

# Double Measurement Regression<sup>1</sup>

STA2053 Fall 2022

---

<sup>1</sup>See last slide for copyright information.

# Overview

- 1 A Small Example
- 2 Computation
- 3 The general model
- 4 Method of Moments
- 5 The BMI study

# Seeking identifiability

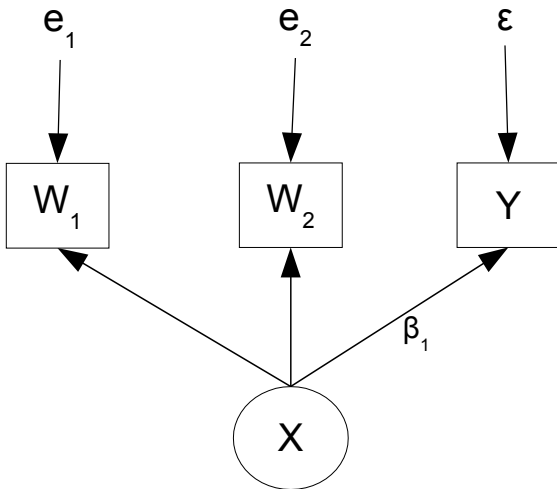
The parameters of this model are not identifiable.

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i \\ W_i &= \nu + X_i + e_i, \end{aligned}$$

- For example,  $X$  might be number of acres planted and  $Y$  might be crop yield.
- Plan the statistical analysis in advance.
- Take 2 independent measurements of the explanatory variable.
- Say, farmer's report and satellite photograph.

# Double measurement

Of the explanatory variable



# Model

Could have written this down based on the path diagram

Independently for  $i = 1, \dots, n$ , let

$$W_{i,1} = \nu_1 + X_i + e_{i,1}$$

$$W_{i,2} = \nu_2 + X_i + e_{i,2}$$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

where

- $X_i$  is normally distributed with mean  $\mu_x$  and variance  $\phi > 0$
- $\epsilon_i$  is normally distributed with mean zero and variance  $\psi > 0$
- $e_{i,1}$  is normally distributed with mean zero and variance  $\omega_1 > 0$
- $e_{i,2}$  is normally distributed with mean zero and variance  $\omega_2 > 0$
- $X_i, e_{i,1}, e_{i,2}$  and  $\epsilon_i$  are all independent.

# Parameter Count Rule

$$W_{i,1} = \nu_1 + X_i + e_{i,1}$$

$$W_{i,2} = \nu_2 + X_i + e_{i,2}$$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

- $\theta = (\nu_1, \nu_2, \beta_0, \mu_x, \beta_1, \phi, \psi, \omega_1, \omega_2)$ : 9 parameters.
- Three expected values, three variances and three covariances: 9 moments.
- Identifiability is possible, but not guaranteed.

# Distribution of the sample data

$$\begin{aligned}W_{i,1} &= \nu_1 + X_i + e_{i,1} \\W_{i,2} &= \nu_2 + X_i + e_{i,2} \\Y_i &= \beta_0 + \beta_1 X_i + \epsilon_i,\end{aligned}$$

The model implies that the triples  $\mathbf{d}_i = (W_{i,1}, W_{i,2}, Y_i)^\top$  are independent multivariate normal with

$$E(\mathbf{d}_i) = E \begin{pmatrix} W_{i,1} \\ W_{i,2} \\ Y_i \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} = \begin{pmatrix} \mu_x + \nu_1 \\ \mu_x + \nu_2 \\ \beta_0 + \beta_1 \mu_x \end{pmatrix},$$

and variance covariance matrix  $\text{cov}(\mathbf{d}_i) = \mathbf{\Sigma} =$

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ & \sigma_{22} & \sigma_{23} \\ & & \sigma_{33} \end{pmatrix} = \begin{pmatrix} \phi + \omega_1 & \phi & \beta_1 \phi \\ & \phi + \omega_2 & \beta_1 \phi \\ & & \beta_1^2 \phi + \psi \end{pmatrix}.$$

# Are the parameters in the covariance matrix identifiable?

Six equations in five unknowns

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ & \sigma_{22} & \sigma_{23} \\ & & \sigma_{33} \end{pmatrix} = \begin{pmatrix} \phi + \omega_1 & \phi & \beta_1 \phi \\ & \phi + \omega_2 & \beta_1 \phi \\ & & \beta_1^2 \phi + \psi \end{pmatrix}.$$

$$\phi = \sigma_{12}$$

$$\omega_1 = \sigma_{11} - \sigma_{12}$$

$$\omega_2 = \sigma_{22} - \sigma_{12}$$

$$\beta_1 = \frac{\sigma_{13}}{\sigma_{12}}$$

$$\psi = \sigma_{33} - \beta_1^2 \phi = \sigma_{33} - \frac{\sigma_{13}^2}{\sigma_{12}}$$

Yes.



