

Skewness and staging: Does the floor effect induce bias in multilevel AR(1) models?

Reproducible codes

MH Manuel Haqiqatkhah

Table of contents

1	Introduction	2
2	Core components	3
2.1	The <i>Simulation</i> component	3
2.1.1	Data-generating models specifications	3
2.1.1.1	The AR(1) model	4
2.1.1.1.1	Construction	4
2.1.1.1.2	Code	5
2.1.1.2	The χ^2 AR(1) model	5
2.1.1.2.1	Construction	5
2.1.1.2.2	Code	5
2.1.1.3	The BinAR(1) model	5
2.1.1.3.1	Construction	6
2.1.1.3.2	Code	6
2.1.1.4	The PoDAR(1) model	6
2.1.1.4.1	Construction	6
2.1.1.4.2	Code	7
2.1.2	General DGM wrappers	7
2.1.2.1	Parameter conversions	7
2.1.2.2	Wrapper around <code>dgm_*</code> functions	7
2.1.3	Dataset generation	8
2.1.3.1	Making individual datasets	8
2.1.3.2	Determining level-2 distribution parameters	8
2.1.3.3	Automate dataset generation	11
2.2	The <i>Analysis</i> component	11
2.2.0.1	Fixed residual variance	12
2.2.0.2	Random residual variance	12

2.2.0.3	Comparing different ANALYSIS parameters	13
2.3	The <i>Harvesting</i> component	14
2.4	The <i>Reporting</i> component	14
3	Pipeline	14
3.1	Book-keeping	14
3.1.0.1	Simulation file references	14
3.1.0.2	Analysis outputs file references	14
3.2	Wrapper functions	15
3.2.0.1	Simulation in parallel	15
3.2.0.2	Analysis and save in parallel	15
3.2.0.3	Harvesting in parallel	15
3.3	Pipeline code	15
3.4	Simulation results	15
4	Additional analyses and figures	15
4.1	Profiles of simulated datasets	15
	COGITO data analysis	15

1 Introduction

This document contains the reproducible code for the manuscript Skewness and staging: Does the floor effect induce bias in multilevel AR(1) models? by M. M. Haqiqatkhah, O. Ryan, and E. L. Hamaker. please cite as:

...

In this study, we simulated multilevel data from three data generating mechanisms (DGMs), namely, the AR(1), χ^2 AR(1), BinAR(1), and PoDAR(1) models with different parameter sets. For details, see the paper.

The simulation was conducted using the following modular pipeline design, inspired by Bien's R package **simulator** (2016), consisting of the following **components**:

- A. *Simulation*: generating the datasets
- B. *Analysis*: modeling the data
- C. *Harvesting*: collecting the relevant parameter estimates
- D. *Reporting*: making tables and plots

And the components were placed in a **pipeline**, that managed:

- 1. Making the simulation design matrix that include all relevant conditions
- 2. Book-keeping data files belonging to each replication of each condition
- 3. Performing simulations in batch

4. Performing Analyses in batch
5. Collecting the data in batch

This document is structured as follows. In Section 2, we explain the four components and the functions used therein. Then, in Section 3 we explain the wrapper functions used in the pipeline, and show how the pipeline was—and can be—executed. Then, in Section 3.4, we discuss how the harvested data was used to make the figures used in the paper (and others that were not included). Finally, in Section 4 we discuss how supplementary plots were made and how the empirical data analysis (on COGITO data) was done.

Before we begin, we need to read the scripts to include them in this document. By default, none of the scripts run here (as they are time consuming). To run the scripts of each component, you can change the following variables to `TRUE` and re-render the document:

```
runComponent_Simulation <- FALSE
runComponent_Analysis <- FALSE
runComponent_Harvesting <- FALSE
runComponent_Reporting <- FALSE
```

In case you want to change the scripts (e.g., to run a smaller portion of the simulation, or try other parameters, etc.), you should look up the `knitr` parameters called in each chunk (with `<<some_param>>`) and find the corresponding code (under `## @knitr some_param`) in the `scripts` folder.

