

A closer look at fixed effects regression in structural equation modeling using lavaan

Henrik Kenneth Andersen^a

^aChemnitz University of Technology, Institute of Sociology, Chair for Empirical Social Research, Thüringer Weg 9, 09126 Chemnitz, Germany

ARTICLE HISTORY

Compiled July 27, 2020

ABSTRACT

This article provides an in-depth look at fixed effects regression in the structural equation modeling (SEM) framework, specifically the application of fixed effects in the `lavaan` package for `R`. It is meant as a applied guide for researchers, covering the underlying model specification, syntax, and summary output. Online supplementary materials further discuss various common extensions to the basic fixed-effect model, demonstrating how to relax model assumptions, deal with measurement error in both the dependent and independent variables, and include time-invariant predictors in a type of hybrid fixed-/ random effects model.

KEYWORDS

Fixed effects, structural equation modeling, lavaan, R, panel analysis

1. Introduction

Several years ago, Curran and Bauer (2011) reflected positively on the growing use of panel studies in empirical social research. Some of the strengths of panel data are well-known, e.g., the ability to establish temporal precedence, increased statistical power and the reduction of potential alternative models. However, perhaps the greatest strength of panel data is that they allow for a more rigorous testing of substantive theories. Panel data, i.e., repeated measures of the same observed units (people, schools, firms, countries, etc.), allow researchers to decompose the error term into a part that stays constant within units and the part that changes over time. The part that does not change over time can be seen as the combined effect of all time-invariant influences (e.g., sex, date of birth, nationality) on the dependent variable. Fixed effects (FE) regression involves controlling for these time-invariant influences via a number of various methods. It thus accounts for a likely and common source of bias.

Structural equation modeling (SEM) is a popular regression framework. One of its main strengths is its flexibility. Not only can complex causal structures with multiple dependent variables be tested simultaneously, but in longitudinal (and, more generally, hierarchical) studies both time-varying and invariant predictors can be included, and effects can easily be allowed to vary over time. Thus researchers can allow for and study effects that increase or fade over time, or that appear only in specific periods. Beyond that, with the use of latent variables, SEM provides a way to deal with measurement error and get closer to the true underlying constructs of interest.

There are a number of articles describing basic concept of panel model regression, and FE regression in SEM (e.g., Allison 2011; Bollen and Brand 2010; Teachman et al. 2001). This article is intended as a *practical guide* for researchers looking for in-depth help with specifying FE models in SEM. It focuses on the `lavaan` (Rosseel 2012) package for R (R Core Team 2017). While `Mplus` (Muthén and Muthén 1998–2017) is arguably the most robust SEM software currently available (in terms of features like alignment, latent variable interactions, for example), the `lavaan` package has many benefits. First, like R it is open source and completely free. For researchers dipping their toes into SEM, there is no financial barrier to try, and no risk if they decide it

is not for them. Second, the implementation of `lavaan` in the larger `R` environment is an enormous advantage. Instead of poring over reams of plain text, copying out coefficients by hand, every part of the `lavaan` output is available as an object. This means that all aspects of the model, from fit indices, to coefficients and standard errors, to the model matrices, can be accessed and easily integrated into tables and plots. Furthermore, `R` can be used for a great deal of applications. It can be used to manage and manipulate as well as simulate data, perform symbolic algebra, run more traditional analyses (e.g., multiple regression, logistic regression, principal component analysis), etc. Once one is comfortable using `R`, there is no longer any need to switch between different software for data preparation and analysis.

The following article outlines the basic idea of panel regression, the particularities of panel regression SEM, and shows its implementation in `lavaan`. The focus on FE panel regression stems from the fact that the assumptions needed to justify other common models, e.g., random effects, are often implausible. As we will see, the practical differences between the two are, in any event, very small. Using simulated data, it demonstrates and annotates the code for the most basic FE model and provides an overview of the summary output. One of the main strengths of panel SEM compared to the more traditional methods of panel analysis is its flexibility. Therefore, a number of potential extensions to the basic model, including relaxing various assumptions, dealing with measurement error in both the independent and dependent variables, as well as the inclusion of time-invariant predictors in the form of a hybrid fixed-/ random effects model, are shown in detail in the form of online supplementary materials.

2. Panel models

To begin, let us start by reviewing a general panel model (Bollen and Brand 2010), also referred to as the ‘unobserved effects model’ (Wooldridge 2012; Croissant and Millo 2008) (we will return to this model in the online supplementary materials when we discuss loosening assumptions)

$$y_{it} = \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{z}_i\boldsymbol{\gamma} + \alpha_i + \varepsilon_{it} \tag{1}$$

where y_{it} is the dependent variable for unit i , $i = 1, \dots, N$ at time t , $t = 1, \dots, T$, \mathbf{x}_{it} is a $1 \times K$ vector of time-varying covariates (which could include a constant) linked to the dependent variable by the $K \times 1$ vector of coefficients $\boldsymbol{\beta}$. \mathbf{z}_i is a $1 \times M$ vector of time-invariant covariates linked to the dependent variable by the $M \times 1$ vector of coefficients in $\boldsymbol{\gamma}$, α_i represents the combined effect of all unobserved time-constant variables affecting the dependent variable and ε_{it} is the idiosyncratic error.

We can make stating some of the model assumptions easier by rewriting it in matrix notation

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\gamma} + \boldsymbol{\iota}_T \alpha_i + \boldsymbol{\varepsilon}_i$$

where \mathbf{y}_i and $\boldsymbol{\varepsilon}_i$ are $T \times 1$ vectors, \mathbf{X}_i and \mathbf{Z}_i are $T \times K$ and $T \times M$ matrices, respectively, $\boldsymbol{\iota}_T$ is a $T \times 1$ vector of ones and α_i is a scalar. $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are unchanged from Equation (1).

Consistency of the following models requires the assumption of strict exogeneity, although what constitutes strict exogeneity differs between the random and fixed effects setups. Each assumption will be discussed shortly. Apart from that, we typically make the following assumptions about this model (see, e.g., Wooldridge 2002; Schmidheiny 2019):

- Linearity: the model is linear in its parameters.
- Independence: the observations are independent across individuals (assured by random sampling in the cross-section), but not necessarily across time.
- The usual rank condition: we have more observations than independent variables and there is no perfect collinearity between any of the independent variables.

2.1. *Random effects*

Somewhat unintuitively, the modern approach to FE regression treats the unobserved effect α_i as a *random* variable, rather than a fixed parameter to be estimated for each of the cross sectional units (Wooldridge 2002). And in fact, this view carries over to the way in which FE regression is applied in SEM. So, in order to help facilitate the discussion on the application of FE-SEM later on, it makes sense to spend some

time discussing the basic idea of random effects regression (again, as we will see, the practical difference between the two models is very small).

