

The General Structural Equation Model¹

STA2101 F19

¹See last slide for copyright information.

Features of Structural Equation Models

- Multiple equations.
- All the variables are random.
- An explanatory variable in one equation can be the response variable in another equation.
- Models are represented by path diagrams.
- Identifiability is always an issue.
- The statistical models are explicitly models of influence. They are often called *causal models*.

Correlation versus Causation

- The path diagrams deliberately imply influence. If $A \rightarrow B$, we are saying A *contributes* to B , or partly *causes* it.
- Data are usually observational. The correlation-causation issue does not go away.
- You may be able to argue on theoretical grounds that $A \rightarrow B$ is more believable than $B \rightarrow A$.
- A statistical model cannot “prove” causality, but if you have a causal model, you may be able to test whether it’s compatible with the data.

Vocabulary

- **Exogenous variable:** In the regression-like equations of a structural equation model, the exogenous variables are ones that appear *only* on the right side of the equals sign, and never on the left side in any equation. If you think of Y being a function of X , this is one way to remember the meaning of **exogenous**.
- **Endogenous variable:** Endogenous variables are those that appear on the left side of at least one equals sign. Endogenous variables depend on the exogenous variables, and possibly other endogenous variables. Think of an arrow from an exogenous variable to an endogenous variable. The **end** of the arrow is pointing at the **endogenous** variable.
- **Factor:** This term has a meaning that actually conflicts with its meaning in mainstream Statistics, particularly in experimental design. A *factor* is an underlying trait or characteristic that cannot be measured directly, like intelligence. It is a latent variable, period.

$$Y_{i,1} = \alpha_1 + \gamma_1 X_{i,1} + \gamma_2 X_{i,2} + \epsilon_{i,1}$$

$$Y_{i,2} = \alpha_2 + \beta Y_{i,1} + \epsilon_{i,2}$$

- Regression coefficients (links between exogenous variables and endogenous variables) are now called gamma instead of beta.
- Betas are used for links between endogenous variables.
- Intercepts are alphas but they will soon disappear.

Losing the intercepts and expected values

- Mostly the intercepts and expected values are not identifiable anyway, as in multiple regression with measurement error.
- We have a chance to identify a *function* of the parameter vector – the parameters that appear in the covariance matrix Σ of an observable data vector. $\Sigma = cov(\mathbf{D}_i)$.
- Denote the vector of parameters that appear in Σ by θ .
- Re-parameterize. The new parameter vector is (θ, μ) , where $\mu = E(\mathbf{D}_i)$.
- Estimate μ with $\bar{\mathbf{D}}$, forget it, and concentrate on θ .
- From this point on the models *seemingly* have zero means, and no intercepts.

A General Two-Stage Model

Stage 1 is the latent variable model and Stage 2 is the measurement model.

$$\mathbf{Y}_i = \beta \mathbf{Y}_i + \Gamma \mathbf{X}_i + \epsilon_i$$

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

$$\mathbf{D}_i = \Lambda \mathbf{F}_i + \mathbf{e}_i$$

- \mathbf{D}_i (the data) are observable. All other variables are latent.
- $\mathbf{Y}_i = \beta \mathbf{Y}_i + \Gamma \mathbf{X}_i + \epsilon_i$ is called the *Latent Variable Model*.
- The latent vectors \mathbf{X}_i and \mathbf{Y}_i are collected into a *factor* \mathbf{F}_i . This is *not* a categorical explanatory variable, the usual meaning of “factor” in experimental design.
- $\mathbf{D}_i = \Lambda \mathbf{F}_i + \mathbf{e}_i$ is called the *Measurement Model*.

$$\mathbf{Y}_i = \boldsymbol{\beta}\mathbf{Y}_i + \boldsymbol{\Gamma}\mathbf{X}_i + \boldsymbol{\epsilon}_i \quad \mathbf{F}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix} \quad \mathbf{D}_i = \boldsymbol{\Lambda}\mathbf{F}_i + \mathbf{e}_i$$

- \mathbf{Y}_i is a $q \times 1$ random vector.
- $\boldsymbol{\beta}$ is a $q \times q$ matrix of constants with zeros on the main diagonal.
- \mathbf{X}_i is a $p \times 1$ random vector.
- $\boldsymbol{\Gamma}$ is a $q \times p$ matrix of constants.
- $\boldsymbol{\epsilon}_i$ is a $q \times 1$ random vector.
- \mathbf{F}_i (F for Factor) is just \mathbf{X}_i stacked on top of \mathbf{Y}_i . It is a $(p + q) \times 1$ random vector.
- \mathbf{D}_i is a $k \times 1$ random vector. Sometimes, $\mathbf{D}_i = \begin{pmatrix} \mathbf{W}_i \\ \mathbf{V}_i \end{pmatrix}$.
- $\boldsymbol{\Lambda}$ is a $k \times (p + q)$ matrix of constants: “factor loadings.”
- \mathbf{e}_i is a $k \times 1$ random vector.
- \mathbf{X}_i , $\boldsymbol{\epsilon}_i$ and \mathbf{e}_i are independent.

