

PSYC 402: Structural Equation Modelling

By Jack Zhou

1 Linear Algebra Review

TBD

2 Three classes of models

Structural Equations can be classified by 3 different classes.

2.1 Path Model

A path models can also be referred to as a multivariate regression model for observed variables, or a causal model. A path model can contain p dependent (*endogenous*) variables and q independent (*exogenous*) variables, denoted as the following:

- $\tilde{\mathbf{y}}$: a $p \times 1$ vector of *endogenous* observed variables
- $\tilde{\mathbf{x}}$: a $q \times 1$ vector of *exogenous* observed variables
- $\tilde{\zeta}$: a $p \times 1$ vector of residuals

As suggested by the match of dimensionality between $\tilde{\mathbf{y}}$ and $\tilde{\zeta}$, each of the entries in $\tilde{\mathbf{y}}$ corresponds to a value in $\tilde{\zeta}$ with the same index. That is, for each dependent variable \tilde{y}_i , we have a corresponding residual $\tilde{\zeta}_i$ that the model is unable to account for.

It is important to note that path models do not contain any latent (unobserved) variables, as we can see in the following equation:

$$\tilde{\mathbf{y}} = [\mathbf{B}\tilde{\mathbf{y}} + \mathbf{\Gamma}\tilde{\mathbf{x}}] + \tilde{\zeta}$$

The general interpretation is that with this model, the set of endogenous variables are to be described in terms of a sum between the structural model $[\mathbf{B}\tilde{\mathbf{y}} + \mathbf{\Gamma}\tilde{\mathbf{x}}]$ and residuals $\tilde{\zeta}$.

The model consists of two parts, the $\mathbf{B}\tilde{\mathbf{y}}$ term describes the regression of endogenous variables onto themselves, and the $\mathbf{\Gamma}\tilde{\mathbf{x}}$ term describes the regression of exogenous variables onto the endogenous variables, in combined effort to best account for the endogenous variables.

One glaring issue is that in order to solve for $\tilde{\mathbf{y}}$, it should only be on one side of the equation, not both. We can address this with algebra:

$$\begin{aligned}\tilde{\mathbf{y}} &= [\mathbf{B}\tilde{\mathbf{y}} + \mathbf{\Gamma}\tilde{\mathbf{x}}] + \tilde{\zeta} \\ \tilde{\mathbf{y}} - \mathbf{B}\tilde{\mathbf{y}} &= \mathbf{\Gamma}\tilde{\mathbf{x}} + \tilde{\zeta} \\ (\mathbf{I} - \mathbf{B})\tilde{\mathbf{y}} &= \mathbf{\Gamma}\tilde{\mathbf{x}} + \tilde{\zeta} \\ \tilde{\mathbf{y}} &= (\mathbf{I} - \mathbf{B})^{-1}(\mathbf{\Gamma}\tilde{\mathbf{x}} + \tilde{\zeta})\end{aligned}$$

The equation $\tilde{\mathbf{y}} = (\mathbf{I} - \mathbf{B})^{-1}(\mathbf{\Gamma}\tilde{\mathbf{x}} + \tilde{\zeta})$ is known as the reduced form, which is easier to work with mathematically. Below are descriptions to all the matrices in a path model:

- \mathbf{B} : $p \times p$ regression coefficients for $\tilde{\mathbf{y}}$
- $\mathbf{\Gamma}$: $q \times p$ regression coefficients for $\tilde{\mathbf{x}}$
- $\mathbf{\Phi}$: $q \times q$ covariance matrix of $\tilde{\mathbf{x}}$
- $\mathbf{\Psi}$: $p \times p$ covariance matrix of $\tilde{\zeta}$

Note that $\mathbf{\Phi}$ and $\mathbf{\Psi}$ are matrices that aren't seen in the model.

For any covariance matrix \mathbf{M} , any \mathbf{M}_{ii} entry is the variance of the i th variable and any \mathbf{M}_{ij} entry is the covariance between the i th and j th variable.

$\mathbf{\Phi}$ can be considered a parameter for $\tilde{\mathbf{y}}$. Therefore, we consider it as a starting point.

Often times we see correlations between residuals (ex. in longitudinal models), so $\mathbf{\Psi}$ does not necessarily have to be diagonal.