For the random effects (RE) model, we define a composite error term: $\nu_{it} = \alpha_i + \varepsilon_{it}$ and rewrite the model in Equation (1) as

$$\begin{aligned} y_{it} &= \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{z}_i\boldsymbol{\gamma} + \nu_{it}, \text{ or} \\ \mathbf{y}_i &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\boldsymbol{\gamma} + \boldsymbol{\nu}_i, \end{aligned}$$

where $\boldsymbol{\nu}_i = \alpha_i \boldsymbol{\iota}_T + \boldsymbol{\varepsilon}_i$ and $\boldsymbol{\iota}_T$ is a $T \times 1$ vector of ones. The strict exogeneity assumption in the RE model implies

$$\begin{aligned} \mathbb{E}[\varepsilon_{it} | \mathbf{X}_i, \mathbf{z}_i, \alpha_i] &= 0, \\ \mathbb{E}[\alpha_i | \mathbf{X}_i, \mathbf{z}_i] &= \mathbb{E}[\alpha_i] = 0, \end{aligned}$$

where $\mathbf{X}_i = \mathbf{x}_{i1}, \dots, \mathbf{x}_{iT}$. For both parts, the assumption that the unconditional expectations are 0 is unproblematic as long as a constant is included in the regression. The first part says the idiosyncratic errors at each timepoint are assumed to be independent of the explanatory variables at *all* timepoints which is stronger than just assuming that they are *contemporaneously* independent. This implies that they are also uncorrelated, i.e., $\mathbb{E}[\mathbf{x}_{is}^\top \varepsilon_{it}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{z}_i^\top \varepsilon_{it}] = \mathbf{0}$, $\forall s, t = 1, \dots, T$ (Wooldridge 2002; Brüderl and Ludwig 2015). We assume the idiosyncratic errors are further independent of the individual effects, which implies $\mathbb{E}[\alpha_i \varepsilon_{it}] = 0$.

The second part is the potentially controversial assumption: it states that the individual effects are uncorrelated with the independent variables. We can use an intuitive concrete example to show why this is often controversial: If we are interested in the question of whether married men earn more than unmarried men, then the second part of the strict exogeneity assumption means that a man's marriage status would have to be uncorrelated with all the time-invariant characteristics that could potentially make that man an attractive marriage candidate in the first place; e.g., looks, personality, family's status, profession, etc. (Brüderl and Ludwig 2015).¹

From what we have discussed so far, the $T \times T$ covariance matrix of the errors $\boldsymbol{\Omega}_i = \mathbb{E}[\boldsymbol{\nu}_i \boldsymbol{\nu}_i^\top]$ can be constructed. However, the standard random effects model adds

¹ Assuming, for the sake of argument, that these characteristics are constant over time.

the additional assumptions

$$\begin{aligned}\mathbb{E}[\varepsilon_{it}^2 | \mathbf{X}_i, \mathbf{z}_i, \alpha_i] &= \mathbb{E}[\varepsilon_{it}^2] = \sigma_\varepsilon^2, \quad t = 1, \dots, T, \\ \mathbb{E}[\varepsilon_{it}\varepsilon_{is} | \mathbf{X}_i, \mathbf{z}_i, \alpha_i] &= \mathbb{E}[\varepsilon_{it}\varepsilon_{is}] = 0, \quad \forall t \neq s\end{aligned}$$

i.e., the idiosyncratic errors are conditionally homoscedastic and serially uncorrelated and

$$\mathbb{E}[\alpha_i^2 | \mathbf{X}_i, \mathbf{z}_i] = \mathbb{E}[\alpha_i^2] = \sigma_\alpha^2$$

i.e., the individual effects are conditionally homoscedastic (they are necessarily serially correlated as long as $\sigma_\alpha^2 > 0$). From that, we arrive at the typical random effects structure of the $NT \times NT$ matrix $\mathbf{\Omega}$:

$$\mathbf{\Omega} = \begin{pmatrix} \mathbf{\Omega}_1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{\Omega}_i & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & \mathbf{\Omega}_N \end{pmatrix}$$

with $T \times T$ typical elements

$$\mathbf{\Omega}_i = \begin{pmatrix} \sigma_\nu^2 & \sigma_\alpha^2 & \dots & \sigma_\alpha^2 \\ \sigma_\alpha^2 & \sigma_\nu^2 & \dots & \sigma_\alpha^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_\alpha^2 & \sigma_\alpha^2 & \dots & \sigma_\nu^2 \end{pmatrix} \quad (2)$$

where $\sigma_\nu^2 = \sigma_\alpha^2 + \sigma_\varepsilon^2$ (Wooldridge 2002; Schmidheiny 2019). This means that in the conditional covariance matrix of the errors, given the time-varying and -invariant covariates, units over time will be correlated due to the individual effects. We should keep the covariance structure of the errors in mind as it will help make sense of the use of latent variables to decompose the dependent variable into between- and within-variance components, discussed below in Section 3.

Estimation of the RE model can be done using feasible generalized least squares (GLS) in which the two unknowns in $\mathbf{\Omega}$, σ_α^2 and σ_ν^2 , are first estimated using pooled

ordinary least squares (pooled OLS or POLS), where²

$$\hat{\sigma}_\nu^2 = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \hat{\nu}_{it}^2,$$

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{NT - N} \sum_{i=1}^N \sum_{t=1}^T (\hat{\nu}_{it} - \bar{\hat{\nu}}_i)^2$$

and $\hat{\nu}_{it} = y_{it} - \mathbf{x}_{it}\hat{\boldsymbol{\beta}}_{POLS} - \mathbf{z}_i\hat{\boldsymbol{\gamma}}_{POLS}$, $\bar{\hat{\nu}}_i = T^{-1} \sum_{t=1}^T \hat{\nu}_{it}$ and $\hat{\sigma}_\alpha^2 = \hat{\sigma}_\nu^2 - \hat{\sigma}_\varepsilon^2$. Then, the coefficients are estimated using those estimates in the variance matrix $\hat{\boldsymbol{\Omega}}$:

$$\begin{pmatrix} \hat{\boldsymbol{\beta}}_{RE} \\ \hat{\boldsymbol{\gamma}}_{RE} \end{pmatrix} = (\mathbf{W}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{W})^{-1} \mathbf{W}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{y}$$

where $\mathbf{W} = \begin{pmatrix} \mathbf{X} & \mathbf{Z} \end{pmatrix}$ and \mathbf{X} is $NT \times K$ and \mathbf{Z} is $NT \times M$ and \mathbf{y} is $NT \times 1$.

In practice, however, computational problems can arise with large cross-sectional samples, where it can become difficult to invert the $\hat{\boldsymbol{\Omega}}$ matrix. One solution is to use ‘partial-demeaning’ to transform the data before performing simple POLS:

$$(y_{it} - \theta \bar{y}_i) = (\mathbf{x}_{it} - \theta \bar{\mathbf{x}}_i) \boldsymbol{\beta} + (\mathbf{z}_i - \theta \bar{\mathbf{z}}_i) \boldsymbol{\gamma} + (\varepsilon_{it} - \theta \bar{\varepsilon}_i) \quad (3)$$

where $\theta = 1 - [\sigma_\alpha^2 / (\sigma_\alpha^2 + T\sigma_\varepsilon^2)]^{1/2}$, and $\bar{y}_i = T^{-1} \sum_{t=1}^T y_{it}$, $\bar{\mathbf{x}}_i = T^{-1} \sum_{t=1}^T \mathbf{x}_{it}$, $\bar{\mathbf{z}}_i = T^{-1} \sum_{t=1}^T \mathbf{z}_i$ and $\bar{\varepsilon}_i = T^{-1} \sum_{t=1}^T \varepsilon_{it}$ (Croissant and Millo 2008).