2 Core components

2.1 The *Simulation* component

The Simulation component consists of three sets of functions:

- i. Functions that implement the DGMs and generate univariate ($N = 1$) time series of length T from the parameters given to them;
- ii. Wrappers that interface the DGM functions;
- iii. A wrapper to generate datasets (consisting on N time series of length T) with a given DGM

2.1.1 Data-generating models specifications

First we define functions for each data-generating models (DGMs) that can produce univariate, single-subject ($N = 1$) time series of desired length T (default: `T = 100`) with the two canonical parameters and a given random seed (default: `seed = 0`). All model(-implied) parameters are saved in a list (called `pa`).

For each model, the first observation (X_1) is randomly drawn from the model-implied marginal distribution, to eliminate the need for removing the burn-in window in the beginning of the data. After the data is generated, in case the argument `only.ts` is set to be `TRUE`, the raw data (as a vector of length `T`) is returned. Otherwise, the function calculates empirical dynamic (ϕ) and marginal (μ , σ^2 , and γ) parameters based on the simulated data, and save it in a list (`Empirical.Parameters`). Furthermore, two L^AT_EX-ready strings (`Model.Description` and `Model.Description.Short`) are made which include a summary of the model parameters (that can be used, e.g., in plots). Finally, in case `only.ts != TRUE`, the function returns a list consisting of the time series (stored in `x`), verbal description of the dataset (`Model.Description` and `Model.Description.Short`), theoretical (i.e., model-implied) parameters (`Model.Parameters`), and empirical (i.e., sample) estimated parameters (`Empirical.Parameters`).

2.1.1.1 The AR(1) model

The canonical parameters of the AR(1) model with normally distributed residuals (which we referred to as `NAR(1)` in the simulation) are the autoregressive parameter ϕ (default: `phi = 0.4`), mean μ (default: `Mean = 50`), and the marginal variance σ^2 (default: `var.marginal = 4`). Based on the marginal variance, the residual variance (`var.resid`) is calculated via $\sigma_\epsilon^2 = \sigma^2(1 - \phi^2)$.

2.1.1.1.1 Construction

The time series is constructed by first generating a zero-centered time series \tilde{X}_t (`x_cent`). To do so, first the initial observation in the time series (`x_cent[1]`) is sampled from normal distribution with mean zero and a variance equal to the marginal variance of the model:

$$\tilde{X}_1 \sim \mathcal{N}(0, \sigma^2)$$

Then, the remainder of the time series is generated using the definition of the AR(1) model (not that the here the residual variance is used in the normal distribution):

$$\begin{aligned}\tilde{X}_t &= \phi \tilde{X}_{t-1} + \epsilon_t \\ \epsilon_t &\sim \mathcal{N}(0, \sigma_\epsilon^2)\end{aligned}$$

Finally, the mean is added to the centered zero-centered time series to reach the final time series with mean μ :

$$X_t = \tilde{X}_t + \mu$$

2.1.1.1.2 Code

2.1.1.2 The $\chi^2\text{AR}(1)$ model

The canonical parameters of the $\chi^2\text{AR}(1)$ model (which we referred to as **ChiAR(1)** in the simulation) are the autoregressive parameter ϕ (default: `phi = 0.4`), and degrees of freedom ν (default: `nu = 3`). We set the intercept to zero (`c = 0`).¹

2.1.1.2.1 Construction

Similar to the $\text{AR}(1)$ model, we need to sample the first observation of the $\chi^2\text{AR}(1)$ model from its marginal distribution. However, since this model does not have a closed-form marginal distribution, as an approximation, we instead sample $\mathbf{x}[1]$ from a χ^2 distribution with ν degrees of freedom:

$$X_1 \sim \chi^2(\nu)$$

Then, we generate the remainder of the time series using the definition of the $\chi^2\text{AR}(1)$ model:

$$\begin{aligned} X_t &= c + \phi X_{t-1} + a_t \\ a_t &\sim \chi^2(\nu). \end{aligned}$$

2.1.1.2.2 Code

2.1.1.3 The **BinAR(1)** model

The canonical parameters of the **BinAR(1)** model (which we referred to as **BinAR(1)** in the simulation) are the survival probability α (default: `alpha = 0.5`) and the revival probability β (default: `beta = 0.4`). By default, the maximum value on scale k was set to `k = 10`.

