

# Structural Equation Models: The General Case

STA431: Spring 2015

See last slide for copyright information

# An Extension of Multiple Regression

- More than one regression-like equation
- Includes latent variables
- Variables can be explanatory in one equation and response in another
- Modest changes in notation
- Vocabulary
- Path diagrams
- No intercepts, all expected values zero
- Serious modeling (compared to ordinary statistical models)
- Parameter identifiability

# Variables can be response in one equation and explanatory in another

- Variables (IQ = Intelligence Quotient):

- $X_1$  = Mother's adult IQ

- $X_2$  = Father's adult IQ

- $Y_1$  = Person's adult IQ

- $Y_2$  = Child's IQ in Grade 8

$$Y_1 = \alpha_1 + \gamma_1 X_1 + \gamma_2 X_2 + \epsilon_1$$

$$Y_2 = \alpha_2 + \beta Y_1 + \epsilon_2$$

- Of course all these variables are measured with error.
- We will lose the intercepts very soon.

# Modest changes in notation

- Regression coefficients are now called gamma instead of beta
- Betas are used for links between Y variables
- Intercepts are alphas but they will soon disappear.
- We feel free to drop the subscript  $i$ ; implicitly, everything is independent and identically distributed for  $i = 1, \dots, n$ .

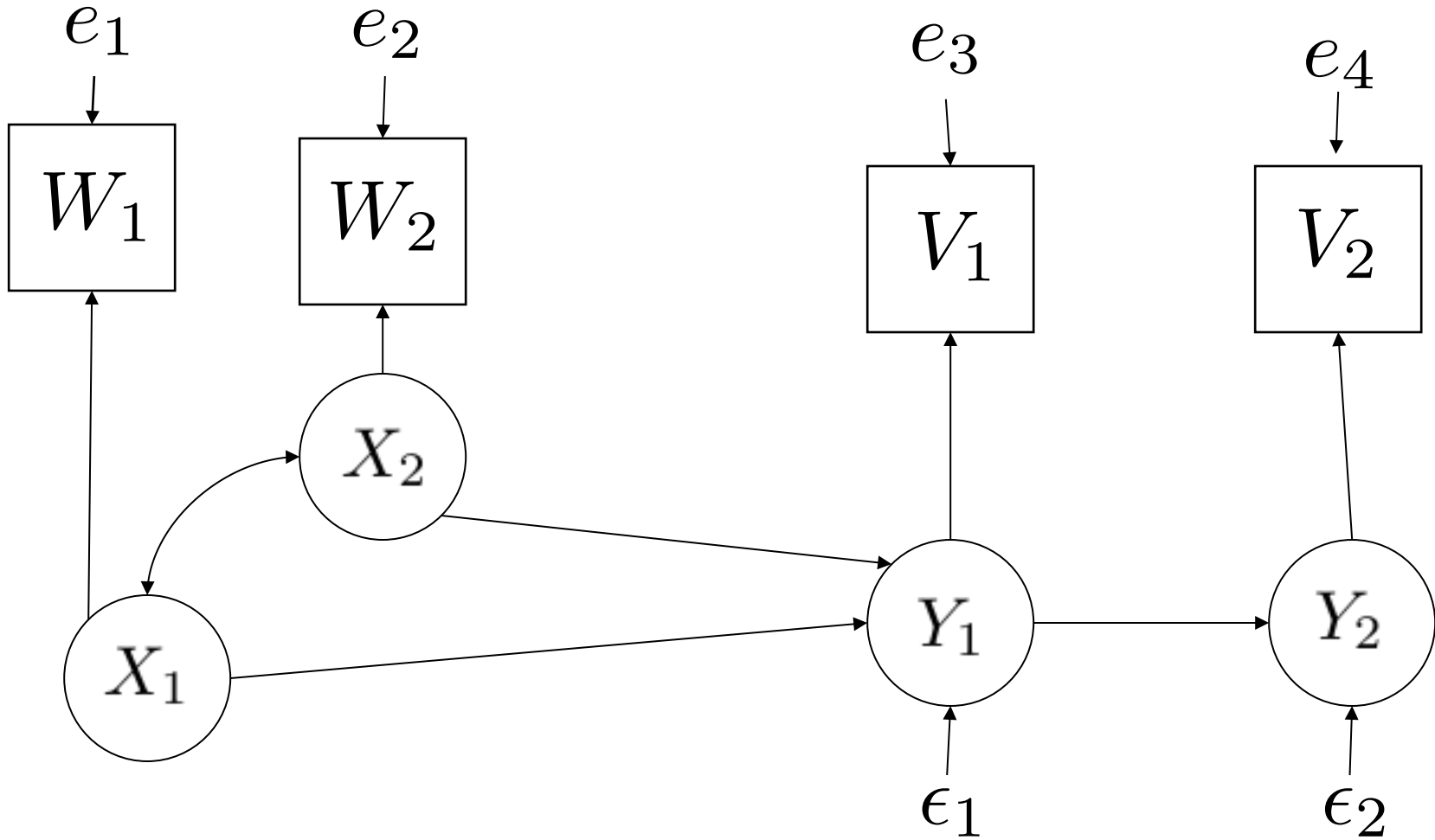
$$Y_1 = \alpha_1 + \gamma_1 X_1 + \gamma_2 X_2 + \epsilon_1$$

$$Y_2 = \alpha_2 + \beta Y_1 + \epsilon_2$$

# Vocabulary

- Variables can be Latent or Manifest.
  - Manifest means observable
  - All error terms are latent
- Variables can be Exogenous or Endogenous
  - **Exogenous** variables appear only on the right side of the = sign.
    - Think “X” for explanatory variable.
    - All error terms are exogenous
  - **Endogenous** variables appear on the left of at least one = sign.
    - Think “end” of an arrow pointing from exogenous to endogenous
    - Betas link endogenous variables to other endogenous variables.

