

Double Measurement Regression, Part Two¹

STA 2101 Fall 2019

¹See last slide for copyright information.

Overview

- 1 The general model
- 2 The BMI study

Double measurement

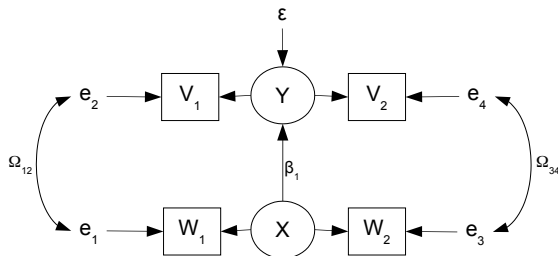
- We have studied an example where two independent measurements of a latent explanatory variable made all the model parameter identifiable.
- Extend the model.
- Double measurement can also help with correlated measurement error.

Correlated measurement error

- We are “measuring” exercise and snack food consumption by self-report.
- A simple additive model: What people report is the truth, plus a piece of noise that pushes the number up or down by a random amount that is different for each person.
- Is it reasonable to assume the error term for snack food is independent of the error term for exercise?
- This is another case of omitted variables.

- Acres planted by farmer’s report and aerial photograph is a different story.
- Double measurement can help with correlated measurement error.

The general double measurement design



These are all matrices.

- The main idea is that \mathbf{X} and \mathbf{Y} are each measured twice, perhaps at different times using different methods.
- Measurement errors V may be correlated within but not between sets of measurements.

Double Measurement Regression: A Two-Stage Model

Setting up a two-stage proof of identifiability

$$\mathbf{Y}_i = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{X}_i + \boldsymbol{\epsilon}_i$$

$$\mathbf{F}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix}$$

$$\mathbf{D}_{i,1} = \boldsymbol{\nu}_1 + \mathbf{F}_i + \mathbf{e}_{i,1}$$

$$\mathbf{D}_{i,2} = \boldsymbol{\nu}_2 + \mathbf{F}_i + \mathbf{e}_{i,2}$$

Observable variables are $\mathbf{D}_{i,1}$ and $\mathbf{D}_{i,2}$: both are $(p + q) \times 1$.

$E(\mathbf{X}_i) = \boldsymbol{\mu}_x$, $cov(\mathbf{X}_i) = \boldsymbol{\Phi}_x$, $cov(\boldsymbol{\epsilon}_i) = \boldsymbol{\Psi}$, $cov(\mathbf{e}_{i,1}) = \boldsymbol{\Omega}_1$,
 $cov(\mathbf{e}_{i,2}) = \boldsymbol{\Omega}_2$. Also, \mathbf{X}_i , $\boldsymbol{\epsilon}_i$, $\mathbf{e}_{i,1}$ and $\mathbf{e}_{i,2}$ are independent.

Measurement errors may be correlated

Look at the measurement model

$$\mathbf{F}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix}$$

$$\mathbf{D}_{i,1} = \boldsymbol{\nu}_1 + \mathbf{F}_i + \mathbf{e}_{i,1}$$

$$\mathbf{D}_{i,2} = \boldsymbol{\nu}_2 + \mathbf{F}_i + \mathbf{e}_{i,2}$$

$$\text{cov}(\mathbf{e}_{i,1}) = \boldsymbol{\Omega}_1 = \left(\begin{array}{c|c} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \hline \boldsymbol{\Omega}_{12}^\top & \boldsymbol{\Omega}_{22} \end{array} \right)$$

$$\text{cov}(\mathbf{e}_{i,2}) = \boldsymbol{\Omega}_2 = \left(\begin{array}{c|c} \boldsymbol{\Omega}_{33} & \boldsymbol{\Omega}_{34} \\ \hline \boldsymbol{\Omega}_{34}^\top & \boldsymbol{\Omega}_{44} \end{array} \right)$$