Covariance matrices

All assumed positive definite unless otherwise specified

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

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

$$\mathbf{D}_i = \mathbf{\Lambda} \mathbf{F}_i + \mathbf{e}_i$$

$$\text{cov}(\mathbf{X}_i) = \mathbf{\Phi}_x$$

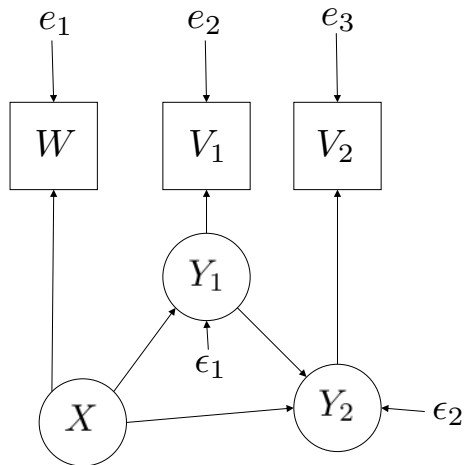
$$\text{cov}(\boldsymbol{\epsilon}_i) = \mathbf{\Psi}$$

$$\text{cov}(\mathbf{F}_i) = \mathbf{\Phi} = \begin{pmatrix} \text{cov}(\mathbf{X}_i) & \text{cov}(\mathbf{X}_i, \mathbf{Y}_i) \\ \text{cov}(\mathbf{Y}_i, \mathbf{X}_i) & \text{cov}(\mathbf{Y}_i) \end{pmatrix} = \begin{pmatrix} \mathbf{\Phi}_{11} & \mathbf{\Phi}_{12} \\ \mathbf{\Phi}_{12}^\top & \mathbf{\Phi}_{22} \end{pmatrix}$$

$$\text{cov}(\mathbf{e}_i) = \mathbf{\Omega}$$

$$\text{cov}(\mathbf{D}_i) = \mathbf{\Sigma}$$

Example: A Path Model with Measurement Error



$$Y_{i,1} = \gamma_1 X_i + \epsilon_{i,1}$$

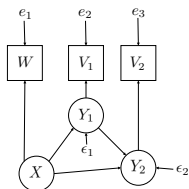
$$Y_{i,2} = \beta Y_{i,1} + \gamma_2 X_i + \epsilon_{i,2}$$

$$W_i = X_i + e_{i,1}$$

$$V_{i,1} = Y_{i,1} + e_{i,2}$$

$$V_{i,2} = Y_{i,2} + e_{i,3}$$

Matrix Form



$$Y_{i,1} = \gamma_1 X_i + \epsilon_{i,1}$$

$$Y_{i,2} = \beta Y_{i,1} + \gamma_2 X_i + \epsilon_{i,2}$$

$$W_i = X_i + e_{i,1}$$

$$V_{i,1} = Y_{i,1} + e_{i,2}$$

$$V_{i,2} = Y_{i,2} + e_{i,3}$$

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

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

$$\mathbf{D}_i = \boldsymbol{\Lambda} \mathbf{F}_i + \mathbf{e}_i$$

$$\begin{pmatrix} \mathbf{Y}_i \\ Y_{i,1} \\ Y_{i,2} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\beta} \\ 0 & 0 \\ \beta & 0 \end{pmatrix} \begin{pmatrix} \mathbf{Y}_i \\ Y_{i,1} \\ Y_{i,2} \end{pmatrix} + \begin{pmatrix} \boldsymbol{\Gamma} \\ \gamma_1 \\ \gamma_2 \end{pmatrix} \mathbf{X}_i + \begin{pmatrix} \boldsymbol{\epsilon}_i \\ \epsilon_{i,1} \\ \epsilon_{i,2} \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{D}_i \\ W_i \\ V_{i,1} \\ V_{i,2} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Lambda} \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{F}_i \\ X_i \\ Y_{i,1} \\ Y_{i,2} \end{pmatrix} + \begin{pmatrix} \mathbf{e}_i \\ e_{i,1} \\ e_{i,2} \\ e_{i,3} \end{pmatrix}$$

