

Path diagrams as a notational formalism (Noodling around with pictures)

Vince Wiggins

Vice President, Scientific Development

StataCorp LP

2013 London Stata Users Group Meeting

Sometimes path diagrams are used in lieu of mathematical notation to define and explain models.

How far can they be pushed?

Mathematical notation for a covariate

Mathematical notation for a covariate

X

Mathematical notation for a covariate

X

Path notation for a covariate

Mathematical notation for a covariate

X

Path notation for a covariate



Mathematical notation for a dependent variable

Mathematical notation for a dependent variable

y

Mathematical notation for a dependent variable

y

Path notation for a dependent variable

Mathematical notation for a dependent variable

y

Path notation for a dependent variable



Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

β

Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

β

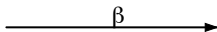
Path notation for a parameter

Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

β

Path notation for a parameter

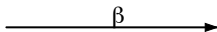


Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

$$\beta$$

Path notation for a parameter



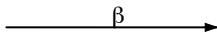
Paths are truly parameters in a linear form, and must connect observed or latent variables. So better is

Maths and Paths — Parameters/Coefficients

Mathematical notation for a parameter

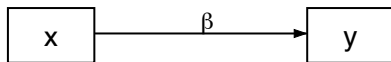
$$\beta$$

Path notation for a parameter



Paths are truly parameters in a linear form, and must connect observed or latent variables. So better is

$$y = \beta x + \dots$$



Mathematical notation for a latent variable

Mathematical notation for a latent variable

μ_j *unobserved*

$\mu_j \sim \text{Normal}(0, \sigma_\mu)$

Mathematical notation for a latent variable

μ_j *unobserved*

$\mu_j \sim \text{Normal}(0, \sigma_\mu)$

Path notation for a latent variable

Mathematical notation for a latent variable

μ_j *unobserved*

$\mu_j \sim \text{Normal}(0, \sigma_\mu)$

Path notation for a latent variable



A digression on learning to love latent variables

They can be:

- an error term
- a random intercept
- part of a random slope
- a frailty
- a creator/modeler of endogeneity
- a creator/modeler of sample selection
- a creator/modeler of endogenous treatment effects
- part of a switching model
- or any other latent/unobserved thingy

Linear regression

Mathematical notation for a linear regression

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

$$\epsilon_j \sim \text{Normal}(0, \sigma_\epsilon)$$

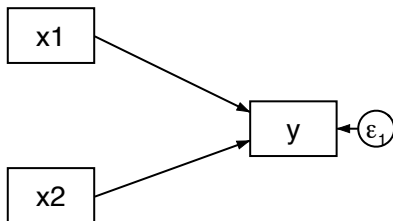
Linear regression

Mathematical notation for a linear regression

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

$$\epsilon_j \sim \text{Normal}(0, \sigma_\epsilon)$$

Path notation for a linear regression



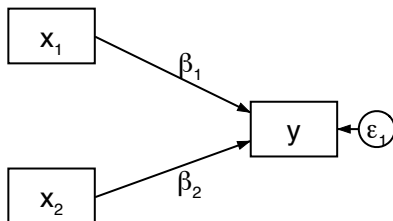
Linear regression

Mathematical notation for a linear regression

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

$$\epsilon_j \sim \text{Normal}(0, \sigma_\epsilon)$$

Path notation for a linear regression — with coefficients



Linear multilevel model

Mathematical notation for a multilevel model

$$y_{ji} = \beta_0 + \beta_1 x_{1ji} + (\beta_2 + \mu_{1j}) x_{2ji} + \mu_{0j} + \epsilon_{ji}$$

$$\mu_{0j} \sim \text{Normal}(0, \sigma_{\mu_0})$$

$$\mu_{1j} \sim \text{Normal}(0, \sigma_{\mu_1})$$

Linear multilevel model

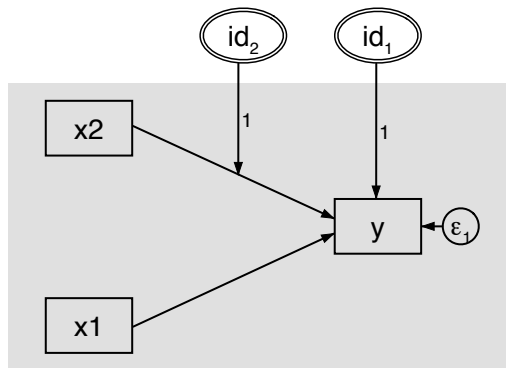
Mathematical notation for a multilevel model