The RE model can also be estimated in the maximum likelihood framework, where in the associated literature on panel models are generally referred to as either mixed models, hierarchical models or longitudinal models. The typical RE model discussed here is the equivalent to a mixed model with random intercepts and fixed slopes (Croissant and Millo 2008). Under the assumption of normality, along with homoscedasticity and serially uncorrelated errors, the maximum likelihood estimator is the same as the OLS estimator. For more on the topic of mixed models, see for example Bates et al. (2015), Bates (2010).

²Normally a degrees-of-freedom correction is applied by subtracting off the number of independent variables that is negligible in large N samples. It is ignored here for the sake of simplicity.

2.2. Fixed effects

For fixed effects, we assume that the individual effects are *not independent* of the model covariates, i.e., $\mathbb{E}[\alpha_i | \mathbf{X}_i, \mathbf{z}_i] \neq \mathbb{E}[\alpha_i] \neq 0$. Under this assumption, grouping the individual effects in with the composite error will cause the coefficients of interest, here specifically β to be inconsistent (Wooldridge 2002, 2012). We write the model therefore again in terms of the general panel model in Equation (1) with separate individual effects and idiosyncratic error terms. In order to drop assumptions involving the individual effects, a number of methods are available (e.g., differencing, least squares dummy variable regression), but the most common approach is to *demean* the equation (Bröderl and Ludwig 2015). Demeaning is the same as the transformation applied in Equation (3) in the special case where $\theta = 1$. I.e., demeaning involves subtracting the per-unit, over-time average from each of the model terms, i.e.,

$$\begin{aligned} (y_{it} - \bar{y}_i) &= (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\beta + (\mathbf{z}_i - \bar{\mathbf{z}}_i)\gamma + (\alpha_i - \bar{\alpha}_i) + (\varepsilon_{it} - \bar{\varepsilon}_i) \\ \ddot{y}_{it} &= \ddot{\mathbf{x}}_{it}\beta + \ddot{\varepsilon}_{it}, \text{ or} \\ \ddot{\mathbf{y}}_i &= \ddot{\mathbf{X}}_i\beta + \ddot{\varepsilon}_i \end{aligned} \tag{4}$$

where the over-time averages are calculated the same as above, and the variables with the dots above them represent the demeaned versions. Because the average of something that does not change is that thing itself, the individual effects, along with any time-invariant predictors, get wiped out by the demeaning. This means that no assumptions about the relatedness of the model covariates and the unit-specific portion of the error are needed. Consistency of the estimates is related solely to the strict exogeneity assumption imposed on the idiosyncratic errors, i.e., $\mathbb{E}[\ddot{\varepsilon}_{it} | \ddot{\mathbf{x}}_{it}] = \mathbb{E}[\ddot{\varepsilon}_{it}] = 0$ which also implies $\mathbb{E}[\ddot{\mathbf{x}}_{is}^\top \ddot{\varepsilon}_{it}] = \mathbf{0}$, $\forall s, t = 1, \dots, T$ (Bröderl and Ludwig 2015; Wooldridge 2002).

Having demeaned the data, the typical FE estimator is POLS on the transformed data

$$\beta_{FE} = (\ddot{\mathbf{X}}^\top \ddot{\mathbf{X}})^{-1} \ddot{\mathbf{X}}^\top \ddot{\mathbf{y}}$$

(Bröderl and Ludwig 2015). The downside to this approach is that no time-invariant

predictors can be included in the model. However, there are alternative approaches in the random effects and mixed model frameworks that allow them to be included. These models are sometimes referred to as ‘within-between’ or ‘hybrid’ models, often based on the Chamberlain (1980) and Mundlak (1978) approaches, see for example Bell, Fairbrother, and Jones (2018); Allison (2011); Schunck (2013); Enders and Tofighi (2007). In the online appendix, it will be discussed how to also get around this restriction using SEM.

3. Fixed effects in structural equation modeling

Moving from the conventional methods outlined above to SEM, we must state the FE model in a different way. We turn to latent variables to account for time-invariant unobserved heterogeneity. In fact, besides accounting for measurement error and the representation of abstract hypothetical concepts, unobserved heterogeneity has historically been one of the main uses of latent variables in SEM (Skrondal and Rabe-Hesketh 2004).

3.1. *Modeling time-invariant unobserved heterogeneity as a latent variable*

We first need to convert the data from stacked, long-format vectors of length NT into T individual vectors of length N . To see why this is necessary, consider what effect this has on the vector of responses y_{it} . Let us, for a minute ignore any covariates and focus just on the dependent variable (a so-called ‘intercept-only’ or ‘null’ model) so that we have $y_{it} = \alpha_i + \varepsilon_{it}$. When we convert the data to wide-format, we get T individual equations,

$$\mathbf{y}_t = \boldsymbol{\alpha} + \boldsymbol{\varepsilon}_t$$

$$\begin{bmatrix} y_{1t} \\ y_{2t} \\ \vdots \\ y_{Nt} \end{bmatrix} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_N \end{bmatrix} + \begin{bmatrix} \varepsilon_{1t} \\ \varepsilon_{2t} \\ \vdots \\ \varepsilon_{Nt} \end{bmatrix} \quad (5)$$

for each $t = 1, 2, \dots, T$. Because the idiosyncratic errors are assumed to be uncorrelated

across units and across time, the covariance between any two of the new wide vectors $\text{Cov}(y_{ti}, y_{si}) = \text{Var}(\alpha_i)$, $t \neq s$. Otherwise, when $t = s$, the covariance $\text{Cov}(y_{ti}, y_{ti}) = \text{Var}(\alpha_i) + \text{Var}(\varepsilon_{ti})$. This is the structure we saw above in a typical element of $\mathbf{\Omega}$.

And in fact this is exactly how a latent variable is used to account for time-invariant unobserved heterogeneity. The dependent variable at each timepoint is regressed onto the latent variable, see Figure 1. Here, the regression weights or ‘factor loadings’ are fixed to one to represent our assumption that the effect of the time-invariant unobserved heterogeneity is constant over time.³ It also means that the estimated variance of the latent variable is equal to the *average covariance between the wide-format columns of the dependent variable over time*. If $y_{it} = \alpha_i + \varepsilon_{it}$ is the true data generating process, then the relationship between two units over time is just $\text{Var}(\alpha)$, regardless of the time distance. Referring back to the random effects structure of $\mathbf{\Omega}_i$ in Equality (2) for a generic unit i , we see the covariance on all of the off-diagonals is σ_α^2 . And, as we know, the average of something that does not change is that thing itself. I.e., if $T(T-1)/2 = h$ is the number of elements on either the upper- or lower triangle of $\mathbf{\Omega}_i$, then we have $h^{-1} \sum_{i=1}^h \sigma_\alpha^2 = \frac{h\sigma_\alpha^2}{h} = \sigma_\alpha^2$.