Observable variables in the “latent” variable model

$$\mathbf{Y}_i = \beta \mathbf{Y}_i + \Gamma \mathbf{X}_i + \epsilon_i$$

Fairly common

- These present no problem.
- Let $P(e_j = 0) = 1$, so $Var(e_j) = 0$.
- And $Cov(e_i, e_j) = 0$
- Because if $P(e_j = 0) = 1$,

$$\begin{aligned}Cov(e_i, e_j) &= E(e_i e_j) - E(e_i)E(e_j) \\ &= E(e_i \cdot 0) - E(e_i) \cdot 0 \\ &= 0 - 0 = 0\end{aligned}$$

- In $\Omega = cov(\mathbf{e}_i)$, column j (and row j) are all zeros.
- Ω singular, no problem.

What should you be able to do?

- Given a path diagram, write the model equations and say which exogenous variables are correlated with each other.
- Given the model equations and information about which exogenous variables are correlated with each other, draw the path diagram.
- Given either piece of information, write the model in matrix form and say what all the matrices are.
- Calculate model covariance matrices.
- Check identifiability.

Recall the notation

$$\mathbf{Y}_i = \beta \mathbf{Y}_i + \Gamma \mathbf{X}_i + \epsilon_i$$

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

$$\mathbf{D}_i = \Lambda \mathbf{F}_i + \mathbf{e}_i$$

$$\text{cov}(\mathbf{X}_i) = \Phi_x$$

$$\text{cov}(\epsilon_i) = \Psi$$

$$\text{cov}(\mathbf{F}_i) = \Phi = \begin{pmatrix} \text{cov}(\mathbf{X}_i) & \text{cov}(\mathbf{X}_i, \mathbf{Y}_i) \\ \text{cov}(\mathbf{Y}_i, \mathbf{X}_i) & \text{cov}(\mathbf{Y}_i) \end{pmatrix} = \begin{pmatrix} \Phi_{11} & \Phi_{12} \\ \Phi_{12}^\top & \Phi_{22} \end{pmatrix}$$

$$\text{cov}(\mathbf{e}_i) = \Omega$$

$$\text{cov}(\mathbf{D}_i) = \Sigma$$

Calculate a general expression for $\Sigma(\theta)$.

For the latent variable model, calculate $\Phi = cov(\mathbf{F}_i)$

Have $cov(\mathbf{X}_i) = \Phi_x$, need $cov(\mathbf{Y}_i)$ and $cov(\mathbf{X}_i, \mathbf{Y}_i)$

$$\mathbf{Y}_i = \beta \mathbf{Y}_i + \Gamma \mathbf{X}_i + \epsilon_i$$

$$\Rightarrow \mathbf{Y}_i - \beta \mathbf{Y}_i = \Gamma \mathbf{X}_i + \epsilon_i$$

$$\Rightarrow \mathbf{I} \mathbf{Y}_i - \beta \mathbf{Y}_i = \Gamma \mathbf{X}_i + \epsilon_i$$

$$\Rightarrow (\mathbf{I} - \beta) \mathbf{Y}_i = \Gamma \mathbf{X}_i + \epsilon_i$$

$$\Rightarrow (\mathbf{I} - \beta)^{-1} (\mathbf{I} - \beta) \mathbf{Y}_i = (\mathbf{I} - \beta)^{-1} (\Gamma \mathbf{X}_i + \epsilon_i)$$

$$\Rightarrow \mathbf{Y}_i = (\mathbf{I} - \beta)^{-1} (\Gamma \mathbf{X}_i + \epsilon_i)$$

So,

$$\begin{aligned} cov(\mathbf{Y}_i) &= (\mathbf{I} - \beta)^{-1} cov(\Gamma \mathbf{X}_i + \epsilon_i) (\mathbf{I} - \beta)^{-1 \top} \\ &= (\mathbf{I} - \beta)^{-1} (cov(\Gamma \mathbf{X}_i) + cov(\epsilon_i)) (\mathbf{I} - \beta^{\top})^{-1} \\ &= (\mathbf{I} - \beta)^{-1} (\Gamma \Phi_x \Gamma^{\top} + \Psi) (\mathbf{I} - \beta^{\top})^{-1} \end{aligned}$$

Does $(\mathbf{I} - \beta)^{-1}$ exist?

