# Longitudinal Data Analysis
# Using Structural Equation Modeling

Paul Allison, Ph.D.

146 ☐ **Allow Patents to Affect Later R&D**

147 ☐ **Allow Patents to Affect Later Patents**

148 ☐ **Include Science and LogSize**

149 ☐ **Summary**

# Longitudinal Data Analysis
# Using **sem**

Paul D. Allison, Ph.D.

February 2016

# Causal Inference

How do you make causal inferences with non-experimental panel data?

What's panel data?

- Data in which variables are measured at multiple points in time for the same individuals.
  - Response variable $y_{it}$ with $t = 1, 2,…, T; i = 1,…, N$
  - Vector of predictor variables $x_{it}$.
  - Some of the predictors may vary with time, others may not.
- Assume that time points are the same for everyone in the sample.
  - Nice if they are equally spaced, but not essential.
  - We will eventually allow for drop out and other kinds of missing data.

# Causal Inference

How do you demonstrate that *x* causes *y*?
– Show that *x* is correlated with *y*.
– Show that the correlation is not produced by other variables that affect both *x* and *y*.
– Show that *y* is not causing *x*.
– Show that the correlation is not due to chance alone.

Randomized experiments are great at all these things.

Panel data make it possible to
– Control for unobserved variables.
– Estimate the effect of *x* on *y*, even if *y* is also affecting *x*.

No standard methods can do both of these things simultaneously.

# Fixed Effects Methods

To control for unobservables, we can used *fixed effects* methods
- These control for all unchanging variables whether observed or not.
- For linear models, the most common way to do fixed effects is to express all variables as deviations from individual-specific means.  But there are several alternative approaches.
    - Implement FE with **xtreg** or **xtlogit** in Stata
- Downsides:
    – Standard errors go up (because you're only using within-individual variation).
    – Many methods don't produce estimates for time-invariant predictors.

# Some References



2009



2005

# Cross-Lagged Linear Models

To allow for reciprocal causation, estimate 2-wave, 2-variable panel model ([OD Duncan 1969](#)) by ordinary least squares:

$$y_2 = b_0 + b_1 y_1 + b_2 x_1 + \varepsilon_2$$
$$x_2 = a_0 + a_1 y_1 + a_2 x_1 + \varepsilon_1$$



Inclusion of lagged dependent variable is intended to control for past characteristics of the individual.

Among those with the same value of $y_1$, $b_2$ is the effect of $x_1$ on $y_2$.

# Our Goal

To be able to estimate models that combine fixed effects with cross-lags using structural equation modeling software. The models look like this:

Cross-lagged Effect

Fixed Effect

$$y_{it} = \mu_t + \beta_1 x_{i,t-1} + \beta_2 y_{i,t-1} + \delta_1 w_{it} + \gamma_1 z_i + \alpha_i + \varepsilon_{it}$$

$$x_{it} = \tau_t + \beta_3 x_{i,t-1} + \beta_4 y_{i,t-1} + \delta_2 w_{it} + \gamma_2 z_i + \eta_i + \upsilon_{it}$$

Cross-lagged Effect

Fixed Effect

To get there, we'll
- Review models with cross-lagged effects using SEM.
- Review conventional fixed effects
- See how to do fixed effects with SEM
- Combine the two methods

# Path Analysis of Observed Variables

In the SEM literature, it's common to represent a linear model by a path diagram.

– A diagrammatic method for representing a system of linear equations. There are precise rules so that you can write down equations from looking at the diagram.

– Single equation:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

## Some Rules and Definitions

→  Direct causal effect

⌒  Correlation
(no causal assumptions)

Why the curved double-headed arrow in the diagram?
Because omitting it implies no correlation between $x_1$ and $x_2$.

Endogenous variables:  Variables caused by other variables in the system.  These variables have straight arrows leading into them.

Exogenous variables:  Variables not caused by others in the system.  No straight arrows leading into them.

Not the same as dependent and independent because a variable that is dependent in one equation and independent in another equation is still endogenous.

Curved double-headed arrows can only link *exogenous* variables.

9

## Three Predictor Variables



The fact that there are no curved arrows between $\varepsilon$ and the x's implies that $\rho_{1\varepsilon} = 0$, $\rho_{2\varepsilon} = 0$, and $\rho_{3\varepsilon} = 0$.  We make this assumption in the usual linear regression model.