## What about the expected values?

Model equations again:

$$W_{i,1} = \nu_1 + X_i + e_{i,1}$$

$$W_{i,2} = \nu_2 + X_i + e_{i,2}$$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

Expected values:

$$\mu_1 = \nu_1 + \mu_x$$

$$\mu_2 = \nu_2 + \mu_x$$

$$\mu_3 = \beta_0 + \beta_1 \mu_x$$

Four parameters appear only in the expected values:  $\nu_1, \nu_2, \mu_x, \beta_0$ .

- Three equations in four unknowns, even with  $\beta_1$  identified from the covariance matrix.
- Parameter count rule applies.

## Re-parameterize

$$\mu_1 = \nu_1 + \mu_x \quad \mu_2 = \nu_2 + \mu_x \quad \mu_3 = \beta_0 + \beta_1 \mu_x$$

- Absorb  $\nu_1, \nu_2, \mu_x, \beta_0$  into  $\boldsymbol{\mu}$ .
- Parameter was  $\boldsymbol{\theta} = (\nu_1, \nu_2, \beta_0, \mu_x, \beta_1, \phi, \psi, \omega_1, \omega_2)$
- Now it's  $\boldsymbol{\theta} = (\mu_1, \mu_2, \mu_3, \beta_1, \phi, \psi, \omega_1, \omega_2)$ .
- Dimension of the parameter space is now one less.
- We haven't lost much, especially because the model was already re-parameterized.

# Re-parameterization

- Re-parameterization makes maximum likelihood possible.
- Otherwise the maximum is not unique and it's a mess.
- Estimate  $\boldsymbol{\mu}$  with  $\bar{\mathbf{d}}$  and it simply disappears from

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = |\boldsymbol{\Sigma}|^{-n/2} (2\pi)^{-np/2} \exp -\frac{n}{2} \left\{ \text{tr}(\hat{\boldsymbol{\Sigma}}\boldsymbol{\Sigma}^{-1}) + (\bar{\mathbf{d}} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{d}} - \boldsymbol{\mu}) \right\}$$

- This step is so common it becomes silent.
- Model equations are often written in centered form.

# Back to the covariance structure equations

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ & \sigma_{22} & \sigma_{23} \\ & & \sigma_{33} \end{pmatrix} = \begin{pmatrix} \phi + \omega_1 & \phi & \beta_1\phi \\ & \phi + \omega_2 & \beta_1\phi \\ & & \beta_1^2\phi + \psi \end{pmatrix}.$$

- Notice that the model dictates  $\sigma_{1,3} = \sigma_{2,3}$ .
- There are two ways to solve for  $\beta_1$ :  
 $\beta_1 = \frac{\sigma_{13}}{\sigma_{12}}$  and  $\beta_1 = \frac{\sigma_{23}}{\sigma_{12}}$ .
- Does this mean the solution for  $\beta_1$  is not “unique?”

# Testing goodness of fit.

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ & \sigma_{22} & \sigma_{23} \\ & & \sigma_{33} \end{pmatrix} = \begin{pmatrix} \phi + \omega_1 & \phi & \beta_1 \phi \\ & \phi + \omega_2 & \beta_1 \phi \\ & & \beta_1^2 \phi + \psi \end{pmatrix}.$$

- $\sigma_{1,3} = \sigma_{2,3}$  is a *model-induced constraint* upon  $\Sigma$ .
- It's a testable null hypothesis.
- If rejected, the model is called into question.
- Likelihood ratio test comparing this model to a completely unrestricted multivariate normal model:

$$G^2 = -2 \log \frac{L(\bar{\mathbf{d}}, \Sigma(\hat{\theta}))}{L(\bar{\mathbf{d}}, \hat{\Sigma})}$$

- Valuable even if the data are not normal.