Expected values of the observable variables

$$\mathbf{D}_{i,1} = \boldsymbol{\nu}_1 + \mathbf{F}_i + \mathbf{e}_{i,1} \text{ and } \mathbf{D}_{i,2} = \boldsymbol{\nu}_2 + \mathbf{F}_i + \mathbf{e}_{i,2}$$

$$E(\mathbf{D}_{i,1}) = \begin{pmatrix} \boldsymbol{\mu}_{1,1} \\ \boldsymbol{\mu}_{1,2} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu}_{1,1} + E(\mathbf{X}_i) \\ \boldsymbol{\nu}_{1,2} + E(\mathbf{Y}_i) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu}_{1,1} + \boldsymbol{\mu}_x \\ \boldsymbol{\nu}_{1,2} + \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{\mu}_x \end{pmatrix}$$

$$E(\mathbf{D}_{i,2}) = \begin{pmatrix} \boldsymbol{\mu}_{2,1} \\ \boldsymbol{\mu}_{2,2} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu}_{2,1} + E(\mathbf{X}_i) \\ \boldsymbol{\nu}_{2,2} + E(\mathbf{Y}_i) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\nu}_{2,1} + \boldsymbol{\mu}_x \\ \boldsymbol{\nu}_{2,2} + \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \boldsymbol{\mu}_x \end{pmatrix}$$

- $\boldsymbol{\nu}_1$, $\boldsymbol{\nu}_2$, $\boldsymbol{\beta}_0$ and $\boldsymbol{\mu}_x$ parameters appear only in expected value, not covariance matrix.
- \mathbf{X}_i is $p \times 1$ and \mathbf{Y}_i is $q \times 1$.
- Even with $\boldsymbol{\beta}_1$ identified from the covariance matrix, have $2(p+q)$ equations in $3(p+q)$ unknown parameters.
- Identifying the expected values and intercepts is impossible.
- Re-parameterize, absorbing them into $\boldsymbol{\mu} = E \begin{pmatrix} \mathbf{D}_{i,1} \\ \mathbf{D}_{i,2} \end{pmatrix}$.

Losing the intercepts and expected values by re-parameterization

- We cannot identify ν_1 , ν_2 , β_0 and μ_x separately.
- Swallow them into μ .
- Estimate μ with $\bar{\mathbf{D}}$.
- And it disappears from $L(\mu, \Sigma) = |\Sigma|^{-n/2} (2\pi)^{-np/2} \exp -\frac{n}{2} \left\{ \text{tr}(\hat{\Sigma}\Sigma^{-1}) + (\bar{\mathbf{D}} - \mu)^\top \Sigma^{-1} (\bar{\mathbf{D}} - \mu) \right\}$.
- And forget it. It's no great loss.
- Concentrate on the parameters that appear only in the covariance matrix of the observable data.
- Try to identify $\theta = (\beta_1, \Phi_x, \Psi, \Omega_1, \Omega_2)$ from $\Sigma = \text{cov} \left(\begin{array}{c} \mathbf{D}_{i,1} \\ \mathbf{D}_{i,2} \end{array} \right)$.

Stage One: The latent variable model

$$\theta = (\beta_1, \Phi_x, \Psi, \Omega_1, \Omega_2)$$

$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{X}_i + \epsilon_i$, where

- $cov(\mathbf{X}_i) = \Phi_x$
- $cov(\epsilon_i) = \Psi$
- \mathbf{X}_i and ϵ_i are independent.

Vector of “factors” is $\mathbf{F}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix}$.

- Let $\Phi = cov(\mathbf{F}_i)$.
- We know that Φ_x , β_1 and Ψ are functions of Φ .
- We’ve already shown it; this is a regression model.

That’s Stage One. Parameters of the latent variable model are functions of Φ .