¹The $\chi^2\text{AR}(1)$, in a more general form, can have an intercept ($X_t = c + \phi X_{t-1} + a_t$, $a_t \sim \chi^2(\nu)$). Since the intercept was set to zero in the simulation study, we discussed a zero-intercept version of this model ($c = 0$) in the paper. See the Supplemental Materials for more details.

2.1.1.3.1 Construction

We first calculate the θ parameter, which characterizes the marginal distribution of the BinAR(1) model:

$$\theta = \frac{k\beta}{1 - (\alpha - \beta)}$$

Then we draw X_1 (`x[1]`) from the marginal distribution of the model:

$$X_1 \sim \text{Binom}(k, \theta)$$

The rest of time series is generated sequentially, for each time point t , by drawing values for the number of survived (`S_t[t]`) and revived (`R_t[t]`) elements of the BinAR(1) model based on the previous observations (X_{t-1}), and then adding them:

$$\begin{aligned} S_t &\sim \text{Binom}(X_{t-1}, \alpha) \\ R_t &\sim \text{Binom}(k - X_{t-1}, \beta) \\ X_t &= S_t + R_t \end{aligned}$$

2.1.1.3.2 Code

2.1.1.4 The PoDAR(1) model

The canonical parameters of the PoDAR(1) model (which we referred to as PoDAR(1) in the simulation) are the persistence probability τ (default: `tau = 0.7`) and the average rate λ (default: `lambda = 0.5`).

2.1.1.4.1 Construction

To generate the time series, we first draw the first observation X_1 (`x[1]`) from a Poisson distribution with rate λ :

$$X_1 \sim \text{Poisson}(\lambda)$$

And generate the rest of the time series by first drawing Z_t from a Poisson distribution with rate λ and P_t from a binomial distribution with size probability of success τ (that is equivalent to a Bernoulli distribution with probability τ). Then, we calculate X_t based on the previous observation (`x[t-1]`) and values of Z_t (`Z_t[t]`) and P_t (`P_t[t]`), using the definition of the PoDAR(1) model:

$$\begin{aligned}
Z_t &\sim \text{Poisson}(\lambda) \\
P_t &\sim \text{Binom}(1, \tau) \\
X_t &= P_t X_{t-1} + (1 - P_t) Z_t
\end{aligned}$$

2.1.1.4.2 Code

2.1.2 General DGM wrappers

Given that, in each model, two canonical parameters characterize the dynamic and marginal features of the generated time series, and given that we have analytic formulas that link the canonical parameters to the model-implied ϕ , μ , σ^2 , and γ , we use a function (`dgm_parameterizer`) to calculate canonical parameters from two given parameters, and make a complete list of parameters (called `pa`). This list also includes non-parameter variables, importantly, the time series length T (saved in `pa$T`) and the random seed used in the `dgm_*` functions (saved in `pa$seed`). A wrapper function (`dgm_generator`) is used as an interface to all `dgm_*` functions, which first makes sure the given parameters are sufficient for data generation, makes a complete parameter list `pa` with the help of `dgm_parameterizer`, and passes `pa` to the respective DGM generating function.

2.1.2.1 Parameter conversions

The function `dgm_parameterizer` calculates canonical/model-implied parameters of a given DGM (specified using the `Model` argument) based on the parameters given to it as arguments, and saves them in a list of parameters (`pa`), which is returned by the function. The function makes sure that the set of parameters provided are sufficient to characterize the dynamic parameter of the model (i.e., the autoregression ϕ) and at least one of the marginal parameters (importantly, the mean μ) but giving default values to some parameters.