10

# Two-Equation System

$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon_1$

$x_2 = \alpha_0 + \alpha_1 x_1 + \varepsilon_2$

The diagram is now



Note: The diagram goes further than the equations by asserting that

$$\rho_{\varepsilon_1 \varepsilon_2} = 0, \ \rho_{\varepsilon_1 x_1} = 0, \ \rho_{\varepsilon_1 x_2} = 0, \ \rho_{x_1 \varepsilon_2} = 0$$

# Cross-Lagged Linear Models

$y_2 = b_0 + b_1 y_1 + b_2 x_1 + \varepsilon_2$

$x_2 = a_0 + a_1 y_1 + a_2 x_1 + \varepsilon_1$



- This model can be estimated by ordinary least squares for each equation separately.
- Other predictors could also be included in each equation.
- Presumes no simultaneous causation.

# 3 Wave-2 Variable Model



- Can extend to more waves
- Each of the 4 equations could be estimated by OLS
- Can estimate simultaneously via SEM
  - Can constrain coefficients to be the same across waves.
  - Can test the overall fit of the model (OK to omit lag-2 effects?)
  - Can handle missing data by full information maximum likelihood

# NLSY Data Set

581 children interviewed in 1990, 1992, and 1994 as part of the National Longitudinal Survey of Youth (NLSY).

<u>Time-varying variables (measured at each of the three time points)</u>:

ANTI       antisocial behavior, measured with a scale from 0 to 6.

SELF       self-esteem, measured with a scale ranging from 6 to 24.

POV       poverty status of family, coded 1 for family in poverty, otherwise 0.

<u>Time-invariant variables</u>:

| | |
|---|---|
| BLACK | 1 if child is black, otherwise 0 |
| HISPANIC | 1 if child is Hispanic, otherwise 0 |
| CHILDAGE | child's age in 1990 |
| MARRIED | 1 if mother was currently married in 1990, otherwise 0 |
| GENDER | 1 if female, 0 if male |
| MOMAGE | mother's age at birth of child |
| MOMWORK | 1 if mother was employed in 1990, otherwise 0 |

Data are in the "wide form": one record for each child , with different names for the variables at each time point, e.g., ANTI90, ANTI92 and ANTI94.

# Estimating a Cross-Lagged Model

- We'll estimate the 3W-2V panel model with SEM to answer the question, does antisocial behavior affect self-esteem, or does self-esteem affect antisocial behavior?
- Other variables could be included, but we'll leave them out for simplicity.
- Cross-sectionally, these variables are significantly correlated at about -.15.
- Important to allow for correlated errors. Why? Other factors affecting both variables are not included.
- No missing data in this data set.
- We'll see how to do it with the **sem** command.

# Software for SEMs

**LISREL** – Karl Jöreskog and Dag Sörbom

**EQS** – Peter Bentler

**PROC CALIS (SAS)** – W. Hartmann, Yiu-Fai Yung

**OpenMX** (**R**) – Michael Neale

**Amos** – James Arbuckle

**Mplus** – Bengt Muthén

**sem, gsem (Stata)**

**lavaan (R)** – Yves Rosseel

# Stata Program

```
use "C:\data\nlsy.dta", clear
sem (anti94 <- anti92 self92) ///
    (anti92 <- anti90 self90) ///
    (self94 <- anti92 self92) ///
    (self92 <- anti90 self90), ///
  cov(e.anti94*e.self94 e.anti92*e.self92)
```

- Stata is case sensitive
- <- means "is regressed on"
- e.anti94 refers to the error term for anti94
- The **cov** option allows for covariances (and therefore correlations) between pairs of variables.
- /// goes to a new a line within a single command, in a DO file.

