

A Closer Look at Fixed-Effects Regression in Structural Equation Modeling Using lavaan

Henrik Kenneth Andersen

2020-06-15

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. It further discusses 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 fixed-/random-effects hybrid 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, the perhaps greatest strength of panel data is that they allow for a more rigorous testing of substantive theories as a tool for *causal analysis* (Greene 2012). 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. This is the essence of fixed-effects (FE) regression. It drastically reduces the number of potential confounders of the relationship between variables of interest by allowing one to control for all time-invariant influences. Thus, by using each unit at an earlier point in time as its own control, and by examining solely *within-unit changes*, we come closer to the *counterfactual* or *potential outcomes* ideal: what would the outcome have been, had the observational unit not been exposed to the treatment (Angrist and Pischke 2009)?

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.

This article is intended as a 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, n.d.) is unarguably 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, it and R are 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 fitted `lavaan` object, 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 FE regression, the particularities of FE in SEM, and shows its implementation in `lavaan`. Using simulated data (found in the Appendix), it demonstrates and annotates the code for the most

basic FE model and provides an overview of the summary output. After that, a number of potential extensions are discussed and demonstrated, from loosening assumptions, to dealing with measurement error, to specifying hybrid fixed- and random effects models in order to include time-invariant predictors.

2 Fixed-effects

The most basic fixed-effects model can be expressed as

$$y_{it} = \mathbf{x}_{it}\boldsymbol{\beta} + \alpha_i + \varepsilon_{it}, \quad (1)$$

where y_{it} is the observed outcome variable for unit i at time t , \mathbf{x}_{it} is a $1 \times k$ vector of covariates linked to the outcome by the $k \times 1$ vector of coefficients $\boldsymbol{\beta}$. All stable, unit-specific characteristics such as date of birth, country of origin, sex, etc. are captured by the time-constant part of the error α_i , whereas ε_{it} is the idiosyncratic error term that varies between units and over time (Brüderl and Ludwig 2015).

The point of FE regression is to be able to drop assumptions regarding the relatedness of the model covariates and the unobserved time-constant effects, i.e., $\mathbb{E}[\mathbf{x}_{it}^\top \alpha_i]$. We can do so by manipulating Equation (2) so that the individual effect, α_i is eliminated. In doing so, we no longer need to worry about potential bias due to unit-specific characteristics. This can be achieved in a number of ways (e.g., differencing, least squares dummy variable regression), but the most common method involves subtracting the over-time unit mean from each of the terms:

$$\begin{aligned} (y_{it} - \bar{y}_i) &= (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + (\alpha_i - \bar{\alpha}_i) + (\varepsilon_{it} - \bar{\varepsilon}_i) \\ \check{y}_{it} &= \check{\mathbf{x}}_{it}\boldsymbol{\beta} + \check{\varepsilon}_{it}. \end{aligned}$$

Because the average of something that does not change is that thing itself, the individual effects 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. The unbiasedness of the estimate is related solely to the strict exogeneity assumption imposed on the indiosyncratic errors, i.e., $\mathbb{E}[\check{\varepsilon}_{it} | \check{\mathbf{x}}_{it}] = 0$ which also implies $\mathbb{E}[\check{\mathbf{x}}_{is}^\top \check{\varepsilon}_{it}] = \mathbf{0}$, $\forall s, t = 1, \dots, T$ (Brüderl and Ludwig 2015; Wooldridge 2002). Strict exogeneity not only rules out contemporaneous correlations between the covariates and the error, but it also implies that for the estimates $\boldsymbol{\beta}$ to be unbiased, the covariates at all timepoints must not be correlated with the error at all timepoints.

2.1 Fixed-effects in structural equation modeling

Analytically, Equation (2) applies as well to FE models in SEM as we will see shortly. However, there are a number of differences we must address. First, SEM is a covariance-based approach. This means that the vectors of individual observations are only of interest insofar as they can be converted into observed covariances. Second, the way in which we use a latent variable to represent the individual effects forces us to convert the data from stacked, long-format vectors of length NT to T wide-format vectors of length N . Lastly, the notation in

SEM differs from traditional least squares-based models. 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):

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

where $\boldsymbol{\eta}^+ = (\mathbf{y}, \mathbf{x}, \boldsymbol{\eta}, \boldsymbol{\xi})^\top$, $\boldsymbol{\zeta}^+ = (\boldsymbol{\varepsilon}, \boldsymbol{\delta}, \boldsymbol{\zeta}, \boldsymbol{\xi})^\top$, $\mathbf{y}^+ = (\mathbf{y}, \mathbf{x})^\top$. Furthermore, \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 (i.e., ‘endogenous’) and independent (i.e., ‘exogenous’), $\boldsymbol{\zeta}^+$ 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 $\boldsymbol{\zeta}^+$ 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).