To elaborate on this concept some more, consider the following matrix equation of the variances and the nonredundant covariances in a three-wave intercept-only model that follows directly from Equation (5) (assuming $\text{Cov}(\varepsilon_{ti}, \varepsilon_{si}) = 0$, $t \neq s$), and which we can solve easily with least squares:

$$\mathbf{Ax} = \mathbf{b}$$

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \psi \\ \phi_1 \\ \phi_2 \\ \phi_3 \end{pmatrix} = \begin{pmatrix} \text{Var}(y_1) \\ \text{Cov}(y_2, y_1) \\ \text{Cov}(y_3, y_1) \\ \text{Var}(y_2) \\ \text{Cov}(y_3, y_2) \\ \text{Var}(y_3) \end{pmatrix}$$

where $\psi = \text{Var}(\alpha)$, $\phi_t = \text{Var}(\varepsilon_t)$. We can solve this equation to show

³In fact, the initial FE-SEM setup shown in the main article mimics the POLS methods described above in that it assumes constant effects and error variances over time. These assumptions can be loosened and tested, as will be shown in the supplementary materials. For now, for the sake of simplicity and comparability, we retain the assumptions associated with the ‘pooled’ models for the most part.

$$\begin{aligned}
\mathbf{A}\mathbf{x} &= \mathbf{b} \\
(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{A}\mathbf{x} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\
\hat{\mathbf{x}} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\
\begin{pmatrix} \hat{\psi} \\ \hat{\phi}_1 \\ \hat{\phi}_2 \\ \hat{\phi}_3 \end{pmatrix} &= \begin{pmatrix} \frac{1}{3} \text{Cov}(y_2, y_1) + \frac{1}{3} \text{Cov}(y_3, y_1) + \frac{1}{3} \text{Cov}(y_3, y_2) \\ \text{Var}(y_1) - \hat{\psi} \\ \text{Var}(y_2) - \hat{\psi} \\ \text{Var}(y_3) - \hat{\psi} \end{pmatrix}.
\end{aligned}$$

So if, in fact the covariance between any two wide-format columns of y is $\text{Cov}(y_t, y_s) = \text{Var}(\alpha)$, $\forall s \neq t$, then $\hat{\psi} = \frac{1}{3} \text{Var}(\alpha) + \frac{1}{3} \text{Var}(\alpha) + \frac{1}{3} \text{Var}(\alpha) = \frac{3 \text{Var}(\alpha)}{3} = \text{Var}(\alpha)$. This shows that if our assumption about the underlying DGP is correct, i.e., $y_{it} = \alpha_i + \varepsilon_{it}$ and $\text{Cov}(\varepsilon_t, \varepsilon_s) = 0$, $\forall t \neq s$, then the estimated variance of α is just what it should be: the average covariance between units of y over time. Once we add in observed covariates, the estimated covariance of α then become the *conditional* covariance of y over time, given those covariates.

3.2. Model notation

Having explained how a latent variable is used to estimate the individual effects, we need to state the FE model using SEM-compatible matrix notation. There are a number of different model notations (see, for example Bollen (1989) for an overview), but the one that will serve us best is one that was proposed by Graff (1979):

$$\begin{aligned}
\mathbf{y}^+ &= \mathbf{\Lambda}_y^+ \boldsymbol{\eta}^+, \\
\boldsymbol{\eta}^+ &= \mathbf{B} \boldsymbol{\eta}^+ + \boldsymbol{\zeta}^+,
\end{aligned}$$

where $\boldsymbol{\eta}^+ = (\mathbf{y}, \mathbf{x}, \boldsymbol{\eta}, \boldsymbol{\xi})^\top$, $\boldsymbol{\zeta}^+ = (\boldsymbol{\varepsilon}, \boldsymbol{\delta}, \boldsymbol{\zeta}, \boldsymbol{\xi})^\top$, $\mathbf{y}^+ = (\mathbf{y}, \mathbf{x})^\top$. \mathbf{y} is a vector of observed dependent variables and \mathbf{x} is a vector of observed independent variables. $\boldsymbol{\eta}$ is a vector of the latent dependent variables and $\boldsymbol{\xi}$ is a vector of latent independent variables. $\boldsymbol{\varepsilon}$ and $\boldsymbol{\delta}$ are vectors of the errors of the observed dependent and independent variables, respectively, and $\boldsymbol{\zeta}$ is a vector of the errors, or disturbances, of the latent variables. Notice the $^+$ symbol is just meant to differentiate the vectors with them from those without them. That means, $\boldsymbol{\eta}^+$ is a vector that holds the observed and latent variables,

both dependent (in SEM they are referred to as ‘endogenous’) and independent (i.e., ‘exogenous’)⁴, ζ^+ holds the errors for the observed variables and the disturbances of the latent variables. \mathbf{y}^+ holds just the observed variables, both dependent and independent, and $\mathbf{\Lambda}_y^+$ is a matrix of ones and zeros that selects the observed variables from $\boldsymbol{\eta}^+$. Lastly, \mathbf{B} is a matrix that holds the regression coefficients.

If we say that p and q stand for the number of observed dependent and independent variables, respectively, and m and n stand for the number of latent dependent and independent variables, respectively, then $\boldsymbol{\eta}^+$ and ζ^+ are $p + q + m + n$, \mathbf{y}^+ is $p + q$, $\mathbf{\Lambda}_y^+$ is $(p + q) \times (p + q + m + n)$ and \mathbf{B} is $(p + q + m + n) \times (p + q + m + n)$ (Bollen 1989).

This notation may be confusing at first, but it has advantages. First, it allows us the flexibility we need for the models. For example, it allows observed \mathbf{x} to directly influence observed \mathbf{y} (more common notation assumes that substantive effects occur only between latent variables, observed ones are only used as indicators, see for example Bollen (1989), Kline (2016)). It also allows $\boldsymbol{\xi}$, i.e., any latent exogenous variables, to influence \mathbf{y} directly. These two scenarios cover the traditional FE model with observed variables, and one in which latent variables are used to account for measurement error in the independent variables. It is also consistent with the notation used for these models in `lavaan`. In fact, `lavaan` switches automatically between matrix notations depending on the specified model. That means the matrix representation of the model one sees if they type in `lavInspect(model, what = "est")` after specifying their model in `lavaan` will match up with the notation used here.⁵ It does have a potential disadvantage however. Besides being less intuitive than the typical $\boldsymbol{\eta} = \mathbf{B}\boldsymbol{\eta} + \mathbf{\Gamma}\boldsymbol{\xi} + \zeta$ notation, it means that by including the observed covariates in the stacked long vector \mathbf{y}^+ , they are treated as another response variable with variances and covariances to be estimated by the model (instead of just using the sample statistics). This means that the assumption of multivariate normality (otherwise just imposed on the dependent variables) also applies to the independent ones (Skrondal and Rabe-Hesketh 2004, p. 75). This can be problematic for noncontinuous independent variables like sex,

⁴For our purposes, the terms endogenous and dependent, on the one hand, and exogenous and independent, on the other, can be used interchangeably.