# Path diagrams



# Path Diagram Rules

- Latent variables are enclosed by ovals.
- Observable (manifest) variables are enclosed by rectangles.
- Error terms are not enclosed
  - Sometimes the arrows from the error terms seem to come from nowhere. The symbol for the error term does not appear in the path diagram.
  - Sometimes there are no arrows for the error terms at all. It is just assumed that such an arrow points to each endogenous variable.
- Straight, single-headed arrows point from each variable on the right side of an equation to the endogenous variable on the left side.
  - Sometimes the coefficient is written on the arrow, but sometimes it is not.
- A curved, double-headed arrow between two variables (always exogenous variables) means they have a non-zero covariance.
  - Sometimes the symbol for the covariance is written on the curved arrow, but sometimes it is not.

# Causal Modeling (cause and effect)

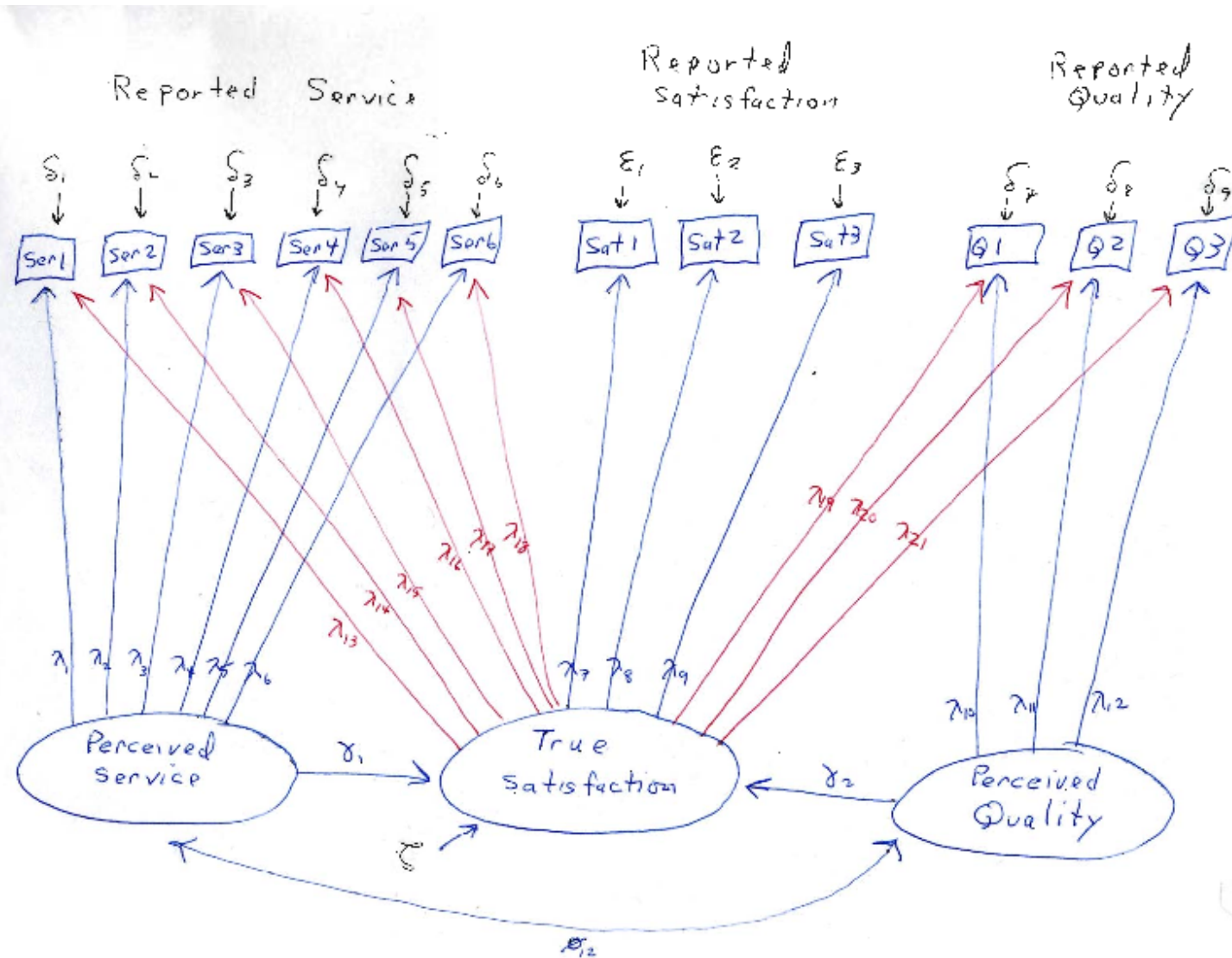
- The arrows deliberately imply that if  $A \rightarrow B$ , we are saying *A contributes* to B, or partly *causes* it.
- There may be other contributing variables. All the ones that are unknown are lumped together in the error term.
- Are these unknown variables are independent of the variables in the model? Probably not.
- Sometimes we can get around the problem with instrumental variables.