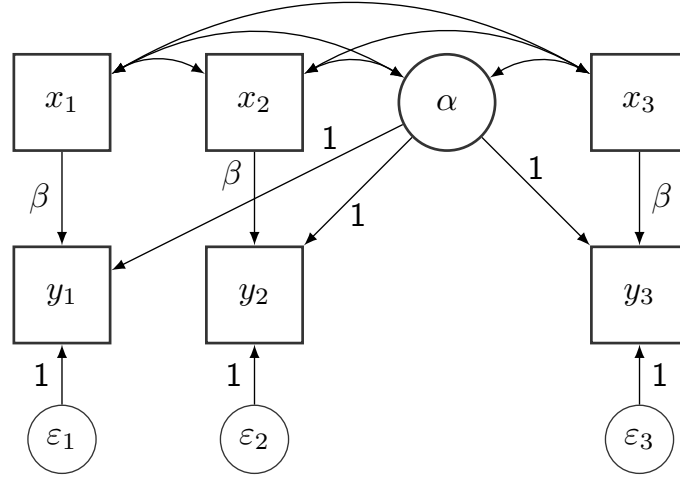


Figure 1: Typical three-wave FE-SEM model with contemporary effects

This notation may be confusing at first, but it has advantages. First, it allows us to be more flexible with our models. For example, it allows observed \mathbf{x} to influence observed \mathbf{y} directly (more common notation assumes that substantive effects occur only between latent variables, observed ones are only used as indicators). 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, which will be discussed later in this article. It is also consistent with the notation used in popular SEM software like `lavaan` and `Mplus`. 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. Let us, however, make things more concrete and take a look at a simple, three-wave version of the typical FE model in SEM using this notation (see Figure 1). For that, we have:

$$\mathbf{y}^+ = \mathbf{\Lambda}_y^+ \boldsymbol{\eta}^+$$

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

$$\boldsymbol{\eta}^+ = \boldsymbol{\eta}^+ \mathbf{B} + \boldsymbol{\zeta}^+$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{pmatrix} \begin{pmatrix} y_1 & y_2 & y_3 & x_1 & x_2 & x_3 & \alpha \\ 0 & 0 & 0 & \beta & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & \beta & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & \beta & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ x_1 = \delta_1 \\ x_2 = \delta_2 \\ x_3 = \delta_3 \\ \alpha = \xi \end{pmatrix}. \quad (2)$$

Notice in $\boldsymbol{\zeta}^+$, for the independent variables it does not matter whether we write, for example x_1 or δ_1 . By treating these variables as exogenous, we can say their entire variance is made up of ‘error’, i.e., it is the combined effect of all unobserved influences.

Admittedly, Equation (2.1) 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.:

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

where \mathbf{I} is the identity matrix. By substituting this back into the equation for

the observed variables we get $\mathbf{y}^+ = \mathbf{\Lambda}_y^+[(\mathbf{I} - \mathbf{B})^{-1}\boldsymbol{\zeta}^+]$, 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 (2).

The traditional FE model in Equation (2) works by eliminating the individual effects by demeaning, i.e., subtracting the over-time unit mean from each of the model variables. In SEM, we specify a latent variable, here we call it α , to represent the combined effect of all time-invariant influences. Colloquially, we can say that this variable represents what the dependent variable has in common over time. For the sake of simplicity, let us ignore the independent variable for a moment, and rewrite Equation (2) as

$$y_{it} = \alpha_i + \varepsilon_{it}.$$

If α and ε are independent, then the variance of y is $\text{Var}(\alpha) + \text{Var}(\varepsilon)$. Once we transform the long-format data into wide-format, then we can see obviously that the covariance between the T columns of y_t is just $\text{Var}(\alpha)$