2.1.2.2 Wrapper around `dgm_*` functions

The function `dgm_generator` gets a set of parameters (either as separate arguments, or a list of parameters, like the one returned by `dgm_parameterizer`), saves them in a list called `pa`. It checks whether ϕ is included in the list (if not, sets the default value `pa$phi = 0.2`), and checks if at least one other parameter (which, together with ϕ , is required to characterize the marginal properties of the DGMs) is calculated for it (if not, it sets the default value `pa$Mean = 5` for μ). Furthermore, if the DGM name, time series length, and the random seed are not provided, it gives them default values (respectively: `Model = "ChiAR(1)"`, `T = 100`, and `seed = 0`) and adds them to `pa`.

Then, it passes the `pa` list to `dgm_parameterizer` to do the necessary conversions to complete the list of canonical and model-implied parameters. Finally, given the model name, it checks if non-canonical parameters k and c are set (otherwise assigns appropriate defaults to them), and passes the complete parameter list to the respective DGM function.

2.1.3 Dataset generation

The machinery described above can be used to generate individual ($N = 1$) time series. However, for the simulation study, we need datasets comprising of multiple ($N = 25, 50, 100$) individuals. As we discussed in the paper, in our study, all individuals in a dataset of a DGM share the same autoregressive parameter ($\phi_i = 0.4$) and the individual differences are only in the individual means ($\mu = [\mu_1, \mu_2, \dots, \mu_N]$). Thus, we write a function (`dgm_make.sample`) that can generate, for each DGM, a dataset of N individuals based on an N -dimensional vector of individual means, all with the same ϕ_i . We then need to find the appropriate parameters for the level-2 distributions (Gaussian and χ^2 distributions) for each DGM, such that we get a considerable proportion of individuals with considerably high skewness while respecting the lower and upper bounds of values supported by each model. Finally, with a wrapper function (`make_datasets`), we facilitate making dataset by automatically generating the means vector suitable for each DGM.

2.1.3.1 Making individual datasets

The function `dgm_make.sample` generates a dataset of time series of length `T` with the autoregressive parameter `phi` from a desired DGM (determined by the `Model` argument) given a vector of means (passed as the argument `Means`). The length of `Means` determine the number of individuals in the dataset (`N <- length(Means)`). If `Means` is not provided, a randomly generated vector of $N = 100$ is used as default. Since each individual time series is generated with a random seed, we need a vector of N unique seeds, which can be provided using the `seeds` argument. In case `seeds` is not provided, it is generated based on the provided means (`seeds.from.means`), and if it is a scalar, the seeds vector is created by adding the scalar to the `seeds.from.means` (which would allow generating different datasets with the same mean distributions).

2.1.3.2 Determining level-2 distribution parameters

For each alternative DGM—the χ^2 AR(1), BinAR(1), and PoDAR(1) models—we should determine appropriate parameters for the level-2 distribution of means such that we have enough skewness in the generated datasets. To do so, we make a function (`Mean.vs.Skewness`) to help us experiment with different values for μ and σ^2 (of the Gaussian level-2 distribution) and ν (of the χ^2 level-2 distribution) for each alternative DGM. Note that we start by generating more than enough samples for each distribution ($10 \times N$) and subsample N values after applying the model-specific lower and upper bounds.

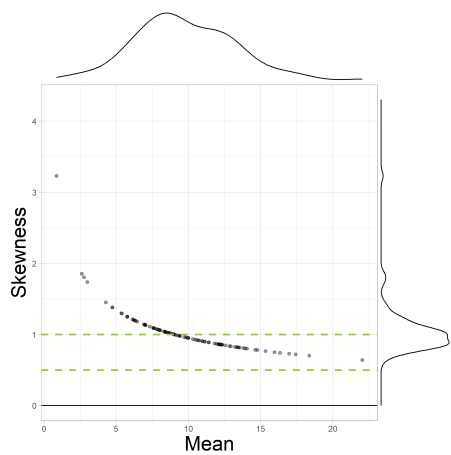
We notice that we get the desired distribution of skewness with the following parameters:

Table 1: Parameters of level-2 distribution of means

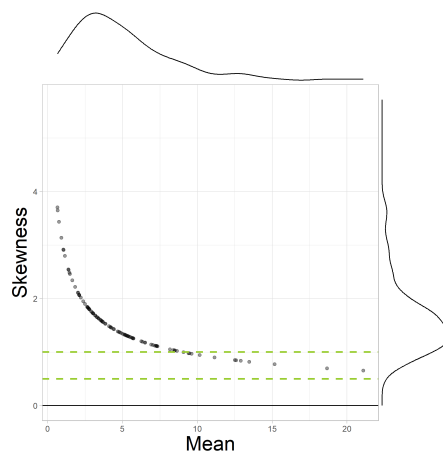
Model	μ	σ^2	ν
$\chi^2\text{AR}(1)$	10	10	5
BinAR(1)	2	1	2.9
PoDAR(1)	4	4	1.5

Giving us the following distributions:

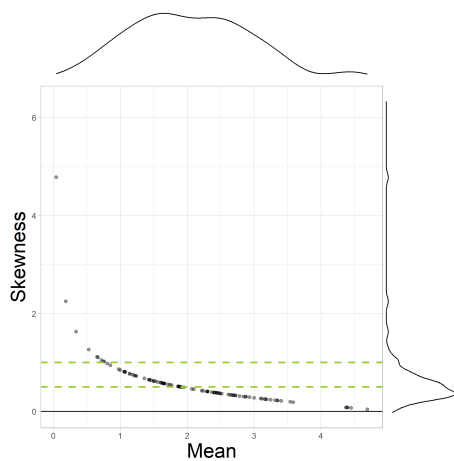
$\chi^2\text{AR}(1)$ with $\mu_i \sim N(10, 10)$



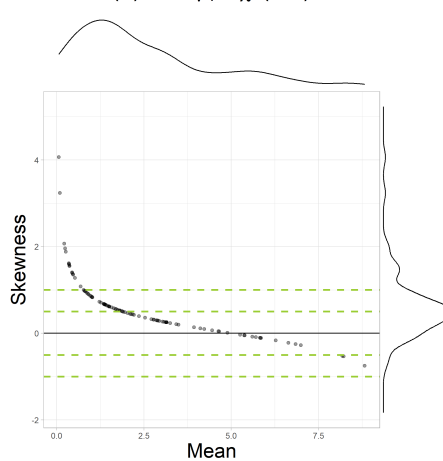
$\chi^2\text{AR}(1)$ with $\mu_i \sim \chi^2(5)$



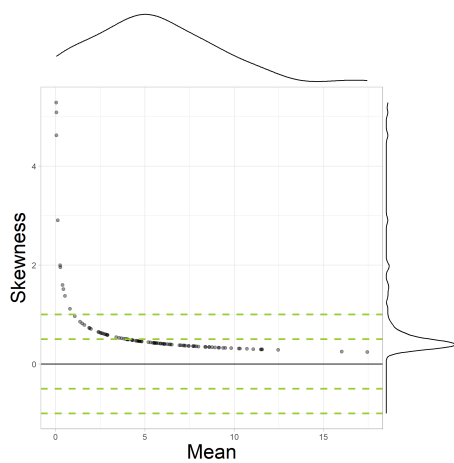
$\text{BinAR}(1)$ with $\mu_i \sim N(2, 1)$



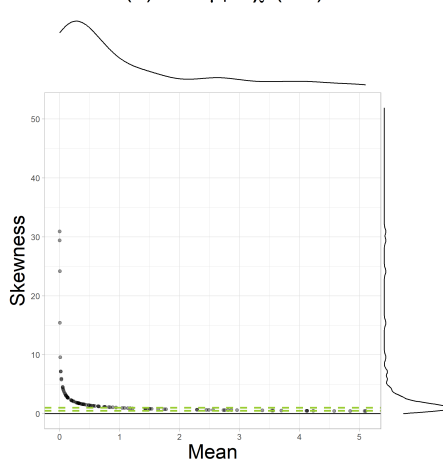
$\text{BinAR}(1)$ with $\mu_i \sim \chi^2(2.9)$



$\text{PoDAR}(1)$ with $\mu_i \sim N(4, 20)$



$\text{PoDAR}(1)$ with $\mu_i \sim \chi^2(1.5)$



2.1.3.3 Automate dataset generation

We then use a wrapper (`make_datasets`) around `dgm_make.sample` that generates datasets for all four DGMs with the appropriate level-2 parameters specified above. Note that here we first generate more than enough (i.e., $2 \times N$) samples of means to make sure we end up with N samples after applying the upper and lower bounds. The datasets are then saved, with some additional variables in separate `.rds` files, using wrapper functions described in Section 3.2.