# Serious Modeling

- Once you accept that model equations are statements about what contributes to what, you realize that structural equation models represent a rough *theory* of the data, with some parts (the parameter values) unknown.
- They are somewhere between ordinary statistical models, which are like one-size-fits-all clothing, and true scientific models, which are like tailor made clothing.
- So they are very flexible and potentially valuable. It is *good* to combine what the data can tell you with what you already know.
- But structural equation models can require a lot of input and careful thought to construct. In this course, we will get by mostly on common sense.
- In general, the parameters of the most reasonable model need not be identifiable. It depends upon the form of the data as well as on the model. Identifiability needs to be checked. Frequently, this can be done by inspection.

# Example: Halo Effects in Real Estate



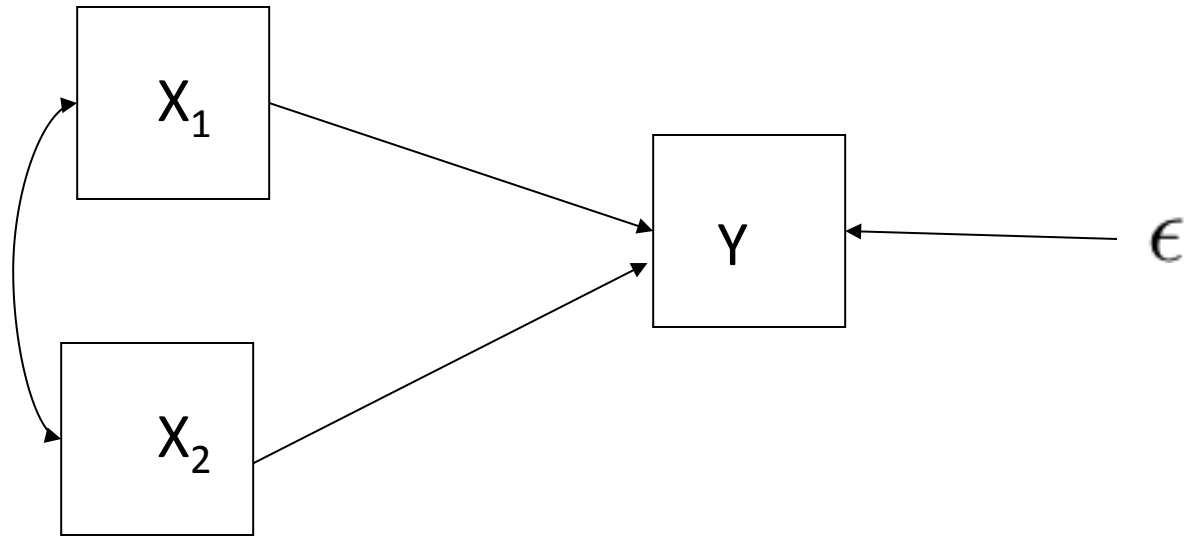
# Losing the intercepts and expected values

- 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  $\Sigma = V(\mathbf{D})$ .
- Re-parameterize. The new parameter vector is the set of parameters in  $\Sigma$ , and also  $\boldsymbol{\mu} = E(\mathbf{D})$ . Estimate  $\boldsymbol{\mu}$  with  $\bar{x}$ , forget it, and concentrate on inference for the parameters in  $\Sigma$ .
- To make calculation of the covariance matrix easier, write the model equations in centered form. The little letters  $c$  over the variables are invisible.

From this point on the models have  
no means and no intercepts.