⁵In `Mplus`, the model matrices can be requested by including `OUTPUT: TECH1` in the input file.

nationality dummies, marriage status (married/unmarried), etc. See Skrondal and Rabe-Hesketh (2004) for more on this topic.

[Figure 1 about here.]

Let us, however, make things more concrete and take a look at a simple, three-wave version of the typical FE-SEM using this notation (shown graphically in Figure 1). For that, we have the following matrix notation (with labels on the outside of the matrices):

$$\mathbf{y}^+ = \mathbf{\Lambda}_y^+ \boldsymbol{\eta}^+ \quad (6)$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{matrix} & \begin{matrix} y_1 & y_2 & y_3 & x_1 & x_2 & x_3 & \alpha \end{matrix} \\ \begin{matrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{pmatrix},$$

$$\boldsymbol{\eta}^+ = \mathbf{B}\boldsymbol{\eta}^+ + \boldsymbol{\zeta}^+ \quad (7)$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{pmatrix} = \begin{matrix} & \begin{matrix} y_1 & y_2 & y_3 & x_1 & x_2 & x_3 & \alpha \end{matrix} \\ \begin{matrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{matrix} & \begin{pmatrix} 0 & 0 & 0 & \beta & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & \beta & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & \beta & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ x_1 = \delta_1 \\ x_2 = \delta_2 \\ x_3 = \delta_3 \\ \alpha = \xi \end{pmatrix}.$$

Notice in $\boldsymbol{\zeta}^+$, for the independent variables we could either write, for example x_t or δ_t . As mentioned above, this is due to the model notation treating the independent variables like dependent variables with variances/covariances to be estimated. For the sake of simplicity, we will ignore this subtlety and refer to the observed variable from now on, keeping in mind that if the multivariate normality assumption holds, the estimated statistics will likely be sufficiently close to the sample ones for it to not make much of a difference.

Admittedly, Equations (6) and (7) may not look like much yet. We can remedy this by first putting the equation for $\boldsymbol{\eta}^+$ in reduced form, i.e., by getting rid of the

dependent variable on the r.h.s.:

$$\begin{aligned}\eta^+ &= B\eta^+ + \zeta^+ \\ \eta^+ - B\eta^+ &= \zeta^+ \\ (I - B)\eta^+ &= \zeta^+ \\ \eta^+ &= (I - B)^{-1}\zeta^+, \end{aligned}$$

where I is the identity matrix. By substituting this back into the equation for the observed variables we get $\mathbf{y}^+ = \Lambda_y^+[(I - B)^{-1}\zeta^+]$, which works out to:

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} \alpha + \beta x_1 + \varepsilon_1 \\ \alpha + \beta x_2 + \varepsilon_2 \\ \alpha + \beta x_3 + \varepsilon_3 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix},$$

which is of course exactly what we should expect given Equation (4)⁶.

3.3. Assumptions

What essentially differentiates an FE from an RE model is our assumption concerning the relationship between the unobserved individual effects and the model covariates (Bollen and Brand 2010). The FE model assumes that $\mathbb{E}[\alpha_t] \neq 0$. As such, if we fail to control for the correlation of the covariate and the time-invariant part of the error, then the coefficient of interest, here β , will be biased. Our assumption regarding whether the individual effects are correlated with the model covariates occurs in $\mathbb{E}[\zeta^+\zeta^{+\top}] = \Psi$, the covariance matrix of the errors

$$\begin{aligned}\mathbf{y}^+\mathbf{y}^{+\top} &= \mathbb{E}[(\Lambda_y^+(I - B)^{-1}\zeta^+)(\Lambda_y^+(I - B)^{-1}\zeta^+)^{\top}] \\ &= \mathbb{E}[(\Lambda_y^+(I - B)^{-1}\zeta^+)(\zeta^{+\top}(I - B)^{-1\top}\Lambda_y^{+\top})] \\ &= \Lambda_y^+(I - B)^{-1}\mathbb{E}[\zeta^+\zeta^{+\top}](I - B)^{-1\top}\Lambda_y^{+\top} \\ &= \Lambda_y^+(I - B)^{-1}\Psi(I - B)^{-1\top}\Lambda_y^{+\top}.\end{aligned}$$

In the case of an FE model, Ψ will reflect our belief that the individual effects are correlated with the model covariates, here again for demonstration the three-wave

⁶We can use the `sympy` package in `python` to verify and show the steps for this and other examples, see the supplementary materials for the code.

model:

$$\Psi = \mathbb{E} \begin{pmatrix} \varepsilon_1 & \varepsilon_2 & \varepsilon_3 & x_1 & x_2 & x_3 & \alpha \\ \varepsilon_1 & \varepsilon_1^2 & 0 & 0 & 0 & 0 & 0 \\ \varepsilon_2 & 0 & \varepsilon_2^2 & 0 & 0 & 0 & 0 \\ \varepsilon_3 & 0 & 0 & \varepsilon_3^2 & 0 & 0 & 0 \\ x_1 & 0 & 0 & 0 & x_1^2 & 0 & 0 \\ x_2 & 0 & 0 & 0 & x_2x_1 & x_2^2 & 0 \\ x_3 & 0 & 0 & 0 & x_3x_1 & x_3x_2 & x_3^2 \\ \alpha & 0 & 0 & 0 & \alpha x_1 & \alpha x_2 & \alpha x_3 \\ & & & & & & \alpha^2 \end{pmatrix}.$$

Knowing this, we can work out the equation for the coefficient of interest, β . For the sake of simplicity, assume here and throughout mean-centered variables:

$$\begin{aligned} \text{Cov}(y_t, x_t) &= \mathbb{E}[y_t x_t] \\ &= \mathbb{E}[(\alpha + \beta x_t + \varepsilon_t)x_t] \\ &= \mathbb{E}[\alpha x_t + \beta x_t^2 + \varepsilon_t x_t] \\ &= \text{Cov}(\alpha, x_t) + \beta \text{Var}(x_t) \\ \hat{\beta}_{FE-SEM} &= \frac{\text{Cov}(y_t, x_t) - \text{Cov}(\alpha, x_t)}{\text{Var}(x_t)}. \end{aligned}$$

This should make intuitive sense. From the observed covariance between the dependent and the independent variable, we are partialling out the part that is due to the covariance between the independent variable and the individual effects per unit, and then dividing by the variance of the independent variable, as usual. For the RE model, we assume $\mathbb{E}[\alpha x_t] = 0$ and the equation reduces to $\hat{\beta}_{RE-SEM} = \text{Cov}(y_t, x_t) / \text{Var}(x_t)$. The rest of the model-implied covariance matrix results from $\mathbf{y}^+ \mathbf{y}^{+\top}$.