$$y_{ji} = \beta_0 + \beta_1 x1_{ji} + (\beta_2 + \mu1_j)x2_{ji} + \mu0_j + \epsilon_{ji}$$

$$\mu0_j \sim \text{Normal}(0, \sigma_{\mu0})$$

$$\mu1_j \sim \text{Normal}(0, \sigma_{\mu1})$$

Path notation for a multilevel model



Linear multilevel model

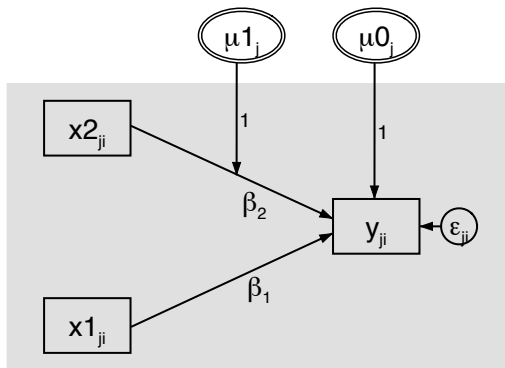
Mathematical notation for a multilevel model

$$y_{ji} = \beta_0 + \beta_1 x1_{ji} + (\beta_2 + \mu1_j)x2_{ji} + \mu0_j + \epsilon_{ji}$$

$$\mu0_j \sim \text{Normal}(0, \sigma_{\mu0})$$

$$\mu1_j \sim \text{Normal}(0, \sigma_{\mu1})$$

Path notation for a multilevel model — with maths



A digression on the SEM Builder

Every diagram in this talk was drawn with it.

They each took 2 or 3 minutes to draw.

They can all be estimated by `sem` or `gsem`.

Let's draw one ...

If you are reading this online, you just had to be there.

Linear multivariate (seemingly unrelated) multilevel model

Mathematical notation

$$y_{1ji} = \beta_{0a} + \beta_1 x_{1ji} + (\beta_2 + \mu_{1j}) x_{2ji} + \mu_{0j} + \epsilon_{1ji}$$

$$y_{2ji} = \beta_{0b} + \beta_3 x_{3ji} + (\beta_4 + \mu_{2j}) x_{2ji} + \mu_{3j} + \epsilon_{2ji}$$

$$\mu_{0j} \sim \text{Normal}(0, \sigma_{\mu 0})$$

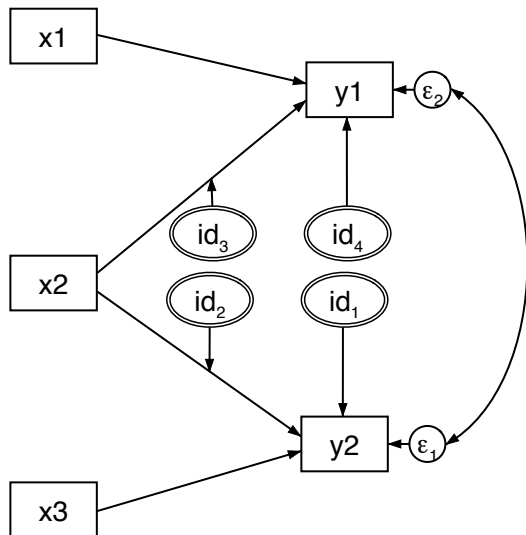
$$\mu_{1j} \sim \text{Normal}(0, \sigma_{\mu 1})$$

$$\mu_{2j} \sim \text{Normal}(0, \sigma_{\mu 2})$$

$$\mu_{3j} \sim \text{Normal}(0, \sigma_{\mu 3})$$

j indexes levels

Linear multivariate (seemingly unrelated) multilevel model



What if the random effects are shared?

- Easy enough in the mathematical notation, but easier to see in the path diagram. Let's edit the path diagram ...
- *I am sorry once again to those following at home.*

Variations on the linear multivariate multilevel model

What if the random effects are shared?

