  $J$  independent variables measured on  $n = 1, \dots, N$  individuals is given by:

$$Y_n = B_0 + B_1 X_{n1} + B_2 X_{n2} + \dots + B_J X_{nJ} + e_n,$$

- Estimated values for the regression coefficients can be obtained by the ordinary least squares estimator

$$\min_{(\hat{B}_0, \dots, \hat{B}_J)} = \sum_{n=1}^N (Y_n - \hat{Y}_n)^2$$

where

$$\hat{Y}_n = B_0 + B_1 X_{n1} + B_2 X_{n2} + \dots + B_J X_{nJ}$$

## Deriving Focal Effect Sizes

- All three effect sizes can be related to the unstandardized regression coefficient.

- 1 Standardized regression coefficient.

$$B_j^* = B_j \times \frac{S_{X_j}}{s_Y}$$

- 2 Standardized regression coefficient.

$$s r_j = B_j^* \times \sqrt{1 - R_{(j)}^2}$$

- 3 Improvement in R-squared.

$$\Delta R_j^2 = R^2 - R_j^2 = s r_j^2$$

where  $1 - R_{(j)}^2$  is the tolerance of the j-th independent variable, and  $R_j^2$  is the reduced model coefficient of determination when  $X_j$  is removed.

- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models**
- 3 Standard Errors in Linear Regression
- 4 Standard Errors in SEM Framework
- 5 Simulation Study
- 6 Brief Summary Results of Simulation Study
- 7 Thanks

# Linear Regression and SEM

- LS estimator for linear regression is equivalent to the maximum likelihood (ML) estimator.
- Any linear regression analysis can therefore be readily expressed as a (saturated) structural equation model (SEM).

$$\Sigma = \Sigma(\theta_B)$$

where the set of freely estimated model parameters equals

$$\theta_B = [\beta_1, \dots, \beta_J, \sigma_\varepsilon^2, \text{vech}(\Sigma_X)]$$

and  $\text{vech}(\bullet)$  is the matrix operator extracting non-duplicated elements of a symmetric matrix into a vector.



# Linear Regression and Correlational Structures

- For linear regression, the SEM specification is straightforward:

$$\Sigma(\theta_{LR}) = \left[ \begin{array}{c|c} \beta' \Sigma_X \beta + \sigma_\varepsilon^2 & \beta' \Sigma_X \\ \hline \Sigma_X \beta & \Sigma_X \end{array} \right]$$

- Many covariance structures can be equivalently estimated as correlational structures (with appropriate constraints).

$$\Sigma_{YX} = \mathbf{D}_\sigma \mathbf{P}(\theta_{LR}^*) \mathbf{D}_\sigma,$$

where

$$\text{Diag}[\mathbf{P}(\theta_{LR}^*)] = \mathbf{I}_k$$

is the imposed constraint function,  $\mathbf{D}_\sigma$  is a diagonal scaling matrix, and  $\theta_{LR}^*$  contains parameters in a standardized metric.

# SEM for Standardized Regression Coefficients

- Standardized regression coefficients can be obtained therefore as:

$$\Sigma_{\gamma X} = \mathbf{D}_\sigma \left[ \begin{array}{c|c} \beta^{*'} P_X \beta^* + \sigma_\varepsilon^2 & \beta^{*'} P_X \\ \hline P_X \beta^* & P_X \end{array} \right] \mathbf{D}_\sigma$$

with the constraint function

$$\sigma_\varepsilon^2 = 1 - \beta^{*'} P_X \beta^*$$

- The set of freely estimated model parameters now equals:

$$\theta_{\beta^*} = [\beta^*, \sigma_\gamma, \sigma_X, \text{vecp}(P_X)],$$

where  $\text{vecp}(\bullet)$  is the matrix operator extracting lower-diagonal elements from a symmetric matrix into a vector.

# SEM for Semipartial Correlation Coefficients

- Using the equivalent constraint function, semipartial correlation coefficients can be obtained as:

$$\Sigma_{YX} = \mathbf{D}_\sigma \left[ \begin{array}{c|c} \beta^{*'} P_X \beta^* + \sigma_\varepsilon^2 & \beta^{*'} P_X \\ \hline P_X \beta^* & P_X \end{array} \right] \mathbf{D}_\sigma$$

where standardized regression coefficients are computed by the user-defined parameter function:

$$\beta^* = \rho_{sr} \times (\mathbf{1} - \text{Diag}[P_X^{-1}])^{-0.5}$$

- The set of freely estimated model parameters equals:

$$\theta_{sr} = [\rho_{sr}, \sigma_Y, \sigma_X, \text{vecp}(P_X)],$$

with the same constraint function on  $\sigma_\varepsilon^2$  being imposed.