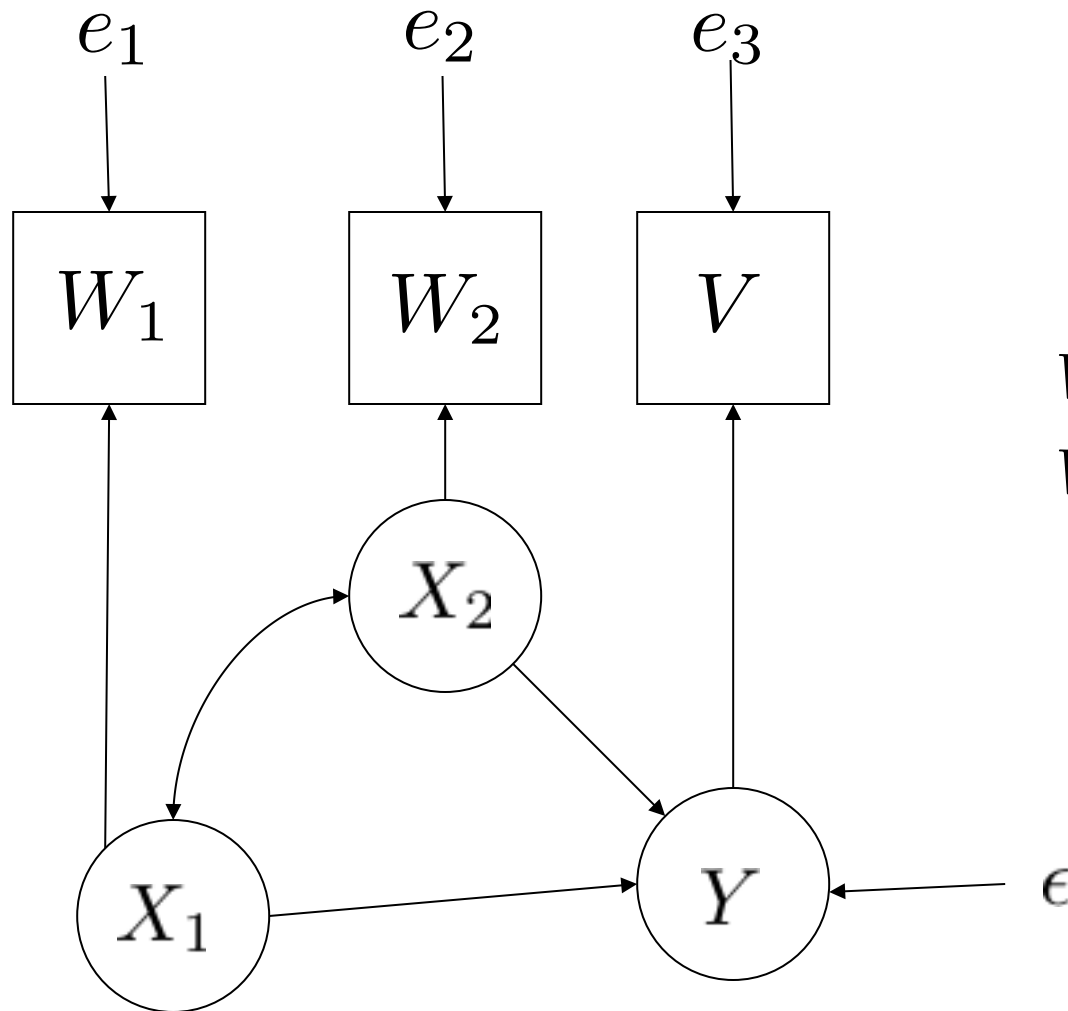
Now more examples

# Multiple Regression



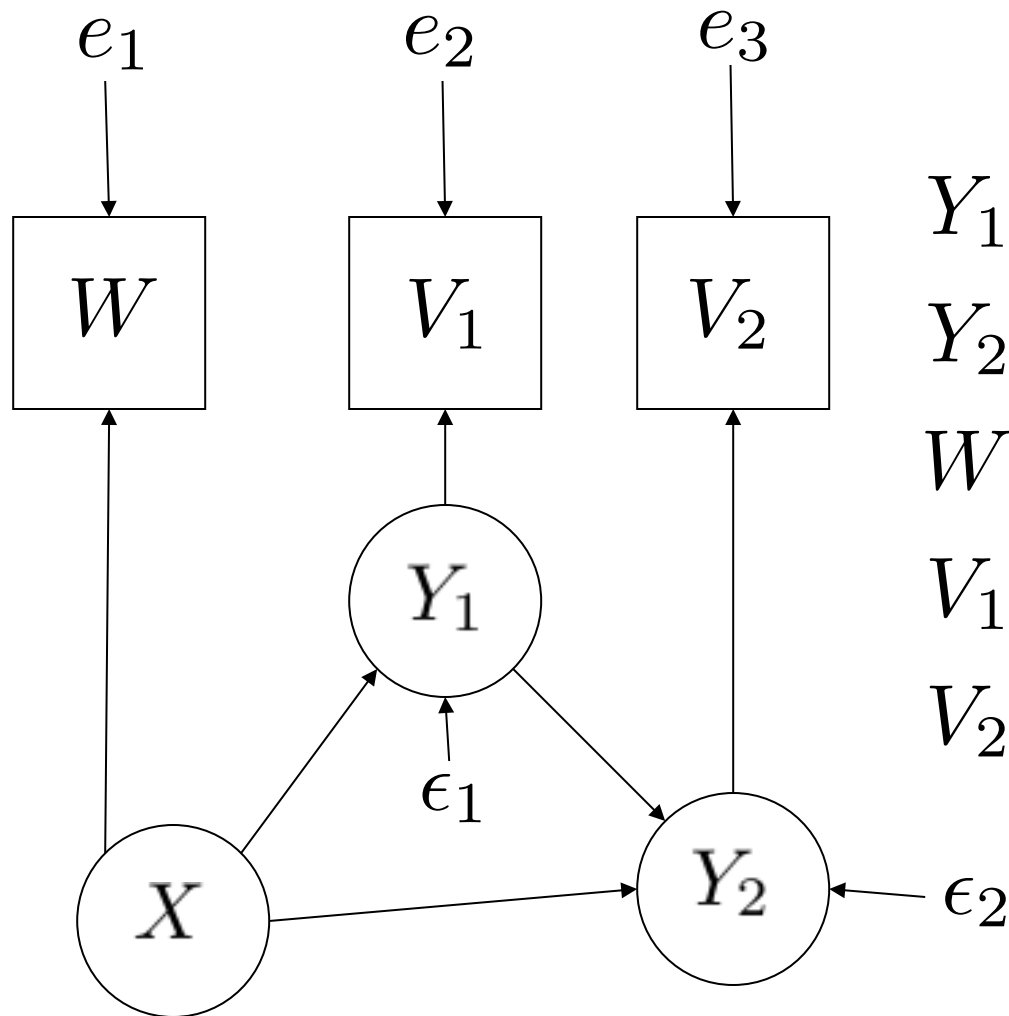
$$Y = \gamma_1 X_1 + \gamma_2 X_2 + \epsilon$$