4. Fixed effects in lavaan

The package `lavaan` needs to be installed once with `install.packages("lavaan")`. To be able to use it, we need to load it for every new R session:

```
library( lavaan)
```

For users unfamiliar with R, SEM analyses can be carried out with almost no knowledge of the language. Typically, someone unfamiliar with R would prepare their data using some other statistical software, and then save the intended dataset as a `.csv`,

.xlsx, .dta, .sav, etc. file. The user must then import the data, preferably as a dataframe, and the rest occurs using the `lavaan` syntax.⁷

Specifying the most basic fixed effects model, like the one shown in Bollen and Brand (2010) (the same model as Equation (4) but with just one time-varying predictor) involves four components. First, we define the latent individual effects variable using the `=~` ‘measured by’ or ‘manifested by’ (Rosseel 2012) operator at the same time constraining the factor loadings at each timepoint to one. I will call the latent variable **a** to stand for α . Constraining all of the factor loadings to one reflects our implicit assumption that the combined effect of the unit-specific unobserved factors is constant over time. This is the default behaviour of traditional POLS-based approaches to FE that use the stacked long-format data.

```
a =~ 1*y1 + 1*y2 + 1*y3 + 1*y4 + 1*y5
```

Second, we regress the dependent variable on the independent variable using the `~` regression operator. With stacked, long-format data, only one regression coefficient is estimated over all observed timepoints. To have our FE-SEM model mimic this behaviour, we need to constrain the the estimated coefficient to equal over time. We do so by adding the same label to the regression coefficient at every time point. We will use the label **b** (this label was chosen arbitrarily, we could have used any letter or string of characters) and have it act as an equality constraint for the regression coefficient of interest β :

```
y1 ~ b*x1
y2 ~ b*x2
y3 ~ b*x3
y4 ~ b*x4
y5 ~ b*x5
```

The key to a FE model, as opposed to an RE model are our assumptions about the relatedness of our covariate and the individual effects, i.e., $\mathbb{E}[x_t\alpha]$. For an FE model, we want to partial out any potential covariance between the independent variable and the individual effects. This accounts for any linear relationship between x_t and

⁷There are many online tutorials for importing data in various formats, see, for example some from datacamp or Quick-R, or any of the many posts on stackoverflow.

the unit-specific characteristics influencing the dependent variable. Further, allowing unrestricted covariances between the independent variable itself over time will not affect how the coefficient β is estimated, but will have an effect on the standard errors. To mimic the behaviour of a conventional FE model, we allow the independent variable to be correlated with the individual effects and itself over time. Covariances (including covariances between a variable and itself, i.e., variances) are specified using the `~~` operator:

```
a  ~~ x1 + x2 + x3 + x4 + x5
x1  ~~ x2 + x3 + x4 + x5
x2  ~~ x3 + x4 + x5
x3  ~~ x4 + x5
x4  ~~ x5
```

The last component of our code involves the variances of the residuals. This component is optional, but we can constrain the residual variances to be equal over time to again mimic the behaviour of a conventional FE model using POLS on stacked data. Here, again, we use labels to make equality constraints. Because y_t is endogenous, the `~~` operator specifies the variances of *residuals*, i.e., ε_t .

```
y1  ~~ e*y1
y2  ~~ e*y2
y3  ~~ e*y3
y4  ~~ e*y4
y5  ~~ e*y5
```

5. A simulated example

To demonstrate the application of FE models in SEM, a dataset can be simulated that embodies the FE assumptions. Again, the code for data simulation can be found in the online supplementary materials.

To show that the latent individual effects variables represent the *combined* effect of all time-invariant characteristics, the dependent variable will be influenced by two separate unit-specific variables, which we can call α_1 and α_2 . We will construct the simulated data such that the independent variable is correlated with both of the time-

invariant variables. This means that approaches that fail to account for this confounding influence, such as POLS or RE, will be biased.

The wide-format equations for the data generating process can be described as:

$$\begin{aligned} \mathbf{x}_t &= \boldsymbol{\alpha}_1 \beta_{x_t, \alpha_1} + \boldsymbol{\alpha}_2 \beta_{x_t, \alpha_2} + \boldsymbol{\delta}_t, \\ \mathbf{y}_t &= \mathbf{x}_t \beta_{y_t, x_t} + \boldsymbol{\alpha}_1 \beta_{y_t, \alpha_1} + \boldsymbol{\alpha}_2 \beta_{y_t, \alpha_2} + \boldsymbol{\varepsilon}_t \end{aligned}$$

where, for the sake of simplicity, $\boldsymbol{\alpha}_1$, $\boldsymbol{\alpha}_2$, $\boldsymbol{\delta}_t$ and $\boldsymbol{\varepsilon}_t$ are $\sim N(0, 1)$.

For the following example, a sample size of 1,000, observed over five waves, was chosen. The unique variance of \mathbf{x} , as well as both the individual-effect variables is also $\sim N(0, 1)$. The coefficient of interest, $\beta_{y, x}$ is set to be equal to 0.3. A correlation between \mathbf{x} and the individual effects is induced through $\beta_{x, \alpha_1} = 0.85$ and $\beta_{x, \alpha_2} = 0.50$. With the variances above set to one, the covariances will be roughly $\text{Cov}(x_t, \alpha_1) = 0.85$ and $\text{Cov}(x_t, \alpha_2) = 0.5$. The dependent variable is also influenced by the individual effects variables with $\beta_{y_t, \alpha_1} = 0.75$ and $\beta_{y_t, \alpha_2} = 0.45$. These values were chosen arbitrarily.

Now, we run the FE-SEM in `lavaan`.

```
fe_sem <- '
# Define individual effects variable
a =~ 1*y1 + 1*y2 + 1*y3 + 1*y4 + 1*y5
# Regressions, constrain coefficient to be equal over time
y1 ~ b*x1
y2 ~ b*x2
y3 ~ b*x3
y4 ~ b*x4
y5 ~ b*x5
# Allow unrestricted correlation between eta and covariates
a ~~ x1 + x2 + x3 + x4 + x5
x1 ~~ x2 + x3 + x4 + x5
x2 ~~ x3 + x4 + x5
x3 ~~ x4 + x5
x4 ~~ x5
# Constrain residual variances to be equal over time
y1 ~~ e*y1
y2 ~~ e*y2
y3 ~~ e*y3
y4 ~~ e*y4
y5 ~~ e*y5
',
fe_sem.fit <- sem(model = fe_sem,
```

```
data = dfw,
estimator = "ML")
```

We can get a summary of the model with `summary()`. The first portion of the summary output gives an overview of some basic information and fit statistics. The maximum likelihood estimator is the default, so it did not have to be explicitly selected in the fitting function call. Other estimators are available, including generalized and unweighted least squares (GLS and ULS, respectively), robust standard errors maximum likelihood (MLM) and several others (see the lavaan online tutorial for more).