## SEM for Improvement in R-squared

- Using the equivalent constraint function, improvement in R-squared values can be obtained as:

$$\Sigma_{YX} = \mathbf{D}_\sigma \left[ \frac{\boldsymbol{\beta}^{*'} \mathbf{P}_X \boldsymbol{\beta}^* + \sigma_\varepsilon^2}{\mathbf{P}_X \boldsymbol{\beta}^*} \mid \frac{\boldsymbol{\beta}^{*'} \mathbf{P}_X}{\mathbf{P}_X} \right] \mathbf{D}_\sigma$$

where standardized regression coefficients are instead computed by the user-defined parameter function:

$$\boldsymbol{\beta}^* = \kappa [\Delta_{R^2} \times (\mathbf{1} - \text{Diag}[\mathbf{P}_X^{-1}])^{-0.5}]$$

and where  $\kappa$  is a sign vector with values equal to +1 or -1.

- The set of freely estimated model parameters equals:

$$\boldsymbol{\theta}_{\Delta_{R^2}} = [\Delta_{R^2}, \sigma_Y, \boldsymbol{\sigma}_X, \text{vecp}(\mathbf{P}_X)],$$

with the same constraint function on  $\sigma_\varepsilon^2$  once again being imposed.

- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models
- 3 Standard Errors in Linear Regression**
- 4 Standard Errors in SEM Framework
- 5 Simulation Study
- 6 Brief Summary Results of Simulation Study
- 7 Thanks

## Standard Errors for Regression Coefficients

- Covariance matrix of regression coefficient parameters under OLS equals (where  $\tilde{\mathbf{X}} = [\mathbf{1}|\mathbf{X}]$ ):

$$\text{Cov}(\hat{\boldsymbol{\beta}})_{OLS} = (N - J)^{-1} \sum_{n=1}^N e_n^2 (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1}$$

- White (1980) proposed a heteroscedastic-consistent estimator when the regression model is misspecified:

$$\text{Cov}(\hat{\boldsymbol{\beta}})_{HC0} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1} (\tilde{\mathbf{X}}' \text{Diag}[e_n^2] \tilde{\mathbf{X}}) (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1}$$

- Hinkley (1977) had earlier and independently proposed a finite sample-adjusted version of HC0 that is less biased:

$$\text{Cov}(\hat{\boldsymbol{\beta}})_{HC1} = \frac{N}{N - J} (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1} (\tilde{\mathbf{X}}' \text{Diag}[e_n^2] \tilde{\mathbf{X}}) (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1}$$

## Two Further Developments

- MacKinnon & White (1985) subsequently proposed two other HC-based estimators to reduce bias further.
- Let  $h_n = \mathbf{x}'_n (\mathbf{X}\mathbf{X}')^{-1} \mathbf{x}_n$  signify a person's leverage statistic.

$$\text{Cov}(\hat{\boldsymbol{\beta}})_{HC2} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1} \left( \tilde{\mathbf{X}}' \text{Diag} \left[ \frac{e_n^2}{(1-h_n)} \right] \tilde{\mathbf{X}} \right) (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1}$$

$$\text{Cov}(\hat{\boldsymbol{\beta}})_{HC3} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1} \left( \tilde{\mathbf{X}}' \text{Diag} \left[ \frac{e_n^2}{(1-h_n)^2} \right] \tilde{\mathbf{X}} \right) (\tilde{\mathbf{X}}\tilde{\mathbf{X}}')^{-1}$$

- Long & Ervin (2000) recommended HC3 as the best general estimator among heteroscedastic-consistent ones.
- All these estimators use raw sample data.

- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models
- 3 Standard Errors in Linear Regression
- 4 Standard Errors in SEM Framework**
- 5 Simulation Study
- 6 Brief Summary Results of Simulation Study
- 7 Thanks



## Normal Standard Errors in SEM

- Under multivariate normality, the covariance matrix of model parameters using ML estimation equals:

$$\text{Cov}(\hat{\theta})_{NT} = N^{-1} (\Delta'_{\theta} \Omega \Delta_{\theta})^{-1}$$

where  $\Omega$  is the asymptotic covariance matrix of population covariances for both the dependent and  $J$  independent variables.

$$\Omega = [2 \times \mathbf{K}'_p (\Sigma \otimes \Sigma) \mathbf{K}_p]^{-1}$$

and  $\mathbf{K}_p$  is a transition matrix.