# Regression with measurement error



$$\begin{aligned} Y &= \gamma_1 X_1 + \gamma_2 X_2 + \epsilon \\ W_1 &= X_1 + e_1 \\ W_2 &= X_2 + e_2 \\ V &= Y + e_3 \end{aligned}$$

# A Path Model with Measurement Error



$$Y_1 = \gamma_1 X + \epsilon_1$$

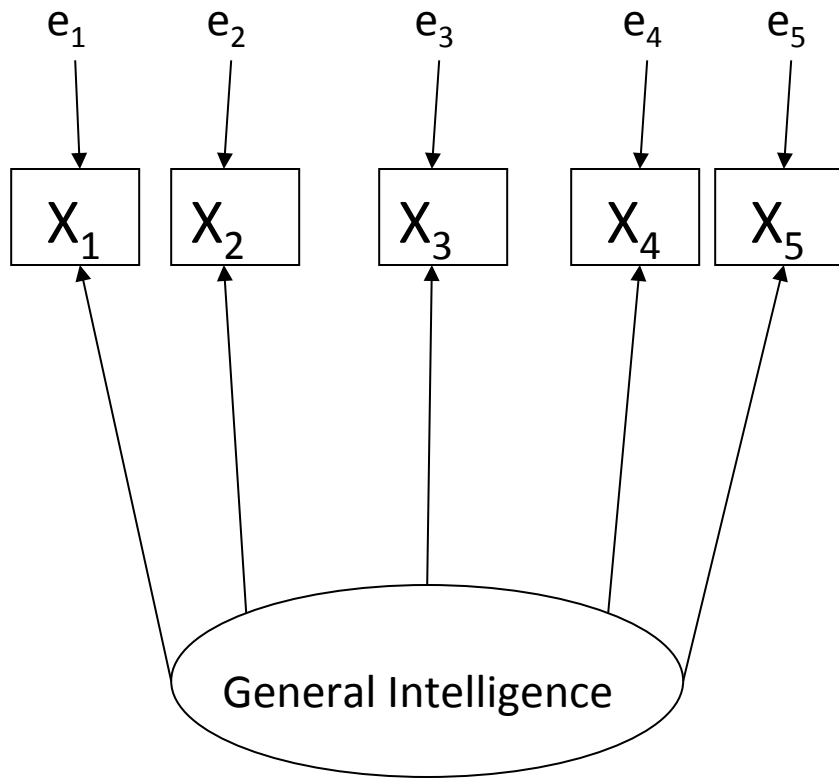
$$Y_2 = \beta Y_1 + \gamma_2 X + \epsilon_2$$