This part of the summary output also tells us that the analysis is based on 1,000 observations (missings would be shown here as well if there were any), and that the χ^2 statistic is 30.138 based on 32 degrees of freedom (55 observed covariances minus 1 error variance, 1 coefficient, 1 latent variable variance, 5 exogenous variable variances and 15 covariances for $55 - 23 = 32$ df). The p-value on the χ^2 statistic is not significant with $p = 0.561$ which tells us the differences between the model-implied and observed covariance matrices are likely due to chance, and that the model fits the data well (given how the data was generated, it would be surprising if this were not the case). Other fit measures including typical comparative fit indices can be requested by either adding `fit.measures = TRUE` as a secondary argument to the `summary()` call, or by asking for a complete list of all available fit statistics using `lavInspect(model, "fit")` where `model` stands for the name of the fitted model, in this case `fe_sem.fit`.

```
summary(fe_sem.fit)
```

```
## lavaan 0.6-6 ended normally after 37 iterations
##
##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of free parameters      31
##      Number of equality constraints    8
##
##      Number of observations          1000
##
## Model Test User Model:
##
##      Test statistic                  30.138
##      Degrees of freedom              32
```

```
## P-value (Chi-square) 0.561
##
## Parameter Estimates:
##
## Standard errors      Standard
## Information          Expected
## Information saturated (h1) model  Structured
...
```

Next the summary output shows the measurement models for the latent variables, if any. In this case the latent variable α is measured by each of the five observed dependent variables with factor loadings fixed to 1.0.

```
...
## Latent Variables:
##           Estimate Std.Err z-value P(>|z|)
## a =~
## y1      1.000
## y2      1.000
## y3      1.000
## y4      1.000
## y5      1.000
...
```

The regressions are shown next. Here, because we have constrained the regression coefficients to be equal over time (the equality constraint label (b) is listed to the left of the estimates), the estimate of $\beta = 0.294$ (0.016) is repeated five times. The corresponding z- and p-values show that the coefficient is, unsurprisingly, significant.

```
...
## Regressions:
##           Estimate Std.Err z-value P(>|z|)
## y1 ~
## x1      (b)  0.294  0.016  18.809  0.000
## y2 ~
## x2      (b)  0.294  0.016  18.809  0.000
## y3 ~
## x3      (b)  0.294  0.016  18.809  0.000
## y4 ~
## x4      (b)  0.294  0.016  18.809  0.000
## y5 ~
## x5      (b)  0.294  0.016  18.809  0.000
...
```

Next, the covariance estimates are listed. First, the covariances between the latent

individual effects variable and the independent variable over time are shown, and then the covariances between the independent variable with itself over time.

One should always take care to double-check that there are no unintended covariances listed here. Like `Mplus`, the `lavaan` package estimates some covariances per default, without the user explicitly having to add them to the model syntax. For example, covariances between latent variables are estimated per default. If one does not wish for them to covary, it must be explicitly stated, e.g., with `f1 ~~ 0*f2`, assuming the latent variables are called `f1` and `f2`, or by overriding the default behaviour for the entire model by adding `orthogonal = TRUE` (which sets the correlation between all latent variables to zero) to the fitting call.⁸

```
...
## Covariances:
##
##      Estimate  Std.Err  z-value  P(>|z|)
##      a ~~
##      x1          0.844    0.055   15.355    0.000
##      x2          0.867    0.056   15.441    0.000
##      x3          0.845    0.055   15.400    0.000
##      x4          0.822    0.053   15.455    0.000
##      x5          0.820    0.053   15.572    0.000
##      x1 ~~
##      x2          0.908    0.070   12.900    0.000
##      x3          0.935    0.069   13.466    0.000
##      x4          0.921    0.067   13.661    0.000
##      x5          0.914    0.067   13.716    0.000
##      x2 ~~
##      x3          0.889    0.070   12.675    0.000
##      x4          0.922    0.069   13.423    0.000
##      x5          0.889    0.068   13.165    0.000
##      x3 ~~
##      x4          0.865    0.067   12.976    0.000
##      x5          0.901    0.066   13.554    0.000
##      x4 ~~
##      x5          0.850    0.064   13.285    0.000
...
```

Finally, the variance estimates are listed. Here, we see that in order to mimic the behaviour of a traditional FE model, the error variances over time were specified to be equal using the equality constraint (`e`). Notice the `.` beside `y1`, `y2`, etc.: this indicates that the listed variance refers to an endogenous variable, and that it is thus an error

⁸This is at least the current behaviour of both the `cfa` and `sem` wrappers. In fact, both wrappers seem to be identical in terms of the default settings, see Rosseel et al. (2020).

variance. In this case, these refer to the variances of ε_t . After that, the variances of the exogenous variables, both observed and unobserved are listed.

```
...
## Variances:
##
```

		Estimate	Std.Err	z-value	P(> z)
##	.y1	(e) 1.022	0.023	44.721	0.000
##	.y2	(e) 1.022	0.023	44.721	0.000
##	.y3	(e) 1.022	0.023	44.721	0.000
##	.y4	(e) 1.022	0.023	44.721	0.000
##	.y5	(e) 1.022	0.023	44.721	0.000
##	x1	1.986	0.089	22.361	0.000
##	x2	2.079	0.093	22.361	0.000
##	x3	1.987	0.089	22.361	0.000
##	x4	1.860	0.083	22.361	0.000
##	x5	1.814	0.081	22.361	0.000
##	a	0.799	0.052	15.310	0.000

6. Conclusion

Fixed effects regression in SEM has been outlined in well-known articles by (Allison 2011; Bollen and Brand 2010; Teachman et al. 2001). This article provides a focused look at the implementation of the basic model using the `lavaan` package in R. The online supplementary materials further discuss common extensions and some tools for evaluating and loosening model assumptions.

The benefits of FE-SEM as opposed to traditional OLS-based FE-models are largely the same ones that apply to the SEM framework in general: for one, SEM allows for a great deal of flexibility. For example, it is easy to loosen model constraints as necessary. Measurement error in both the dependent and independent variables can be dealt with using latent variables to achieve unbiased and more efficient results. Researchers interested in time-invariant predictors can integrate them into a hybrid FE/RE model with ease. Further extensions, like measurement invariance testing (van de Schoot, Lugtig, and Hox 2012; Millsap 2011; Steenkamp and Baumgartner 1998) as well as lagged dependent variables (Bollen and Brand 2010; Allison, Williams, and Moral-Benito 2017) for example, can also be implemented in a straightforward fashion.

The most basic FE-SEM is furthermore the basis for a variety of currently popular extended models, such as Latent Curve Models in general (Curran and Bollen

2001; Bollen and Curran 2004), as well as special implementations like the Dynamic Panel Model (Allison, Williams, and Moral-Benito 2017), the Random-Intercept Cross-Lagged Panel Model (Hamaker, Kuiper, and Grasman 2015) and the Latent Curve Model with Structured Residuals (Curran et al. 2014). For this reason it is all the more important for researchers to have a good grasp on the method of applying panel regression in SEM, and understanding the intuition of controlling for time-invariant confounders. This article is meant to serve as a consolidated resource for researchers looking for concrete advice on specifying FE and more general panel models in SEM.