# Stata Results

| | Coef. | OIM Std. Err. | z | P>\|z\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| **anti94 <-** | | | | | | |
| anti92 | .6606303 | .0365617 | 18.07 | 0.000 | .5889707 | .7322899 |
| self92 | .0214813 | .0161547 | 1.33 | 0.184 | -.0101814 | .0531441 |
| _cons | .2555298 | .3492365 | 0.73 | 0.464 | -.428961 | .9400207 |
| **anti92 <-** | | | | | | |
| anti90 | .6721907 | .0340779 | 19.73 | 0.000 | .6053992 | .7389822 |
| self90 | -.016042 | .0157035 | -1.02 | 0.307 | -.0468202 | .0147362 |
| _cons | .8635114 | .3296234 | 2.62 | 0.009 | .2174614 | 1.509561 |
| **self92 <-** | | | | | | |
| anti90 | -.1297929 | .0946541 | -1.37 | 0.170 | -.3153115 | .0557257 |
| self90 | .3520793 | .0436176 | 8.07 | 0.000 | .2665903 | .4375682 |
| _cons | 13.49922 | .9155554 | 14.74 | 0.000 | 11.70476 | 15.29367 |
| **self94 <-** | | | | | | |
| anti92 | .0190816 | .0819248 | 0.23 | 0.816 | -.1414881 | .1796512 |
| self92 | .3521843 | .0361984 | 9.73 | 0.000 | .2812368 | .4231319 |
| _cons | 13.41623 | .7825436 | 17.14 | 0.000 | 11.88247 | 14.94999 |

# Stata Results (cont.)

```
-------------+-------------------------------------------------------------
var(e.anti94)|    1.81788    .1066577                        1.620406    2.039419
var(e.anti92)|   1.436514    .0842824                        1.280467    1.611577
var(e.self92)|   11.08263    .6502342                        9.878743    12.43324
var(e.self94)|   9.127317    .5355131                        8.135831    10.23963
-------------+-------------------------------------------------------------
cov(e.anti94,|
    e.self94)|  -.6221672    .1709519    -3.64    0.000    -.9572267    -.2871077
cov(e.anti92,|
    e.self92)|   -.656332     .167759    -3.91    0.000    -.9851335    -.3275305
---------------------------------------------------------------------------
LR test of model vs. saturated: chi2(4)    =      48.07, Prob > chi2 = 0.0000
```

The LR (likelihood ratio) test is testing the null hypothesis that all four
two-period lagged paths are 0. Clearly, that must be rejected.

# Path Diagram

# Estimation & Assumptions

By default, **sem** does maximum likelihood (ML) estimation:
- Choose parameter estimates so that the probability of observing what has actually been observed is as large as possible.
- Under most conditions, ML estimators are consistent, asymptotically efficient, and asymptotically normal (if all the assumptions are met).

Assumptions:

- The specified relationships are correct.
- The endogenous variables have a multivariate normal distribution, which implies
  - All variables are normally distributed.
  - All conditional expectation functions are linear.
  - All conditional variance functions are homoscedastic.

Parameter estimates are robust to violations of multivariate normality, but chi-squares may be too large and standard errors too small.

# Chi-Square Test

- If the specified model is correct, the chi-square test has approximately a chi-square distribution. The df is equal to the number of overidentifying restrictions.

- This statistic is a likelihood ratio chi-square comparing the fitted model with a saturated (just-identified) model that perfectly fits the data.  If the chi-square is large and the *p*-value is small, it's an indication that the model should be rejected.

- Note that, although this statistic is properly regarded as a test of the model, it is only testing the overidentifying restrictions.

- This test is sensitive to sample size.  With a large sample, it may be difficult to find any parsimonious model that passes this test.

# Other Measures of Fit

```
. estat gof, stats(all)
---------------------------------------------------------------------------
Fit statistic          |      Value   Description
-----------------------+---------------------------------------------------
Likelihood ratio       |
          chi2_ms(4)   |     48.074   model vs. saturated
             p > chi2  |      0.000
         chi2_bs(14)   |    800.942   baseline vs. saturated
             p > chi2  |      0.000
-----------------------+---------------------------------------------------
Population error       |
               RMSEA   |      0.138   Root mean squared error of
                       |                                  approximation
  90% CI, lower bound  |      0.104
         upper bound   |      0.174
             pclose    |      0.000   Probability RMSEA <= 0.05
```

*We want the RMSEA to be < .05.  Definitely don't want it to be > .10*

# Other Measures of Fit

```
-----------------------+---------------------------------------------------
Information criteria   |
               AIC     |  14925.739   Akaike's information criterion
               BIC     |  15004.304   Bayesian information criterion
-----------------------+---------------------------------------------------
Baseline comparison    |
               CFI     |      0.944   Comparative fit index
               TLI     |      0.804   Tucker-Lewis index
-----------------------+---------------------------------------------------
Size of residuals      |
              SRMR     |      0.034   Standardized root mean squared residual
                CD     |      0.468   Coefficient of determination
---------------------------------------------------------------------------
```

*We want the CFI and TLI to be close to 1, definitely not below .90.*

# Global Goodness of Fit Measures

We want a single number that measures the similarity of $\hat{\Sigma}$ and $S$ , the predicted covariance matrix (based on the model) and the observed covariance matrix.