$$W = X + e_1$$

$$V_1 = Y_1 + e_2$$

$$V_2 = Y_2 + e_3$$

# A Factor Analysis Model



$$X_1 = \lambda_1 F + e_1$$

$$X_2 = \lambda_2 F + e_2$$

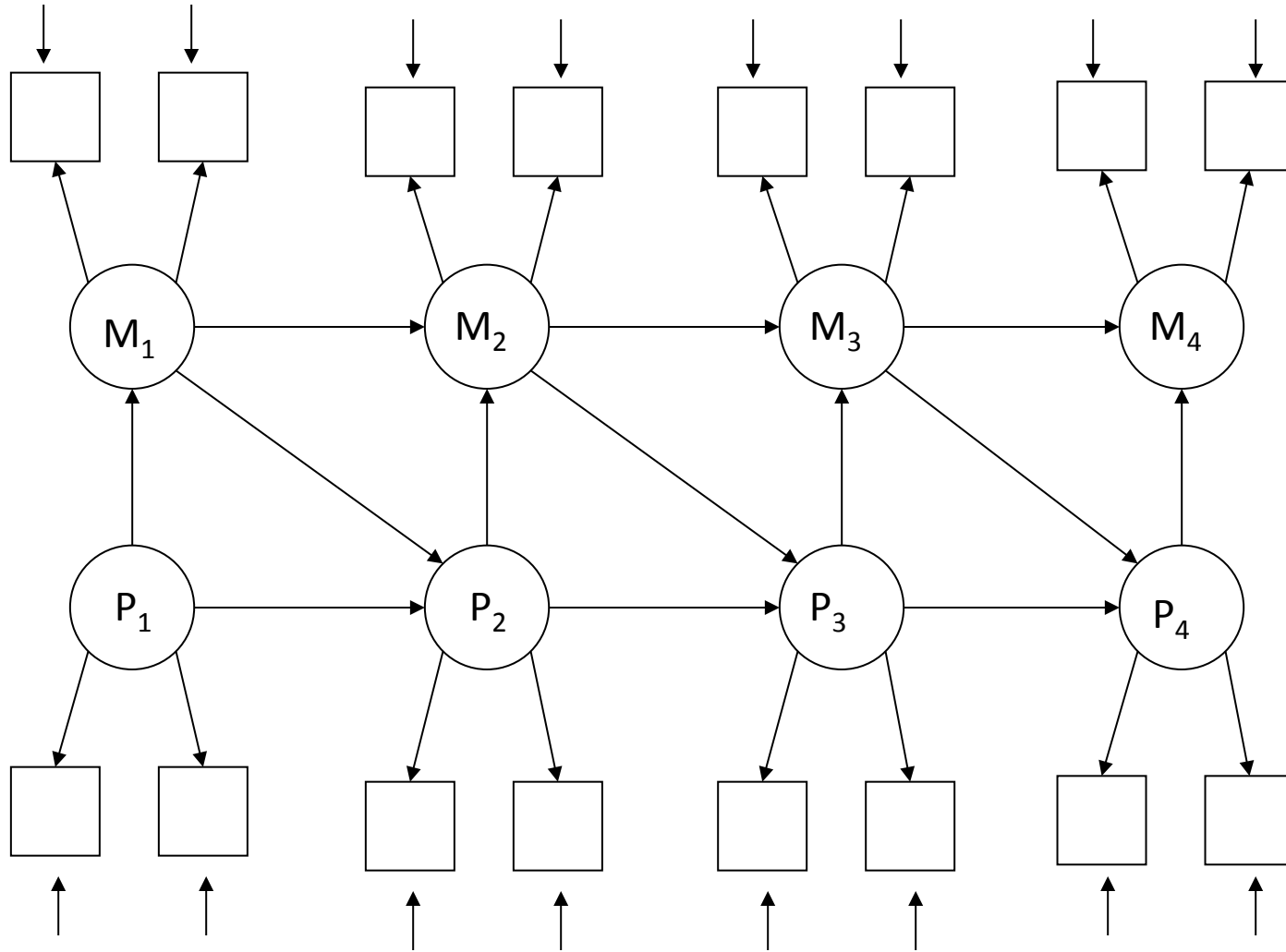
$$X_3 = \lambda_3 F + e_3$$

$$X_4 = \lambda_4 F + e_4$$

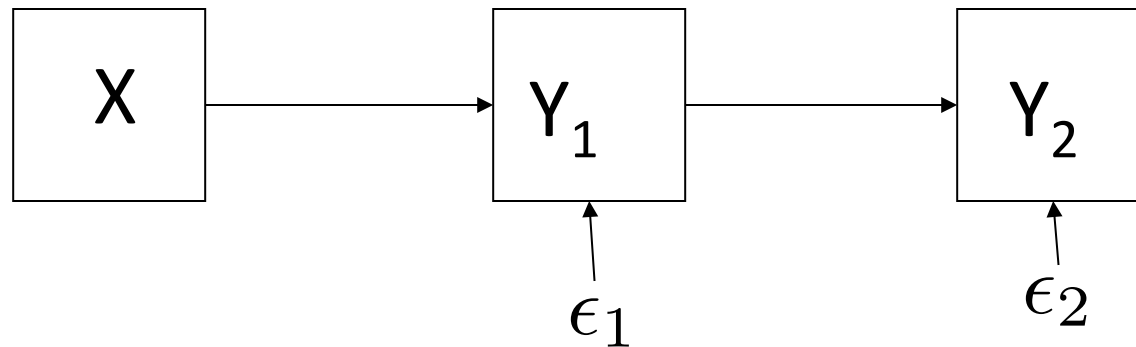
$$X_5 = \lambda_5 F + e_5$$



# A Longitudinal Model



# Estimation and Testing as Before



$$Y_1 = \gamma X + \epsilon_1$$
$$Y_2 = \beta Y_1 + \epsilon_2$$

All expected values equal zero.

$$V(X) = \phi, V(\epsilon_1) = \psi_1, V(\epsilon_2) = \psi_2,$$

$X, \epsilon_1, \epsilon_2$  are all independent.

Everything is normal.

# Distribution of the data

$\begin{pmatrix} X_1 \\ Y_{1,1} \\ Y_{1,2} \end{pmatrix} \dots \begin{pmatrix} X_n \\ Y_{n,1} \\ Y_{n,2} \end{pmatrix}$  are independent normal with mean zero.

and covariance matrix

$$\mathbf{\Sigma} = \begin{pmatrix} \phi & \gamma\phi & \beta\gamma\phi \\ \gamma\phi & \gamma^2\phi + \psi_1 & \beta(\gamma^2\phi + \psi_1) \\ \beta\gamma\phi & \beta(\gamma^2\phi + \psi_1) & \beta^2(\gamma^2\phi + \psi_1) + \psi_2 \end{pmatrix}$$

$$\boldsymbol{\theta} = (\gamma, \beta, \phi, \psi_1, \psi_2)$$

# A General Two-Stage Model

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

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

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

- $\mathbf{D}_i$  (the data) are observable. All other variables are latent.
- $\mathbf{Y}_i = \boldsymbol{\beta}\mathbf{Y}_i + \boldsymbol{\Gamma}\mathbf{X}_i + \boldsymbol{\epsilon}_i$  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 independent variable, the usual meaning of factor in experimental design.
- $\mathbf{D}_i = \boldsymbol{\Lambda}\mathbf{F}_i + \mathbf{e}_i$  is called the *Measurement Model*.