References

- Allison, Paul. 2011. *Fixed Effects Regression Models*. Thousand Oaks: Sage Publications.
- Allison, Paul, Richard Williams, and Enrique Moral-Benito. 2017. “Maximum Likelihood for Cross-lagged Panel Models with Fixed Effects.” *Socius* 3: 1–17.
- Bates, Douglas, Martin Mächler, Ben Bolker, and Steve Walker. 2015. “Fitting Linear Mixed-Effects Models Using lme4.” *Journal of Statistical Software* 67 (1): 1–48.
- Bates, Douglas M. 2010. *lme4: Mixed-effects modeling with R*. Springer.
- Bell, Andrew, Malcolm Fairbrother, and Kelvyn Jones. 2018. “Fixed and random effects models: Making an informed choice.” *Quality and Quantity* 53: 1051–1074.
- Bollen, Kenneth. 1989. *Structural Equations with Latent Variables*. New York, Chichester: Wiley. ISBN 0-471-01171-1.
- Bollen, Kenneth, and Jennie Brand. 2010. “A General Panel Model with Random and Fixed Effects: A Structural Equations Approach.” *Social Forces* 89(1): 1–34.
- Bollen, Kenneth, and Patrick Curran. 2004. “Autoregressive Latent Trajectory (ALT) Models: A Synthesis of Two Traditions.” *Sociological Methods and Research* 32(3): 336–383.
- Brüderl, Josef, and Volker Ludwig. 2015. “Fixed-effects panel regression.” In *The Sage Handbook of Regression Analysis and Causal Inference*, edited by Henning Best and Christof Wolf, Chap. 15, 327–357. London, Thousand Oaks: Sage Publications.
- Chamberlain, Gary. 1980. “Analysis of Covariance with Qualitative Data.” *Review of Economic Studies* 47(1): 225–238.
- Croissant, Yves, and Giovanni Millo. 2008. “Panel Data Econometrics in R: The plm Package.” *Journal of Statistical Software* 27 (2): 1–43.
- Curran, Patrick, and Daniel Bauer. 2011. “The Disaggregation of Within-Person and Between-

- Person Effects in Longitudinal Models of Change.” *Annual Review of Psychology* 62: 583–619.
- Curran, Patrick, and Kenneth Bollen. 2001. “The best of both worlds: Combining autoregressive and latent curve models.” In *New methods for the analysis of change*, edited by L. Collins and A. Sayer, 107–135. Washington, DC: American Psychological Press.
- Curran, Patrick, Andrea Howard, Sierra Bainter, Stephanie Lane, and James McGinley. 2014. “The Separation of Between-Person and Within-Person Components of Individual Change Over Time: A Latent Growth Curve Model With Structured Residuals.” *Journal of Consulting and Clinical Psychology* 82(5): 879–894.
- Enders, Craig, and Davood Tofghi. 2007. “Centering predictor variables in cross-sectional multilevel models: A new look at an old issue.” *Psychological Methods* 12(2): 121–138.
- Graff, J. 1979. “Verallgemeinertes LISREL-Modell.” Unpublished Manuscript.
- Hamaker, Ellen, Rebecca Kuiper, and Raoul Grasman. 2015. “A Critique of the Cross-Lagged Panel Model.” *Psychological Methods* 20(1): 102–116.
- Kline, Rex. 2016. *Principles and Practice of Structural Equation Modeling. Fourth Edition*. New York: The Guilford Press.
- Millsap, Roger. 2011. *Statistical Approaches to Measurement Invariance*. New York, London: Routledge.
- Mundlak, Yair. 1978. “On the Pooling of Time Series and Cross Section Data.” *Econometrica* 46(1): 69–85.
- Muthén, Linda, and Bengt Muthén. 1998–2017. *Mplus User’s Guide. Eighth Edition*. Los Angeles, California.
- R Core Team. 2017. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Rosseel, Yves. 2012. “lavaan: An R Package for Structural Equation Modeling.” *Journal of Statistical Software* 48 (2): 1–36. <http://www.jstatsoft.org/v48/i02/>.
- Rosseel, Yves, Terrence D. Jorgensen, Daniel Oberski, Jarrett Byrnes, Leonard Vanbrabant, Victoria Savalei, Ed Merkle, et al. 2020. *lavaan: Latent Variable Analysis*. <https://CRAN.R-project.org/package=lavaan>.
- Schmidheiny, Kurt. 2019. “Panel Data: Fixed and Random Effects.” Lecture notes.
- Schunck, Reinhard. 2013. “Within and between estimates in random-effects models: Advantages and drawbacks of random effects and hybrid models.” *The Stata Journal* 13(1): 65–76.
- Skrondal, Anders, and Sophia Rabe-Hesketh. 2004. *Generalized Latent Variable Modeling*. Boca

- Raton, London, New York, Washington, D.C.: Chapman & Hall/CRC.
- Steenkamp, Jan-Benedict, and Hans Baumgartner. 1998. "Assessing Measurement Invariance in Cross-National Consumer Research." *Journal of Consumer Research* 25(1): 78–90.
- Teachman, Jay, Greg Duncan, Jean Yeung, and Dan Levy. 2001. "Covariance Structure Models for Fixed and Random Effects." *Sociological Methods and Research* 30(2): 271–288.
- van de Schoot, Rens, Peter Lugtig, and Joop Hox. 2012. "A checklist for testing measurement invariance." *European Journal of Developmental Psychology* 9(4): 486–492.
- Wooldridge, Jeffery. 2002. *Econometric analysis of cross sectional and panel data*. Cambridge, Massachusetts: The MIT Press. ISBN 0-262-23219-7.
- Wooldridge, Jeffery. 2012. *Introductory Econometrics: A Modern Approach, 5th Edition*. Mason, Ohio: Thomson South-Western. ISBN 1-111-53104-8.

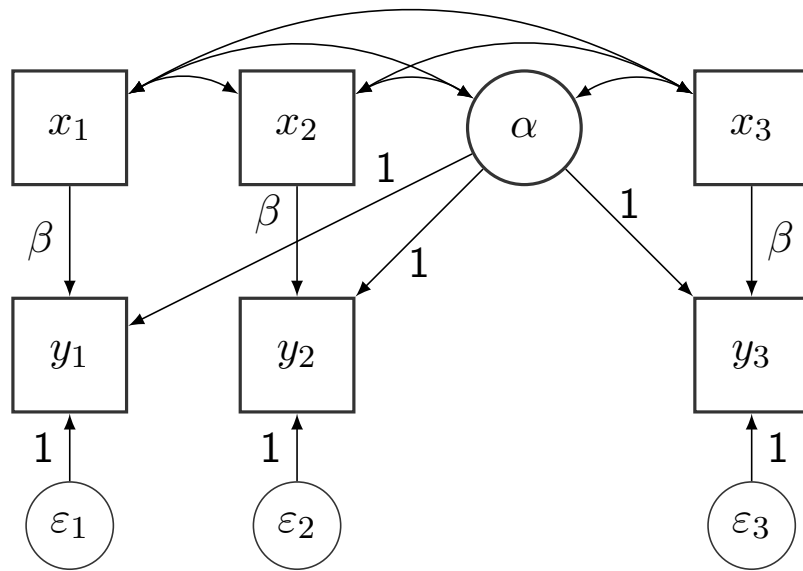


Figure 1. Typical three-wave FE-SEM model with contemporaneous effects