2.2 The *Analysis* component

We analyzed each dataset with AR(1) models with fixed and random residual variance using *Mplus* v. 8.1 (Muthén and Muthén 2017). To interface *Mplus* from R, we used the package *MplusAutomation* (Hallquist and Wiley 2018) and wrote a function (`run_MplusAutomation`) that for each iteration of each condition would save the dataset as a `.dat` file, generate the `.inp` input script for the desired analysis type, and run the model for that dataset. *Mplus* then saves the output files (`.out` and `.gh5`). After the analysis, `run_MplusAutomation` extracts parameter estimates and returns them in an R object, which, with some additional variables, are saved in separate `.rds` files using wrapper functions described in Section 3.2.

In each analysis, we simulated two MCMC chains (`CHAINS = 2`), and to reduce autocorrelation in the estimated parameters, by defining `THIN = 5` we asked *Mplus* to save every 5th sample. By setting `BITERATIONS = 5000(2000)`, we made sure to have between 2000 to 5000 samples (after thinning) for each parameter from each chain. *Mplus* considers the first half of each chain as burn-in samples and discards them, thus, in total, we got at least 2000 “independent” samples from both chains combined. Finally, with `FACTORS = ALL (500)` we asked *Mplus* to draw 500 samples for each individuals when estimating level-1 parameters. We visually inspected the traceplots and autocorrelation plots of parameter estimates and of a sample of analyzed datasets and good convergence was observed. Furthermore, to make sure the number of iterations and thinning used in the analyses provide sufficiently accurate estimates, we re-analyzed two replications of each alternative DGM with Gaussian and χ^2 -distributed means with $N = 100$ and $T = 100$ with `BITERATIONS = 12500(5000)` and `THIN = 20`, which led to estimates of the parameters of interest (`unstd X.WITH.PHI`, `unstd Variances.PHI`, and `stdyx X.WITH.PHI`) almost identical (up to the third decimal) to those estimated with `BITERATIONS = 5000(2000)` and `THIN = 5` (see below).

The generated input files looked like the following. Note that the `TITLE` and `DATA` strings in the `.inp` files are unique to the dataset being analyzed and included the unique dataset seed `uSeed` (passed to `make_datasets` to generate datasets), the number of individuals in the dataset `N`, the length of the time series `T`, the model types used (`resid.fixed` or `resid.random`), and the replication number `Rep` (see Section 3.1 for further details).

2.2.0.1 Fixed residual variance

By estimating the covariances between mean and autoregression (x phi WITH x phi under the BETWEEN% command), the following Mplus script runs an AR(1) with random effect mean and autoregressive parameter but with fixed effect residual variance.

```
TITLE:
fit_uSeed-000000_N-100_T-100_type-resid.fixed_Rep-1

DATA:
FILE = "fit_uSeed-000000_N-100_T-100_type-resid.fixed_Rep-1.dat";

VARIABLE:
NAMES = subject t x;
MISSING=.;
CLUSTER = subject;
LAGGED = x(1);
TINTERVAL = t(1);

ANALYSIS:
TYPE = TWOLEVEL RANDOM;
ESTIMATOR = BAYES;
PROCESSORS = 1;
CHAINS = 2;
THIN = 5;
BITERATIONS = 5000(2000);

MODEL:
%WITHIN%
phi | x ON x&1;
%BETWEEN%
x phi WITH x phi;

OUTPUT:
TECH1 TECH2 TECH3 TECH8 FSCOMPARISON STANDARDIZED STDYX STDY;

PLOT:
TYPE = PLOT3;
FACTORS = ALL (500);
```