$$\begin{aligned}(\mathbf{I} - \beta)\mathbf{Y}_i &= \Gamma\mathbf{X}_i + \epsilon_i \\ \Rightarrow \text{cov}((\mathbf{I} - \beta)\mathbf{Y}_i) &= \text{cov}(\Gamma\mathbf{X}_i + \epsilon_i) \\ \Rightarrow (\mathbf{I} - \beta)\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \beta)^\top &= \Gamma\text{cov}(\mathbf{X}_i)\Gamma^\top + \text{cov}(\epsilon_i) \\ \Rightarrow (\mathbf{I} - \beta)\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \beta)^\top &= \Gamma\Phi_x\Gamma^\top + \Psi\end{aligned}$$

Now let the $q \times 1$ constant vector $\mathbf{a} \neq \mathbf{0}$, and

$$\begin{aligned}\mathbf{a}^\top(\mathbf{I} - \beta)\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \beta)^\top\mathbf{a} &= \mathbf{a}^\top\Gamma\Phi_x\Gamma^\top\mathbf{a} + \mathbf{a}^\top\Psi\mathbf{a} \\ &> 0\end{aligned}$$

Because $\text{cov}(\Gamma\mathbf{X}_i) = \Gamma\Phi_x\Gamma^\top$ is non-negative definite and $\Psi = \text{cov}(\epsilon_i)$ is positive definite. Hence, $(\mathbf{I} - \beta)\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \beta)^\top$ is positive definite.

Have $(\mathbf{I} - \boldsymbol{\beta})\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \boldsymbol{\beta})^\top$ positive definite

- So the $q \times q$ matrix $(\mathbf{I} - \boldsymbol{\beta})\text{cov}(\mathbf{Y}_i)(\mathbf{I} - \boldsymbol{\beta})^\top$ is full rank: rank is q .
- $(\mathbf{I} - \boldsymbol{\beta})$, $\text{cov}(\mathbf{Y}_i)$ and $(\mathbf{I} - \boldsymbol{\beta})^\top$ are all $q \times q$.
- Rank of a product is minimum rank.
- Hence $\text{Rank}(\mathbf{I} - \boldsymbol{\beta}) = q$, and the inverse exists. ■

$(\mathbf{I} - \boldsymbol{\beta})^{-1}$ exists if the rest of the model is correct

- This forms a strange and unexpected hole in the parameter space.
- No need to “assume” it, as all the textbooks do.
- For example, if $\boldsymbol{\beta} = \mathbf{I}$, then

$$\begin{aligned}\mathbf{Y}_i &= \boldsymbol{\beta}\mathbf{Y}_i + \boldsymbol{\Gamma}\mathbf{X}_i + \boldsymbol{\epsilon}_i \\ \Rightarrow \mathbf{Y}_i &= \mathbf{Y}_i + \boldsymbol{\Gamma}\mathbf{X}_i + \boldsymbol{\epsilon}_i \\ \Rightarrow \boldsymbol{\epsilon}_i &= -\boldsymbol{\Gamma}\mathbf{X}_i\end{aligned}$$

Impossible if \mathbf{X}_i and $\boldsymbol{\epsilon}_i$ are independent.

- Summary: The rest of the model places subtle restrictions on $\boldsymbol{\beta}$.

But we were in the middle of a covariance calculation.

For the measurement model, calculate $\Sigma = cov(\mathbf{D}_i)$

$$\begin{aligned}\mathbf{D}_i &= \mathbf{\Lambda}\mathbf{F}_i + \mathbf{e}_i \\ \Rightarrow cov(\mathbf{D}_i) &= cov(\mathbf{\Lambda}\mathbf{F}_i + \mathbf{e}_i) \\ &= cov(\mathbf{\Lambda}\mathbf{F}_i) + cov(\mathbf{e}_i) \\ &= \mathbf{\Lambda}cov(\mathbf{F}_i)\mathbf{\Lambda}^\top + cov(\mathbf{e}_i) \\ &= \mathbf{\Lambda}\mathbf{\Phi}\mathbf{\Lambda}^\top + \mathbf{\Omega} \\ &= \mathbf{\Sigma}\end{aligned}$$

Two-stage Proofs of Identifiability

Stage 1 is the latent variable model and Stage 2 is the measurement model.

- Show the parameters of the latent variable model $(\beta, \Gamma, \Phi_x, \Psi)$ can be recovered from $\Phi = cov(\mathbf{F}_i)$.
- Show the parameters of the measurement model (Λ, Φ, Ω) can be recovered from $\Sigma = cov(\mathbf{D}_i)$.
- This means all the parameters can be recovered from Σ .
- Break a big problem into two smaller ones.
- Develop *rules* for checking identifiability at each stage.
- Just look at the path diagram.

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>