- For linear regression:

$$\text{Cov}(\hat{\theta})_{NT} \equiv \text{Cov}(\hat{\beta})_{OLS}.$$

## Robust Standard Errors in SEM

- SEs based on  $Cov(\hat{\theta})_{NT}$  are too small if data are non-normal. A more robust set can be derived using the sandwich estimator:

$$Cov(\hat{\theta})_{RO} = N^{-1} (\Delta'_{\theta} \Omega \Delta_{\theta})^{-1} (\Delta'_{\theta} \Omega \tilde{\Omega} \Omega \Delta_{\theta}) (\Delta'_{\theta} \Omega \Delta_{\theta})^{-1}$$

- Yuan & Hayashi (2006) show the correct specification of  $\tilde{\Omega}$  is in general given by:

$$\tilde{\Omega} = N^{-1} \sum_{n=1}^N vech[\tilde{\Sigma}_n] vech[\tilde{\Sigma}_n]'$$

where

$$\tilde{\Sigma}_n = (\mathbf{X}_n - \bar{\mathbf{X}}) (\mathbf{X}_n - \bar{\mathbf{X}})' - \Sigma$$

- For linear regression, it can be shown that:

$$Cov(\hat{\theta})_{RO} \equiv Cov(\hat{\beta})_{HCO}.$$

## Extending Robust SEs to the HC3 Estimator

- Robust SEs in SEM do not currently cover MacKinnon & White's (1985) improved heteroscedastic-consistent estimators.
- Performing linear regression within a SEM framework enables this to be accomplished as follows:

$$e_n^2 \equiv \text{vech} \left[ \tilde{\Sigma}_n \right] \text{vech} \left[ \tilde{\Sigma}_n \right]'$$

Therefore, we can define a more robust estimator of the asymptotic covariance matrix of population covariance as:

$$\tilde{\Omega}^* = (1/N) \sum_{n=1}^N \text{vech} \left[ \tilde{\Sigma}_n \right] (1 - h_n)^{-2} \text{vech} \left[ \tilde{\Sigma}_n \right]'$$

which results in

$$\text{Cov}(\hat{\theta})_{R3} = N^{-1} \left( \Delta_{\theta}' \Omega \Delta_{\theta} \right)^{-1} \left( \Delta_{\theta}' \Omega \tilde{\Omega}^* \Omega \Delta_{\theta} \right) \left( \Delta_{\theta}' \Omega \Delta_{\theta} \right)^{-1}$$

## Extending Robust SEs to the HC3 Estimator (cont.)

- This result means that:

$$\text{COV}(\hat{\theta})_{R3} \equiv \text{COV}(\hat{\beta})_{HC3}.$$

- The SEM-based framework, however, is far more flexible and straightforward.
- Two steps are required:
  - 1 Select one of  $\theta_{\beta}$  or  $\theta_{\beta}^*$  or  $\theta_{SR}$  or  $\theta_{\Delta_{R^2}}$  for  $\hat{\theta}$  in  $\text{COV}(\hat{\theta})_{R3}$ .
  - 2 Calculate the appropriate Jacobian matrix  $\Delta_{\theta}$  for the chosen parameter vector.
- The second step can be easily implemented using automatic differentiation algorithms.

- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models
- 3 Standard Errors in Linear Regression
- 4 Standard Errors in SEM Framework
- 5 Simulation Study**
- 6 Brief Summary Results of Simulation Study
- 7 Thanks

# Investigating Confidence Interval Coverage

- ① Used a  $3 \times 5 \times 4 \times 4 \times 4$  experimental design for simulation.
  1. Number of IVs {2, 5, 8}
  2. Semipartial Correlation {0, .05, .15, .35, .55}
  3. Sample size {50, 100, 300, 1000}
  4. Interim R-squared Value {0, .10, .20, .50}
  5. Tolerance of IV {0.20, 0.50, 0.80, 1.00}
- ② 10,000 replications in each of the 960 cells of the design.
- ③ Data for all variables had excess kurtosis equal to 7.0.
- ④ Generated using methods proposed by Headrick & Kowalchuck (2008).
- ⑤ Calculated proportion of replications in which 95% confidence intervals captured the designated population effect size value.
- ⑥ Two outcome measures: *accuracy* and *robustness*.

## Investigating Confidence Interval Coverage (Cont.)

- 7 **Accuracy** was measured by calculating mean square error (MSE), which can itself be further decomposed into (squared) *bias* and *imprecision*:

$$\mathcal{E} \left[ (\hat{\theta} - \theta)^2 \right] = \left( \mathcal{E} [\hat{\theta} - \theta] \right)^2 + \mathcal{E} \left[ (\hat{\theta} - \mathcal{E}[\hat{\theta}])^2 \right]$$

- 8 **Robustness** was assessed by Serlin's (2000) range null hypothesis test, using Bradley's liberal criterion (i.e.,  $\pm 2.5\%$ ).
- 9 The resultant measure was the proportion of replications for which the Serlin's null hypothesis of non-robustness was rejected.
- 10 Measures of *accuracy* and *robustness* are complimentary, in that the most accurate interval may not necessarily be robust whereas a robust interval may not be the most accurate.