2.2.0.2 Random residual variance

By estimating the logarithm of the residual variance at level 1 (by $\log v$ | x under the

%WITHIN% command) and estimating the covariances between the level-2 mean, autoregression, and residual variance (x phi logv WITH x phi logv under the %BETWEEN% command), the following Mplus script runs an AR(1) model with random effect mean, autoregressive parameter, and residual variance.

```
TITLE:
fit_uSeed-000000_N-100_T-100_type-resid.random_Rep-1

DATA:
FILE = "fit_uSeed-000000_N-100_T-100_type-resid.random_Rep-1.dat";

VARIABLE:
NAMES = subject t x;
MISSING=.;
CLUSTER = subject;
LAGGED = x(1);
TINTERVAL = t(1);

ANALYSIS:
TYPE = TWOLEVEL RANDOM;
ESTIMATOR = BAYES;
PROCESSORS = 1;
CHAINS = 2;
THIN = 5;
BITERATIONS = 5000(2000);

MODEL:
%WITHIN%
phi | x ON x&1;
logv | x;
%BETWEEN%
x phi logv WITH x phi logv;

OUTPUT:
TECH1 TECH2 TECH3 TECH8 FSCOMPARISON STANDARDIZED STDYX STDY;

PLOT:
TYPE = PLOT3;
FACTORS = ALL (500);
```

2.2.0.3 Comparing different ANALYSIS parameters

The table below shows the estimated parameters of two replications of with `BITERATIONS = 5000(2000)` and `THIN = 5` (specified in the table with `2k x 5`) and `BITERATIONS = 12500(5000)` and `THIN = 20` (specified in the table with `5k x 20`). The relevant parameters (level-2 correlations `X.WITH.PHI unstd`, covariances `X.WITH.PHI stdyx`, and variance `Variances.PHI unstd`) are highlighted with orange.

2.3 The *Harvesting* component

The harvesting component was directly implemented in the parallel function `do_harvest_doFuture` described in Section 3.2.

2.4 The *Reporting* component

The functions include ...

3 Pipeline

We implemented each of these tasks in separate functions that were essentially wrapper functions (with parallel-computing implementation) around the modular components. Using these wrapper functions, each replication of each simulated condition was saved in a separate `.rds` file. These data files were fed to the analysis wrapper function whose output was saved in separate `.rds` data files. To collect relevant parameter estimates, another wrapper function was used to read the data files and save the desired parameters in a dataframe, which then used in reporting.

3.1 Book-keeping

The outcomes of the components are saved in separate `.rds` files and indexed by unique, descriptive names, and each replication is given a unique numeric identifier that is also used as the random seed used to generate the dataset within each replication. The file names and address are stored in two dataframes along with model parameters used to generate each dataset, and these dataframes are used when reading and writing data files in other components.

3.1.0.1 Simulation file references

The function `make_sim_refs` makes a dataframe including

3.1.0.2 Analysis outputs file references

3.2 Wrapper functions

3.2.0.1 Simulation in parallel

3.2.0.2 Analysis and save in parallel

3.2.0.3 Harvesting in parallel

3.3 Pipeline code

3.4 Simulation results

4 Additional analyses and figures

4.1 Profiles of simulated datasets

COGITO data analysis

Bien, Jacob. 2016. “The Simulator: An Engine to Streamline Simulations.” *arXiv:1607.00021 [Stat]*, June. <http://arxiv.org/abs/1607.00021>.

Hallquist, Michael N., and Joshua F. Wiley. 2018. “*MplusAutomation* : An R Package for Facilitating Large-Scale Latent Variable Analyses in *M Plus*.” *Structural Equation Modeling: A Multidisciplinary Journal* 25 (4): 621–38. <https://doi.org/10.1080/10705511.2017.1402334>.

Muthén, Linda K., and Bengt O. Muthén. 2017. *Mplus User's Guide*. Eighth Edition. Los Angeles, CA: Muthén & Muthén.