$$\begin{aligned} y_t &= \alpha + \varepsilon_t, \\ \text{Cov}(y_t, y_s) &= \mathbb{E}[(\alpha + \varepsilon_t)(\alpha + \varepsilon_s)] \\ &= \alpha^2 \\ &= \text{Var}(\alpha), \end{aligned}$$

assuming $\text{Cov}(\alpha, \varepsilon_t) = 0, \forall t$ and $\text{Cov}(\varepsilon_t, \varepsilon_s) = 0, \forall s, t$. And in fact this is exactly how it works. For example, if we had three waves of observations, we could write the variances and covariances of the observed variables, in this case $\mathbf{y} = (y_1, y_2, y_3)^\top$ as a system of six linear equations in the form of $\mathbf{Ax} = \mathbf{b}$

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

$$\begin{matrix} & \alpha^2 & \varepsilon_1^2 & \varepsilon_2^2 & \varepsilon_3^2 \\ \begin{matrix} y_1^2 \\ y_2 y_1 \\ y_3 y_1 \\ y_2^2 \\ y_3 y_2 \\ y_3^2 \end{matrix} & \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix} & \begin{pmatrix} \alpha^2 \\ \varepsilon_1^2 \\ \varepsilon_2^2 \\ \varepsilon_3^2 \end{pmatrix} & = & \begin{pmatrix} y_1^2 \\ y_2 y_1 \\ y_3 y_1 \\ y_2^2 \\ y_3 y_2 \\ y_3^2 \end{pmatrix} \end{matrix}$$

solve this equation: $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \mathbf{x}$ which works out to

$$\mathbf{x} = \begin{pmatrix} \alpha^2 \\ \varepsilon_1^2 \\ \varepsilon_2^2 \\ \varepsilon_3^2 \end{pmatrix} = \begin{pmatrix} .33(y_2 y_1) + .33(y_3 y_1) + .33(y_3 y_2) \\ y_1^2 - (.33(y_2 y_1) + .33(y_3 y_1) + .33(y_3 y_2)) \\ y_2^2 - (.33(y_2 y_1) + .33(y_3 y_1) + .33(y_3 y_2)) \\ y_3^2 - (.33(y_2 y_1) + .33(y_3 y_1) + .33(y_3 y_2)) \end{pmatrix}$$

which means $\text{Var}(\alpha) = \frac{\text{Cov}(y_2, y_1) + \text{Cov}(y_3, y_1) + \text{Cov}(y_3, y_2)}{3}$ the individual effects is the average covariance between y at the different timepoints. It is assumed, based on the long-format equation $y_{it} = \alpha_i + \varepsilon_{it}$, that $\text{Cov}(y_2, y_1) = \text{Cov}(y_3, y_1) = \text{Cov}(y_3, y_2) = \text{Var}(\alpha)$ which is to say $3 \frac{\text{Var}(\alpha)}{3} = \text{Var}(\alpha)$.

Remember, in Equation (2), y_{it} is a stacked $NT \times 1$ vector which we converted into T individual $N \times 1$ vectors for the SEM. For the long-format data

This explanation may not be satisfying to some, so we can get even more specific and show the intuition behind the use of latent variables to decompose the dependent variable into its between- and within-variance components. For the sake of simplicity, let us ignore the independent variable for a moment, and concentrate just on the dependent variable.

What essentially differentiates an FE from a random effects (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 $\mathbb{E}[\alpha x_t] \neq 0$. 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. The RE model assumes $\mathbb{E}[\alpha x_t] = 0$. The benefit of an RE model over a simple pooled ordinary least squares model (POLS) is that because it accounts for serial correlation in the composite error term, i.e., $\nu_{it} = \alpha_i + \varepsilon_{it}$, we achieve accurate standard errors. The coefficient of interest, however, is unaffected and will be unbiased.