- Easy enough in the mathematical notation, but easier to see in the path diagram. Let's edit the path diagram ...
- *I am sorry once again to those following at home.*

What if the random coefficients on x_2 are the same for both y 's?

- Again, very easy to see (and to change) on the path diagram ...
- Must explicitly constrain paths/coefficients from latent variables.

Variations on the linear multivariate multilevel model

What if the random effects are shared?

- Easy enough in the mathematical notation, but easier to see in the path diagram. Let's edit the path diagram ...
- *I am sorry once again to those following at home.*

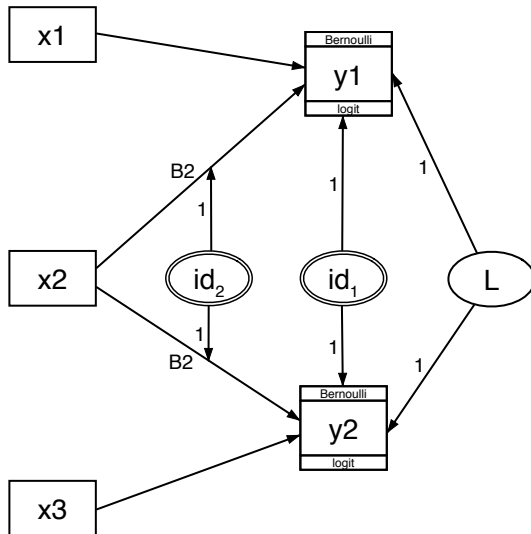
What if the random coefficients on x_2 are the same for both y 's?

- Again, very easy to see (and to change) on the path diagram ...
- Must explicitly constrain paths/coefficients from latent variables.

What if the outcomes were 0/1 and we wanted to model them as logistic?

- We would add some probabilistic statements to our mathematical notation. And likely some verbal explanation.
- We can also designate this on the path diagram ...

Final Variation



Time series — ARMAX model

Mathematical notation

$$y_t = \beta x_t + \epsilon_t$$

$$\epsilon_t = \rho \epsilon_{t-1} + \mu_t$$

$$\mu_t \sim \text{Normal}(0, \sigma_\mu)$$

Time series — ARMAX model

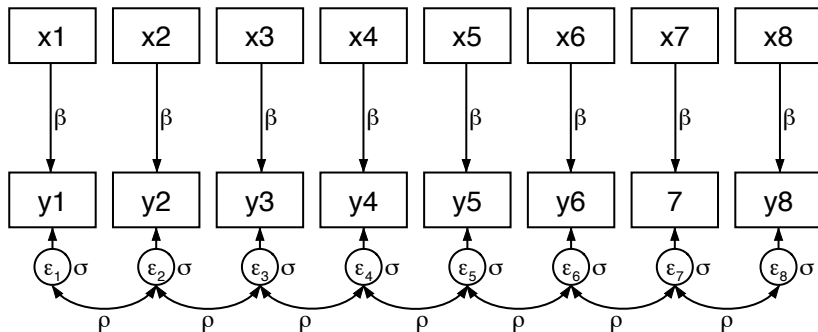
Mathematical notation

$$y_t = \beta x_t + \epsilon_t$$

$$\epsilon_t = \rho \epsilon_{t-1} + \mu_t$$

$$\mu_t \sim \text{Normal}(0, \sigma_\mu)$$

Path notation



Regression with sample selection — Heckman style

Mathematical notation

$$y_i = \beta_0 + \beta_1 x1_i + \beta_2 x2_i + \epsilon_i$$

$$s_i = \lambda_0 + \lambda_2 x2_i + \lambda_3 x3_i + \xi_i$$

$$\text{cov}(\epsilon, \xi) \neq 0$$

y_i observed when $s_i > 0$

Regression with sample selection — Heckman style

Mathematical notation

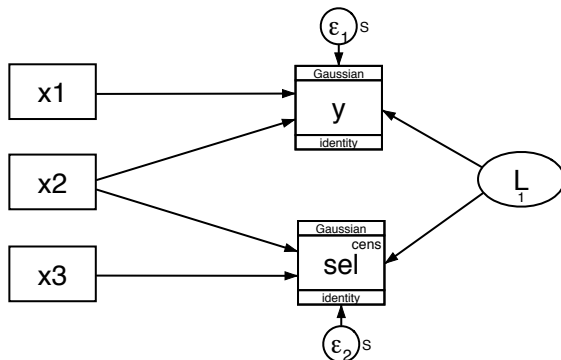
$$y_i = \beta_0 + \beta_1 x1_i + \beta_2 x2_i + \epsilon_i$$