As a general approach to model evaluation, LR chi-square may be too sensitive to sample size. Many alternative statistics have been proposed. Here are some that are reported by Stata.

**Tucker Lewis Index (TLI)**

Also known as Bentler & Bonnet's NonNormed Fit Index

Let $\chi_1^2$ be the chi-square for the fitted model and let $\chi_0^2$ be the chi-square for some baseline model, usually the "independence" model which says all the observed covariances are really 0. (Mplus only considers covariances among endogenous variables and between endogenous and exogenous variables).

$$TLI = \frac{\dfrac{\chi_0^2}{df_0} - \dfrac{\chi_1^2}{df_1}}{\dfrac{\chi_0^2}{df_0} - 1}$$

Adjusts for relative complexity of the two models. Can sometimes be greater than 1, a possible indication of "overfitting".

# Other Global Measures

**Comparative fit index**

$$CFI = \frac{\left(\chi_0^2 - df_0\right) - \left(\chi_1^2 - df_1\right)}{\chi_0^2 - df_0}$$

As with the TLI, models pay a penalty for more parameters. The formula can be greater than 1 or less than 0, in which case the CFI is simply set to 1 or 0.

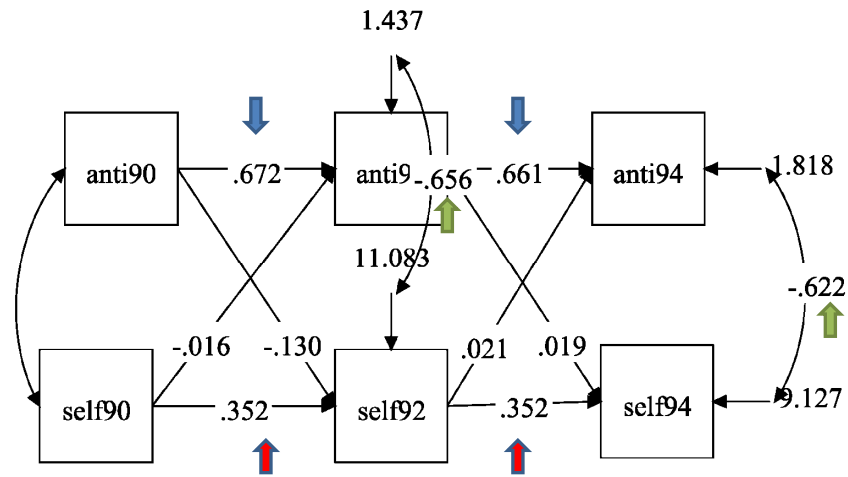**Root mean squared error of approximation:**

$$RMSEA = \sqrt{\frac{\dfrac{\chi_1^2}{df_1} - 1}{N - 1}}$$

Good models have an RMSEA of .05 or less. Models whose RMSEA is .10 or more have poor fit. One nice thing about this statistic is that you can get a confidence interval.

# Equality Constraints

If *T* > 2, consider constraining effects to be equal across time

1.437

anti90 —.672→ anti9 -.656 .661→ anti94 1.818

11.083

-.016 -.130 .021 .019

-.622

self90 —.352→ self92 —.352→ self94 9.127

Why?  Get smaller standard errors -> narrower confidence intervals and more powerful hypothesis tests.  Also easier to interpret.

# Stata Program with Constraints

```
use "C:\data\nlsy.dta", clear
sem (anti94 <- anti92@a self92@b) ///
    (anti92 <- anti90@a self90@b) ///
    (self94 <- anti92@c self92@d) ///
    (self92 <- anti90@c self90@d), ///
   cov(e.anti94*e.self94@e e.anti92*e.self92@e)
```

Compare this model with the last one:

Chi-square:  52.4 – 48.0 = 4.4

DF:  9 – 4 = 5

*p* = .49

So imposing the constraints did not significantly worsen the fit.