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

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

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 model:

$$\Psi = \begin{matrix} & \varepsilon_1 & \varepsilon_2 & \varepsilon_3 & x_1 & x_2 & x_3 & \alpha \\ \begin{matrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ x_1 \\ x_2 \\ x_3 \\ \alpha \end{matrix} & \begin{pmatrix} \varepsilon_1^2 & & & & & & \\ 0 & \varepsilon_2^2 & & & & & \\ 0 & 0 & \varepsilon_3^2 & & & & \\ 0 & 0 & 0 & x_1^2 & & & \\ 0 & 0 & 0 & x_2 x_1 & x_2^2 & & \\ 0 & 0 & 0 & x_3 x_1 & x_3 x_2 & x_3^2 & \\ 0 & 0 & 0 & \alpha x_1 & \alpha x_2 & \alpha x_3 & \alpha^2 \end{pmatrix} \end{matrix}.$$

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:

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

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 $\hat{\beta} = \text{Cov}(y_t, x_t) / \text{Var}(x_t)$. The rest of the model-implied covariance matrix results from $\mathbf{y}^+ \mathbf{y}^{+\top}$.

3 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.¹

Specifying the most basic fixed-effects model, like the one shown in Bollen and Brand (2010) involves four components. First, we define the latent individual effects variable using the `=~` ‘measured by’ 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

¹There are many online tutorials for importing data in various formats, see, for example some from datacamp or Quick-R, or any of the literally thousands of posts on stackoverflow.

could have used any letter) 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 a 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 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
```

4 An example

To demonstrate the application of FE models in SEM, a dataset can be simulated that embodies the FE assumptions. The code for data simulation can be found in Appendix A.

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 effects. This means that approaches that fail to account for this confounding influence, such as pooled ordinary least squares (POLS) or random effects (RE), will be biased.

The wide-format equations for the DGP 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 both $\boldsymbol{\alpha}_1$ and $\boldsymbol{\alpha}_2$ are $\sim N(\mu_{\alpha_j}, \sigma_{\alpha_j})$, $j = 1, 2$. $\boldsymbol{\delta}_t$ is $\sim N(\mu_{\delta}, \sigma_{\delta})$ and $\boldsymbol{\varepsilon}_t$ is $\sim N(0, 1)$.

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 $\sim N(0, 1)$. 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 $\beta_{x, \alpha_2} = 0.50$. 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")
saveRDS(fe_sem, "fe_sem.Rda")
saveRDS(fe_sem.fit, "fe_sem.fit.Rda")
```

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 `a` for 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.²

```
...
## 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
```

²This 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).

```
##      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 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

```

5 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, as well as common extensions using the `lavaan` package in R.

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 (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 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 consolidated guide to a variety of common applications of fixed-effects 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. <https://doi.org/10.1177/2378023117710578>.
- Angrist, Joshua, and Jörn-Steffen Pischke. 2009. *Mostly Harmless Econometrics: An Empiricist’s Companion*. New Jersey: Princeton University Press.
- Bollen, Kenneth. 1989. *Structural Equations with Latent Variables*. New York, Chichester: Wiley.
- 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): 336–83. <https://doi.org/10.1177/0049124103260222>.
- 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, 327–57. London, Thousand Oaks: Sage Publications.
- 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, 107–35. Washington, DC: American Psychological Press. <https://doi.org/10.1037/10409-004>.
- 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): 879–94. <https://doi.org/10.1037/a0035297>.
- Graff, J. 1979. “Verallgemeinertes LISREL-Modell.” Mannheim, Germany.
- Greene, William. 2012. *Econometric Analysis, 7th Edition, International Edition*. Edinburgh Gate, Essex: Pearson Education Limited.
- Hamaker, Ellen, Rebecca Kuiper, and Raoul Grasman. 2015. “A Critique of the Cross-Lagged Panel Model.” *Psychological Methods* 20(1): 102–16. <https://doi.org/10.1037/a0038889>.
- Millsap, Roger. 2011. *Statistical Approaches to Measurement Invariance*. New York, London: Routledge.
- Muthén, Linda, and Bengt Muthén. n.d. *Mplus User’s Guide. Eighth Edition*. Los Angeles, California: Muthén & Muthén.

- 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>.
- Schoot, Rens van de, Peter Lugtig, and Joop Hox. 2012. “A Checklist for Testing Measurement Invariance.” *European Journal of Developmental Psychology* 9(4): 486–92.
- 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): 271–88.
- Wooldridge, Jeffery. 2002. *Econometric Analysis of Cross Sectional and Panel Data*. Cambridge, Massachusetts: The MIT Press.