# More Details

- $\mathbf{Y}_i$  is a  $q \times 1$  random vector.
- $\boldsymbol{\beta}$  is a  $q \times q$  matrix of constants with zeros on the main diagonal.
- $\boldsymbol{\Gamma}$  is a  $q \times p$  matrix of constants.
- $\mathbf{X}_i$  is a  $p \times 1$  random vector.
- $\boldsymbol{\epsilon}_i$  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.
- $\mathbf{D}_i$  is a  $k \times 1$  random vector.
- $\mathbf{e}_i$  is a  $k \times 1$  random vector.
- $\mathbf{X}_i$ ,  $\boldsymbol{\epsilon}_i$  and  $\mathbf{e}_i$  are independent.

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

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

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

- $V(\mathbf{X}_i) = \Phi_x$

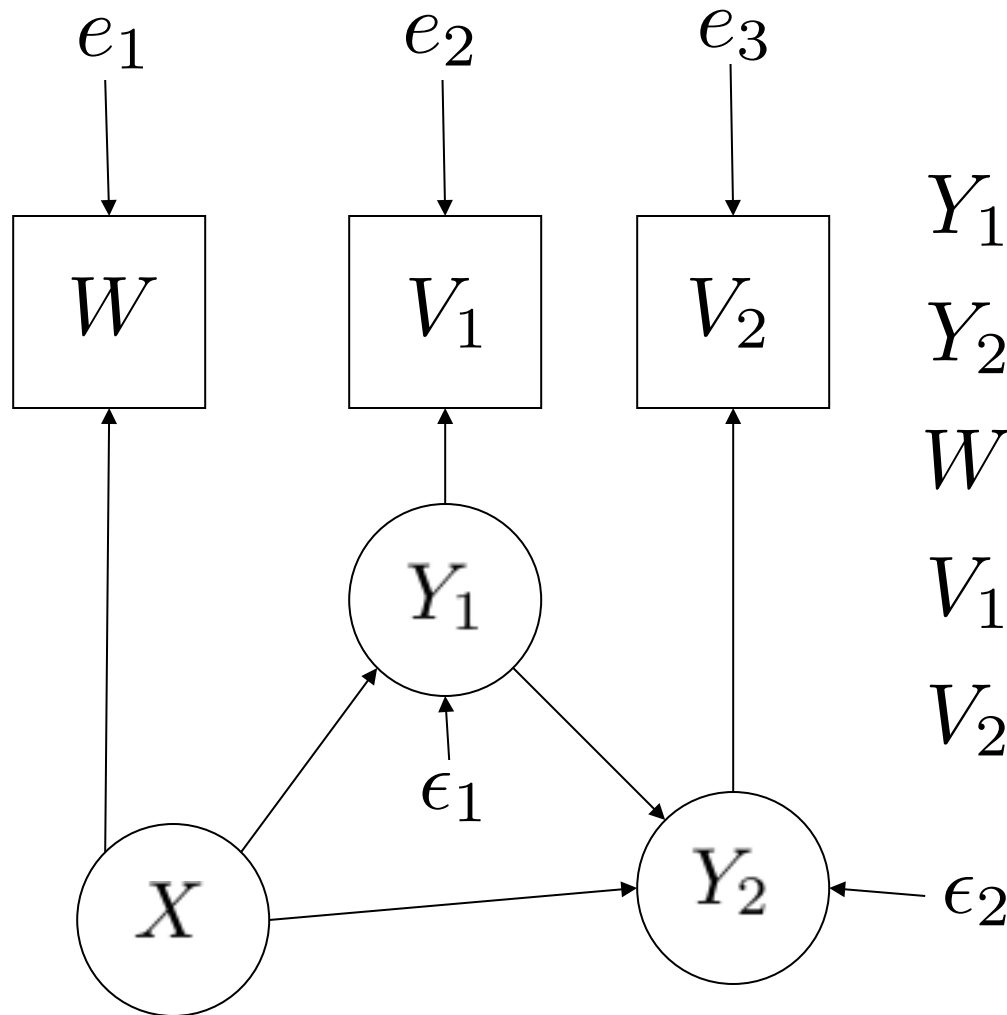
- $V(\boldsymbol{\epsilon}_i) = \Psi$

- $V(\mathbf{F}_i) = V \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix} = \begin{pmatrix} V(\mathbf{X}_i) & C(\mathbf{X}_i, \mathbf{Y}_i) \\ C(\mathbf{Y}_i, \mathbf{X}_i) & V(\mathbf{Y}_i) \end{pmatrix} = \Phi = \begin{pmatrix} \Phi_{11} & \Phi_{12} \\ \Phi'_{12} & \Phi_{22} \end{pmatrix}$