2.1.1 Kinds of Path Models

• Standard Multiple Regression

Suppose we have an endogenous variable y_1 , its corresponding ζ_1 , and a set of exogenous variables $\{x_1, \dots, x_q\}$. A Standard Multiple Regression can be seen as a digraph with the following properties:

- An edge (ζ_1, y_1) exists; the residual impacts the endogenous variable
- For all x_i in the set of exogenous variables, an edge (x_i, y_1) exists; each of the exogenous variables impacts the endogenous variable
- For all $x_i \neq x_j$, an edge (x_i, x_j) exists to indicate correlations between exogenous variables

Let $\tilde{\mathbf{x}}$ be a $p \times 1$ vector of exogenous variables, $\mathbf{\Gamma}$ be a $1 \times q$ vector of regression coefficients for $\tilde{\mathbf{x}}$, then such graph can be written as an equation as follows:

$$y_1 = \mathbf{\Gamma}\tilde{\mathbf{x}} + \zeta_1$$

• Standard Multivariate Regression

In this case, there are greater than $p > 1$ endogenous variables as well as $q > 1$ exogenous variables.

The setting is similar to the Standard Multiple Regression case, but there is more than one y_i , but no edges between any $y_i \neq y_j$ (no correlation between any endogenous variables).

In addition to the notation used in standard multiple regression, we let $\tilde{\mathbf{y}}$ be the $p \times 1$ vector of endogenous variables, and \mathbf{B} be a $p \times p$ matrix as defined before. As the model asserts that there are no correlations between any endogenous variables, \mathbf{B} is a zero matrix, and the covariance matrix of $\tilde{\mathbf{y}}$, $\mathbf{\Psi}$, is taken to be diagonal.

- **First Order Autoregressive Time Series**

A chain of measurements are collected in sequential order, imposing a temporal structure.

The first measurement is considered exogenous, so we would name it x_1 . From there on, each measurement that comes after are all endogenous variables y_i with residuals ζ_i .

Each variable is linearly dependent on the variable before, causing a chain of dependencies.

2.2 Linear Factor/Measurement Models

In the basic treatment, we have the following variables.

- $\tilde{\mathbf{x}}$: a $q \times 1$ vector of *indicators* (observed variables)
- $\tilde{\xi}$: a $n \times 1$ vector of *latent*, common-factor variables
- $\tilde{\delta}$: a $q \times 1$ vector of measurement error variables

In a parallel alternative form, we have these instead:

- $\tilde{\mathbf{y}}$: a $q \times 1$ vector of *indicators* (observed variables)
- $\tilde{\eta}$: a $n \times 1$ vector of *latent*, common-factor variables
- $\tilde{\epsilon}$: a $q \times 1$ vector of measurement error variables

Note that here we have a set of observed variables and a set of latent variables. The observed variables play the role in trying to determine the latent variables. Imagine measuring intelligence (abstract and latent) with test scores (indicators; concrete and observable).

Under Spearman, all $\tilde{\mathbf{x}}_i$ is causally (regression) dependent upon a ξ_i for some i . Moreover, upon regressing $\tilde{\mathbf{x}}$ on $\tilde{\xi}$, within-correlations in $\tilde{\mathbf{x}}$ no longer exist.

In other words, the reason why different $\tilde{\mathbf{x}}_i$ could be correlated is because they arise from a common latent variable; $\tilde{\xi}$ can be taken as an explanation as to why within-correlations in $\tilde{\mathbf{x}}$ exist.

We can summarize this model with the following formula:

$$\tilde{\mathbf{x}} = \mathbf{\Lambda}_{\mathbf{x}} \tilde{\xi} + \tilde{\delta}$$

Where $\mathbf{\Lambda}_{\mathbf{x}}$ is $q \times n$ matrix called factor loadings, linking the indicators to the latent variables.

There is also a $n \times n$ matrix Φ , which is the covariance matrix of $\tilde{\xi}$. Typically there are many restriction on this matrix, for example we would assert it to be a diagonal matrix for orthogonality.

Finally, we also have covariance for $\tilde{\delta}$ as Θ_{δ} , a $q \times q$ matrix. Recall that upon regression, the indicators become uncorrelated, suggesting Θ_{δ} is typically diagonal.

2.2.1 Correlation accross sets of indicator/latent variables

Sometimes we are interested in the correlation between multiple latent variables, and the only access we have towards them are their indicators, so the correlation would have to run on some linear weighted composite of the indicators per grouping.

The resulting correlation, which is used to indicate the correlation between the actual latent variables, is known to be less than the actual correlation between the latent variables.

To account for this issue, a mechanism of disattenuation is used. For two latent variables τ_1, τ_2 , their correlation ρ_{τ_1, τ_2} can be determined by a function of observed composites c_1, c_2 (respectively) as follows:

$$\rho_{\tau_1, \tau_2} = \frac{\rho_{c_1, c_2}}{\sqrt{\rho_{c_1, c'_1} \cdot \rho_{c_2, c'_2}}}$$

Where ρ_{c_1, c'_1} denotes reliability for c_1 , and likewise for c_2 . This is known as the disattenuating correlation formula. With structural equation modelling, we can add an edge to the path diagram that joins the latent variables and work from there, without the need for disattenuation.