# The Reproduced Covariance Matrix

- $\Sigma(\hat{\theta})$  is called the *reproduced covariance matrix*.
- It is the covariance matrix of the observable data, written as a function of the model parameters and evaluated at the MLE.

$$\Sigma(\hat{\theta}) = \begin{pmatrix} \hat{\phi} + \hat{\omega}_1 & \hat{\phi} & \hat{\beta}_1 \hat{\phi} \\ & \hat{\phi} + \hat{\omega}_2 & \hat{\beta}_1 \hat{\phi} \\ & & \hat{\beta}_1^2 \hat{\phi} + \hat{\psi} \end{pmatrix}$$

- The reproduced covariance matrix obeys all model-induced constraints, while  $\hat{\Sigma}$  does not.
- But if the model is right they should be close.

$$G^2 = -2 \log \frac{L(\bar{\mathbf{d}}, \Sigma(\hat{\theta}))}{L(\bar{\mathbf{d}}, \hat{\Sigma})}$$

## General pattern for testing goodness of fit

- Suppose there are  $k$  moment structure equations in  $p$  parameters, and all the parameters are identifiable.
- If  $p < k$ , call the parameter vector *over-identifiable*.
- Only need  $p$  equations to solve for  $\theta$ .
- Substituting the solutions (in terms of  $\sigma_{ij}$ ) back into the unused equations would yield  $k - p$  equality constraints on  $\Sigma$ .
- Test those constraints with  $G^2 = -2 \log \frac{L(\bar{\mathbf{d}}, \Sigma(\hat{\theta}))}{L(\bar{\mathbf{d}}, \hat{\Sigma})}$ .
- $df = k - p$
- Don't need to actually derive the constraints – just count them.

## With the same number of equations and parameters

- If the parameter is identifiable, call it *just identifiable*.
- Parameters are 1-1 with those of an unrestricted multivariate normal.
- Call the model “saturated.”
- There are no equality constraints on  $\Sigma$ .
- No likelihood ratio test because  $G^2 = -2 \log \frac{L(\bar{\mathbf{d}}, \Sigma(\hat{\boldsymbol{\theta}}))}{L(\bar{\mathbf{d}}, \hat{\Sigma})} = 0$ .
- This is what happens in regression with all observed variables.



# Data analysis strategy

- Verify identifiability.
- If the model is over-identified, test goodness of fit.
- If it passes (non-significant), proceed.
- Now think of your model as the “full,” or unrestricted model.
- Compared to some (even more) reduced model that is restricted by a null hypothesis like  $\beta_1 = 0$ .
- Fit the reduced model.
- Subtract goodness of fit ( $G^2$  or “chi-square”) statistics to test  $H_0$ .

# Subtract goodness of fit statistics

$G^2$  tests the full model against the saturated model, and  $G_0^2$  tests the restricted model against the saturated model.

$$\begin{aligned}G_0^2 - G^2 &= -2 \log \frac{L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}_0))}{L(\bar{\mathbf{d}}, \hat{\boldsymbol{\Sigma}})} - -2 \log \frac{L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}))}{L(\bar{\mathbf{d}}, \hat{\boldsymbol{\Sigma}})} \\&= -2 \left( \log L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}_0)) - \log L(\bar{\mathbf{d}}, \hat{\boldsymbol{\Sigma}}) - \log L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}})) \right. \\&\quad \left. + \log L(\bar{\mathbf{d}}, \hat{\boldsymbol{\Sigma}}) \right) \\&= -2 \log \frac{L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}_0))}{L(\bar{\mathbf{d}}, \boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}))}\end{aligned}$$

## Further comments

- Models with non-identifiable parameters can imply testable equality constraints, but testing them is not automatic.
- Models can imply *inequality* constraints on  $\Sigma$ , too.
- Using the solutions

$$\phi = \sigma_{12}$$

$$\omega_1 = \sigma_{11} - \sigma_{12}$$

$$\omega_2 = \sigma_{22} - \sigma_{12}$$

$$\beta_1 = \frac{\sigma_{13}}{\sigma_{12}}$$

$$\psi = \sigma_{33} - \beta_1^2 \phi = \sigma_{33} - \frac{\sigma_{13}^2}{\sigma_{12}}$$

We get four inequality constraints.