- $V(\mathbf{e}_i) = \Omega$

- $V(\mathbf{D}_i) = \Sigma$

# Recall the example



$$Y_1 = \gamma_1 X + \epsilon_1$$

$$Y_2 = \beta Y_1 + \gamma_2 X + \epsilon_2$$

$$W = X + e_1$$

$$V_1 = Y_1 + e_2$$

$$V_2 = Y_2 + e_3$$

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

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

$$\begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ \beta & 0 \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} + \begin{pmatrix} \gamma_1 \\ \gamma_2 \end{pmatrix} X + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \end{pmatrix}$$

$$\begin{pmatrix} W \\ V_1 \\ V_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} X \\ Y_1 \\ Y_2 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \end{pmatrix}$$

- $V(\mathbf{X}) = \mathbf{\Phi}_x = \phi$

- $V(\boldsymbol{\epsilon}) = \mathbf{\Psi} = \begin{pmatrix} \psi_1 & 0 \\ 0 & \psi_2 \end{pmatrix}$

- $V(\mathbf{e}) = \mathbf{\Omega} = \begin{pmatrix} \omega_1 & 0 & 0 \\ 0 & \omega_2 & 0 \\ 0 & 0 & \omega_3 \end{pmatrix}$

- $V(\mathbf{D}) = \mathbf{\Sigma}$



# Observable variables in the latent variable model (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  $P(e_j=0) = 1$

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

- So in the covariance matrix  $\mathbf{\Omega} = V(\mathbf{e})$ , just set  $\omega_{ij} = \omega_{ji} = 0, i=1, \dots, k$

# 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

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

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

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

- $V(\mathbf{X}_i) = \Phi_x$
- $V(\epsilon_i) = \Psi$
- $V(\mathbf{F}_i) = V \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix} = \begin{pmatrix} V(\mathbf{X}_i) & C(\mathbf{X}_i, \mathbf{Y}_i) \\ C(\mathbf{Y}_i, \mathbf{X}_i) & V(\mathbf{Y}_i) \end{pmatrix} = \Phi = \begin{pmatrix} \Phi_{11} & \Phi_{12} \\ \Phi'_{12} & \Phi_{22} \end{pmatrix}$
- $V(\mathbf{e}_i) = \Omega$
- $V(\mathbf{D}_i) = \Sigma$

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

$$\begin{aligned} \mathbf{Y} &= \beta\mathbf{Y} + \Gamma\mathbf{X} + \epsilon \\ \Rightarrow \mathbf{Y} - \beta\mathbf{Y} &= \Gamma\mathbf{X} + \epsilon \\ \Rightarrow \mathbf{I}\mathbf{Y} - \beta\mathbf{Y} &= \Gamma\mathbf{X} + \epsilon \\ \Rightarrow (\mathbf{I} - \beta)\mathbf{Y} &= \Gamma\mathbf{X} + \epsilon \\ \Rightarrow (\mathbf{I} - \beta)^{-1}(\mathbf{I} - \beta)\mathbf{Y} &= (\mathbf{I} - \beta)^{-1}(\Gamma\mathbf{X} + \epsilon) \\ \Rightarrow \mathbf{Y} &= (\mathbf{I} - \beta)^{-1}\Gamma\mathbf{X} + (\mathbf{I} - \beta)^{-1}\epsilon \end{aligned}$$

So,

$$\begin{aligned} V(\mathbf{Y}) &= V((\mathbf{I} - \beta)^{-1}\Gamma\mathbf{X} + (\mathbf{I} - \beta)^{-1}\epsilon) \\ &= V((\mathbf{I} - \beta)^{-1}\Gamma\mathbf{X}) + V((\mathbf{I} - \beta)^{-1}\epsilon) \\ &= (\mathbf{I} - \beta)^{-1}\Gamma V(\mathbf{X}) ((\mathbf{I} - \beta)^{-1}\Gamma)^{\top} + (\mathbf{I} - \beta)^{-1}V(\epsilon)(\mathbf{I} - \beta)^{-1\top} \\ &= (\mathbf{I} - \beta)^{-1}\Gamma \Phi_{11}\Gamma^{\top} (\mathbf{I} - \beta)^{-1\top} + (\mathbf{I} - \beta)^{-1}\Psi(\mathbf{I} - \beta)^{-1\top} \end{aligned}$$

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

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

# Two-stage Proofs of Identifiability

- Show the parameters of the measurement model ( $\Lambda, \Phi, \Omega$ ) can be recovered from  $\Sigma = V(\mathbf{D})$ .
- Show the parameters of the latent variable model ( $\beta, \Gamma, \Phi_{11}, \Psi$ ) can be recovered from  $\Phi = V(\mathbf{F})$ .
- This means *all* the parameters can be recovered from  $\Sigma$ .
- Break a big problem into two smaller ones.
- Develop *rules* for checking identifiability at each stage.

# Copyright Information

This slide show was prepared by Jerry Brunner, Department of Statistics, 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. These Powerpoint slides are available from the course website:

<http://www.utstat.toronto.edu/~brunner/oldclass/431s15>