$$s_i = \lambda_0 + \lambda_2 x2_i + \lambda_3 x3_i + \xi_i$$

$$\text{cov}(\epsilon, \xi) \neq 0$$

y_i observed when $s_i > 0$

Path notation



Logistic model with endogenous covariate

Mathematical notation

$$Pr(y_i) = \text{logit}(\beta_0 + \beta_1 x_{1i} + \beta_2 z_i + UC_i)$$

$$z_i = \gamma_0 + \gamma_1 x_{2i} + UC_i + \nu_i$$

$$\text{cov}(\epsilon, \xi) = 0$$

Logistic model with endogenous covariate

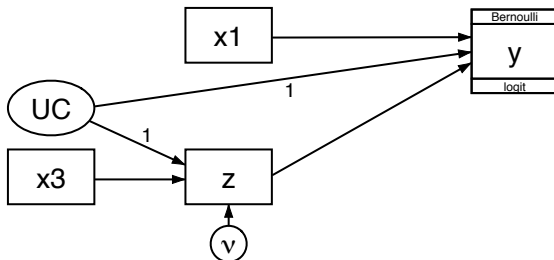
Mathematical notation

$$Pr(y_i) = \text{logit}(\beta_0 + \beta_1 x1_i + \beta_2 z_i + UC_i)$$

$$z_i = \gamma_0 + \gamma_1 x2_i + UC_i + \nu_i$$

$$\text{cov}(\epsilon, \xi) = 0$$

Path notation



Selection with an endogenous covariate in a multilevel framework

Regression component

$$y_{ji} = \beta_0 + (\beta_1 + \mu_1 \mathbf{1}_j) x_{1ji} + \beta_2 z_{ji} + UC_{ji} + \mu_0 \mathbf{0}_j + \epsilon_{ji}$$

Selection component

$$s_{ji} = \lambda_0 + \lambda_1 z_i + \lambda_2 x_{2i} + \xi_i$$

$$\text{cov}(\epsilon, \xi) \neq 0$$

y_{ji} observed when s_{ji} > 0

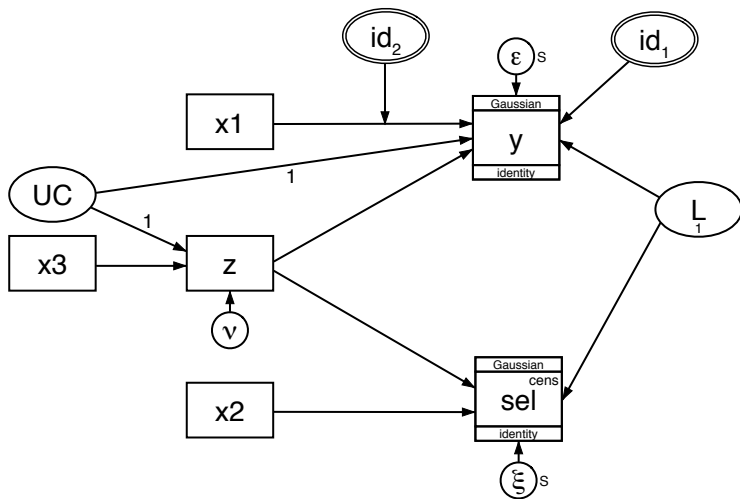
Endogenous component

$$z_{ji} = \gamma_0 + \gamma_1 x_{3ji} + UC_{ji} + \nu_i$$

$$\text{cov}(\epsilon, \nu) = 0$$

$$\text{cov}(\xi, \nu) = 0$$

Selection with an endogenous covariate in a multilevel framework — path diagram



Compared to math, path diagrams are:

- easier to understand
- more attuned to human perception
- less precise
- less flexible
- more amenable to noodling

Don't choose.

Use what works — words, maths, paths.