Four inequality constraints on  $\Sigma$ 

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ & \sigma_{22} & \sigma_{23} \\ & & \sigma_{33} \end{pmatrix} = \begin{pmatrix} \phi + \omega_1 & \phi & \beta_1 \phi \\ & \phi + \omega_2 & \beta_1 \phi \\ & & \beta_1^2 \phi + \psi \end{pmatrix}.$$

$$\phi = \sigma_{12} > 0$$

$$\omega_1 = \sigma_{11} - \sigma_{12} > 0$$

$$\omega_2 = \sigma_{22} - \sigma_{12} > 0$$

$$\psi = \sigma_{33} - \frac{\sigma_{13}^2}{\sigma_{12}} > 0$$

# Inequality constraints

- Inequality constraints arise because variances are positive.
- Or more generally, covariance matrices are positive definite.
- Could inequality constraints be violated in numerical maximum likelihood?
- Definitely.
- But only a little by sampling error if the model is correct.
- So maybe it's not so dumb to test hypotheses like  $H_0 : \omega_1 = 0$ .
- Since the model says  $\omega_1 = \sigma_{11} - \sigma_{12}$  and it might not be true.

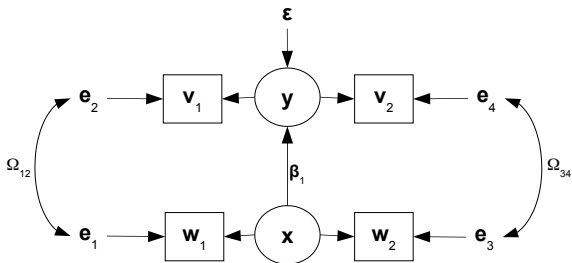
# Computation with lavaan

2053Babydouble.pdf

This link will probably be broken once the term is over. See the course website for another route to the output file:

<http://www.utstat.toronto.edu/brunner/oldclass/2053f22>

# The general double measurement design



These are all matrices.

Double measurement can help solve a big problem: Correlated measurement error.

- The main idea is that  $x$  and  $y$  are each measured twice, perhaps at different times using different methods.
- Measurement errors 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.



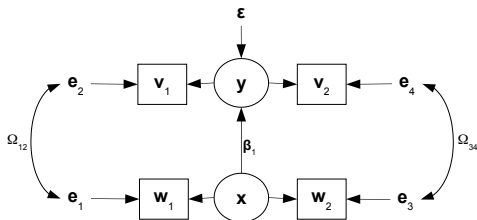
# Measurement errors may be correlated

Look at the measurement model

$$\mathbf{F}_i = \begin{pmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{pmatrix}$$

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

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



$$\text{cov}(\mathbf{e}_{i,1}) = \boldsymbol{\Omega}_1 = \begin{pmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}^\top & \Omega_{22} \end{pmatrix}$$

$$\text{cov}(\mathbf{e}_{i,2}) = \boldsymbol{\Omega}_2 = \begin{pmatrix} \Omega_{33} & \Omega_{34} \\ \Omega_{34}^\top & \Omega_{44} \end{pmatrix}$$

# 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  $\nu_1$ ,  $\nu_2$ ,  $\beta_0$  and  $\mu_x$  separately.
- Swallow them into  $\mu$ .
- Estimate  $\mu$  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  $\epsilon_i$  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  $\Phi_x$ ,  $\beta_1$  and  $\Psi$  are functions of  $\Phi$ .
- We’ve already shown it; this is a regression model.

That’s Stage One. Parameters of the latent variable model are functions of  $\Phi$ .

## 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.

$$cov \left( \begin{array}{c} \mathbf{d}_{i,1} \\ \mathbf{d}_{i,2} \end{array} \right) = \boldsymbol{\Sigma} = \left( \begin{array}{c|c} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \hline \boldsymbol{\Sigma}_{12}^\top & \boldsymbol{\Sigma}_{22} \end{array} \right) = \left( \begin{array}{c|c} \boldsymbol{\Phi} + \boldsymbol{\Omega}_1 & \boldsymbol{\Phi} \\ \hline \boldsymbol{\Phi} & \boldsymbol{\Phi} + \boldsymbol{\Omega}_2 \end{array} \right)$$

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

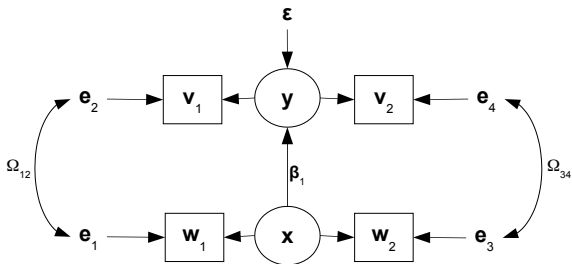
All the parameters in the covariance matrix are identifiable