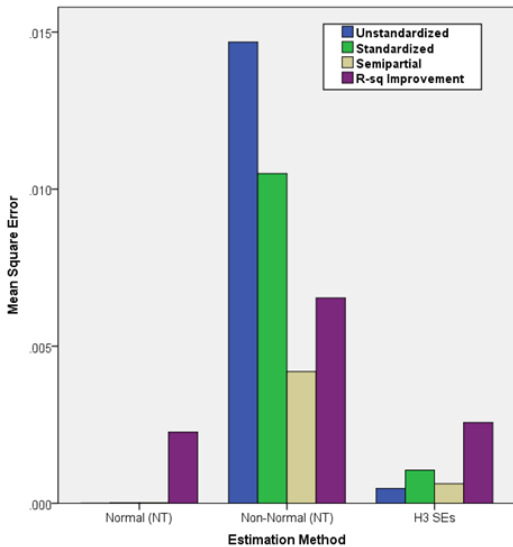
## Brief Explanation about Confidence intervals

- All 95% confidence intervals were calculated using the relevant  $t$  distribution for better small sample performance.
- Confidence intervals for *semipartial correlation* and *improvement in R-squared* were calculated using a method proposed in Browne (1982).
- General approach is to:
  - (a) apply an appropriate transformation on the bounded parameter to make it unbounded,
  - (b) construct the symmetric confidence interval in usual way, and then
  - (c) invert the transformation on both limits to obtain a bounded, asymmetric interval
- All calculations and simulation was undertaken in MATLAB.

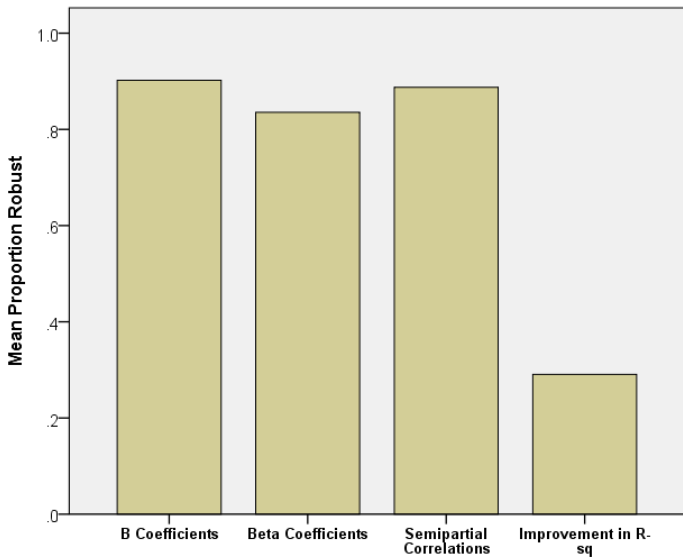


- 1 Effect Sizes in Multiple Linear Regression
- 2 Linear Regression as Structural Equation Models
- 3 Standard Errors in Linear Regression
- 4 Standard Errors in SEM Framework
- 5 Simulation Study
- 6 Brief Summary Results of Simulation Study**
- 7 Thanks

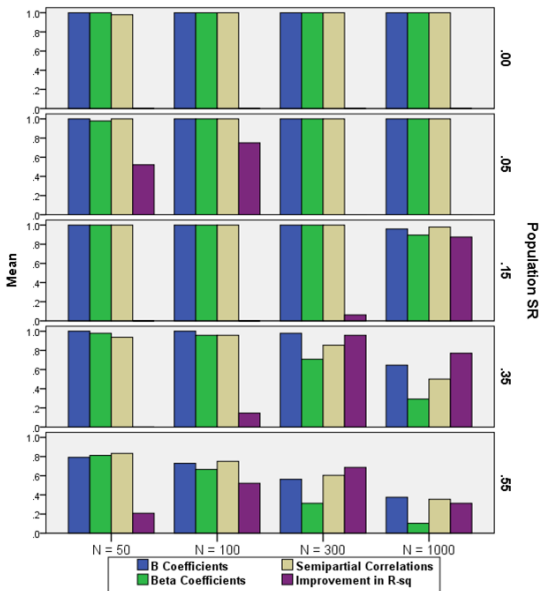
# Accuracy of Confidence Intervals



# Robustness of Confidence Intervals



# Robustness of Confidence Intervals (cont.)



# Summary and Further Potentially Developments

- The robust SEs work far better than normal ones under non-normality for 2 of 3 effect sizes under the majority of conditions.
- Preliminary evidence that the same applies under heteroscedasticity.
- General SEM framework is quite flexible to include other effect sizes such as model R-square and (potentially) relative importance indices.
- Other estimators, such as Yuan & Bentler's (1999) robust covariance methods using Huber weights, could also be investigated.

**And Finally...**

**Thank You for Coming and Listening.**

**Any Questions?**