Stage Two: The measurement model

$$\mathbf{D}_{i,1} = \boldsymbol{\nu}_1 + \mathbf{F}_i + \mathbf{e}_{i,1}$$

$$\mathbf{D}_{i,2} = \boldsymbol{\nu}_2 + \mathbf{F}_i + \mathbf{e}_{i,2}$$

$cov(\mathbf{e}_{i,1}) = \boldsymbol{\Omega}_1$, $cov(\mathbf{e}_{i,2}) = \boldsymbol{\Omega}_2$. Also, \mathbf{F}_i , $\mathbf{e}_{i,1}$ and $\mathbf{e}_{i,2}$ are independent.

$$\boldsymbol{\Sigma} = cov \begin{pmatrix} \mathbf{D}_{i,1} \\ \mathbf{D}_{i,2} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Phi} + \boldsymbol{\Omega}_1 & \boldsymbol{\Phi} \\ \boldsymbol{\Phi} & \boldsymbol{\Phi} + \boldsymbol{\Omega}_2 \end{pmatrix}$$

$\boldsymbol{\Phi}$, $\boldsymbol{\Omega}_1$ and $\boldsymbol{\Omega}_2$ can easily be recovered from $\boldsymbol{\Sigma}$.

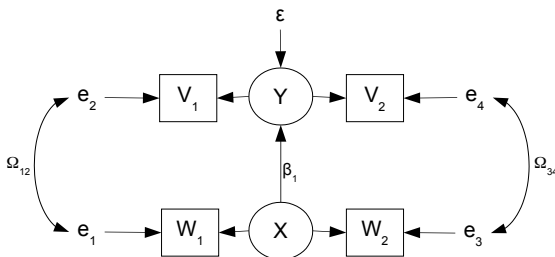
All the parameters in the covariance matrix are identifiable

$$\theta = (\beta_1, \Phi_x, \Psi, \Omega_1, \Omega_2)$$

- Φ_x , β_1 and Ψ are functions of $\Phi = cov(\mathbf{F}_i)$.
- Φ , Ω_1 and Ω_2 are functions of $\Sigma = cov \begin{pmatrix} \mathbf{D}_{i,1} \\ \mathbf{D}_{i,2} \end{pmatrix}$.
- Σ is a function of the probability distribution of the observable data.
- So β_1 , Φ_x , Ψ , Ω_1 , Ω_2 are all functions of the probability distribution of the observable data.
- They are identifiable.

Parameters of the double measurement regression model are identifiable

After re-parameterization



- Correlated measurement error within sets is allowed.
- This is a big plus, because omitted variables are a reality.
- Correlated measurement error between sets must be ruled out by careful data collection.
- No need to do the calculations ever again.

The BMI Health Study

- Body Mass Index: Weight in Kilograms divided by Height in Meters Squared.
- Under 18 means underweight, Over 25 means overweight, Over 30 means obese.
- High BMI is associated with poor health, like high blood pressure and high cholesterol.
- People with high BMI tend to be older and fatter.
- *But*, what if you have a high BMI but are in good physical shape (low percent body fat)?

The Question

- If you control for age and percent body fat, is BMI still associated with indicators for poor health?
- But percent body fat (and to a lesser extent, age) are measured with error. Standard ways of controlling for them with ordinary regression are highly suspect.
- Use the double measurement design.

True variables (all latent)

- $X_1 = \text{Age}$
- $X_2 = \text{BMI}$
- $X_3 = \text{Percent body fat}$
- $Y_1 = \text{Cholesterol}$
- $Y_2 = \text{Diastolic blood pressure}$

Measure twice with different personnel at different locations and by different methods

	Measurement Set One	Measurement Set Two
Age	Self report	Passport or birth certificate
BMI	Dr. Office measurements	Lab technician, no shoes, gown
% Body Fat	Tape and calipers, Dr. Office	Submerge in water tank
Cholesterol	Lab 1	Lab 2
Diastolic BP	Blood pressure cuff, Dr. office	Digital readout, mostly automatic

- Set two is of generally higher quality.
- Correlation of measurement errors is unlikely between sets.

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistical Sciences, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The L^AT_EX source code is available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/2101f19>