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

- $\Phi_x$ ,  $\beta_1$  and  $\Psi$  are functions of  $\Phi = cov(\mathbf{F}_i)$ .
- $\Phi$ ,  $\Omega_1$  and  $\Omega_2$  are functions of  $\Sigma = cov \begin{pmatrix} \mathbf{d}_{i,1} \\ \mathbf{d}_{i,2} \end{pmatrix}$ .
- $\Sigma$  is a function of the probability distribution of the observable data.
- So  $\beta_1$ ,  $\Phi_x$ ,  $\Psi$ ,  $\Omega_1$ ,  $\Omega_2$  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



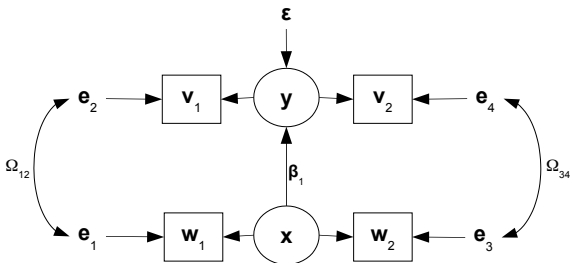
- Correlated measurement error within sets is allowed.
- This is a big plus, because omitted variables are a reality.
- Correlated measurement error between sets must be ruled out by careful data collection.
- No need to do the calculations ever again.

# Method of Moments

- What if the distributions are not normal?
- What's the parameter vector?  
 $\theta = (\beta, \Phi, \Psi, \Omega, F_x, F_\epsilon, F_e)$ .
- We are only interested in  $\beta$  anyway.
- Put hats on solution to covariance structure equations?



## Path diagram again



## Covariance structure equations

$$\begin{aligned}
 \Sigma &= \text{cov} \begin{pmatrix} \mathbf{w}_{i,1} \\ \mathbf{v}_{i,1} \\ \mathbf{w}_{i,2} \\ \mathbf{v}_{i,2} \end{pmatrix} \\
 &= \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} & \Sigma_{14} \\ & \Sigma_{22} & \Sigma_{23} & \Sigma_{24} \\ & & \Sigma_{33} & \Sigma_{34} \\ & & & \Sigma_{44} \end{pmatrix} \\
 &= \begin{pmatrix} \Phi + \Omega_{11} & \Phi\beta^\top + \Omega_{12} & \Phi & \Phi\beta^\top \\ & \beta\Phi\beta^\top + \Psi + \Omega_{22} & \beta\Phi & \beta\Phi\beta^\top + \Psi \\ & & \Phi + \Omega_{33} & \Phi\beta^\top + \Omega_{34} \\ & & & \beta\Phi\beta^\top + \Psi + \Omega_{44} \end{pmatrix}
 \end{aligned}$$

# The interesting equations

$$\text{cov}(\mathbf{w}_{i,1}, \mathbf{w}_{i,2}) = \Sigma_{13} = \Phi$$

$$\text{cov}(\mathbf{v}_{i,1}, \mathbf{w}_{i,2}) = \Sigma_{23} = \beta\Phi$$

$$\text{cov}(\mathbf{w}_{i,1}, \mathbf{v}_{i,2}) = \Sigma_{14} = \Phi\beta^\top$$

## Solutions

$$\Phi = \Sigma_{13}$$

$$\beta = \Sigma_{23}\Phi^{-1} = \Sigma_{14}^\top\Phi^{-1}$$

# MOM estimates

Using the solutions

$$\begin{aligned}\Phi &= \Sigma_{13} = \text{cov}(\mathbf{w}_{i,1}, \mathbf{w}_{i,2}) \\ \beta &= \Sigma_{23}\Phi^{-1} = \Sigma_{14}^\top \Phi^{-1}\end{aligned}$$

$$\hat{\Phi}_M = \frac{1}{2}(\hat{\Sigma}_{13} + \hat{\Sigma}_{13}^\top)$$

$$\hat{\beta}_M = \frac{1}{2}(\hat{\Sigma}_{23} + \hat{\Sigma}_{14}^\top) \hat{\Phi}_M^{-1}$$

Do you agree that the asymptotic distribution of  $\hat{\beta}_M$  is multivariate normal?

# Asymptotic distribution of $\hat{\beta}_M$

- $\hat{\beta}_M$  is approximately multivariate normal with expected value  $\beta$
- And covariance matrix ...

Bootstrap.

# 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?
- 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<sup>A</sup>T<sub>E</sub>X source code is available from the course website:

<http://www.utstat.toronto.edu/brunner/oldclass/2053f22>