

Gaussian distributional structural equation models: A framework for modeling latent heteroscedasticity

Luna Fazio and Paul-Christian Barkner

Department of Statistics, TU Dortmund University, Dortmund, Germany

ABSTRACT

Accounting for the complexity of psychological theories requires methods that can predict not only changes in the means of latent variables – such as personality factors, creativity, or intelligence – but also changes in their variances. Structural equation modeling (SEM) is the framework of choice for analyzing complex relationships among latent variables, but the modeling of latent variances as a function of other latent variables is a task that current methods only support to a limited extent. In this article, we develop a Bayesian framework for Gaussian distributional SEM, which broadens the scope of feasible models for latent heteroscedasticity. We use statistical simulation to validate our framework across four distinct model structures, in which we demonstrate that reliable statistical inferences can be achieved and that computation can be performed with sufficient efficiency for practical everyday use. We illustrate our framework's applicability in a real-world case study that addresses a substantive hypothesis from personality psychology.

KEYWORDS

Structural equation modeling; distributional regression; Bayesian inference; heteroscedasticity; measurement invariance

1. Introduction

Structural equation modeling (SEM) is a widely-used statistical framework that can be regarded as an extension of regression models: it allows modeling multiple dependent variables simultaneously, including relationships among them, as well as the introduction of measurement error and unobserved (latent) variables (for a comprehensive introduction see Bollen, 1989; Kline, 2016). As with regression, the classic formulation of SEM presents an idealized setting where, among other simplifications, it is assumed that the model's parameters (intercepts, coefficients and (co-)variances) all take on values that are constant across people, conditions, etc. Such an assumption often does not hold in practice and this has motivated a rich literature on methods for handling non-invariant parameters (see below). Extending this line of research, we develop a Bayesian framework for Gaussian distributional SEMs, which, compared to past approaches, supports more flexible models of latent heteroscedasticity when dependencies on other latent variables are involved. We demonstrate our framework's statistical validity and usefulness through simulation studies on four distinct structural models

and a real-world case study applied to a research question from personality psychology.

1.1. Related work

The problem of invariance has received attention since the early days of factor analysis, initially focusing on invariance of the covariance matrix of observed data across selected subgroups of some larger population (Thomson & Lederunn, 1939). The introduction of the multiple-group model methodology in Jöreskog (1971) marked the shift in focus to the invariance of model parameters that prevails today. It was followed by the development of *moderated factor analysis* (MFA), which enabled modeling parameter values *via* known functions of observed variables (moderators), including continuous ones; this meant that evaluation of invariance stopped being limited to comparisons over discrete groups (Bauer & Hussong, 2009). A further extension, *local structural equation modeling* (LSEM), fits the model multiple times over the moderators' range in combination with an observation weighting scheme to produce a nonparametric estimate of the moderation functions, thereby avoiding the assumption of a known functional form

CONTACT Luna Fazio fazio@statistik.tu-dortmund.de Department of Statistics, TU Dortmund University, Germany

© 2025 The Author(s). Published with license by Taylor & Francis Group, LLC

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or with their consent.

(Hildebrandt et al., 2016). There are additional approaches that are particularly suited for assessment of non-invariance under specific assumptions of magnitude or structure (for an overview see Leith et al., 2023), but we do not discuss them here as they are less related to our proposed framework.

The above-mentioned techniques already provide a great deal of flexibility for modeling varying parameters within the SEM framework, but they all share the requirement that the moderator be an observed variable. A set of related approaches known as *heteroscedastic factor models* use MFA-like regressions on residual item variance and factor loading parameters together with skewed latent variable distributions (Molenaar et al., 2010, 2011). Another approach introduced in Molenaar (2015) is to use latent skewed distributions to allow the model to account for the effects of continuous latent moderators on the latent trait of interest. However, it achieves so by effectively marginalizing over the moderator and hence is not applicable when one wishes to include a measurement model for the latent moderator. In his discussion, Molenaar mentions this limitation and notes that models with explicit latent moderators would constitute a useful addition to the literature, citing methods for investigating latent heteroscedasticity (Molenaar et al., 2012) and latent variable interactions (Klein & Moosbrugger, 2000) as examples.

Indeed, one can already find some developments toward the use of latent predictors for latent variances. For instance, the works of Nestler (2020) and Martin and Rast (2022), motivated from the perspective of measurement reliability, provide techniques for modeling the variance of measurement errors (which can be conceived of as a special type of latent variable) as dependent on other latent variables. The original formulation of MFA (Bauer & Hussong, 2009) also received an extended treatment in Bauer (2017), where it is emphasized that the method can be used to assess measurement invariance and differential item functioning, including the case of both observed and latent moderators of item-level residual variances. Moving beyond heteroscedastic errors, modeling of the residual variance of a structural latent variable has also been demonstrated under a frequentist framework in a simple latent regression setting (i.e. one exogenous latent variable predicting one endogenous latent variable; de Kort et al., 2017).

One key challenge in maximum likelihood estimation of SEM is that latent variables can be regarded as *incidental* (Neyman & Scott, 1948) or *nuisance* (Basu, 1977) parameters, which means that they must be

marginalized out of the likelihood in order for consistent estimates to be obtainable. When latent heteroscedasticity is introduced, a closed-form expression of the marginal likelihood will generally not be available and numerical integration has to be performed at each step of the maximization procedure (e.g. Hessen & Dolan, 2009). Such an approach is sometimes called Marginal Maximum Likelihood (MML) and corresponds to an application of the more general expectation-maximization (EM) algorithm (Bock & Aitkin, 1981). As the quadrature methods that are commonly used to approximate the integral do not scale well with dimension (which in turn grows with the number of latent variables), de Kort et al. (2017) have suggested that Bayesian procedures could provide a viable alternative to MML for estimation of larger models with latent heteroscedasticity. To our knowledge, a systematic assessment of such an approach has not yet been conducted.

1.2. Our contributions

We develop and validate a Bayesian approach to support latent moderators of latent variances, which works by including latent variables as parameters to sample from instead of marginalizing over them. Such an approach has been termed *conditional likelihood* in the latent variable literature (e.g. Merkle et al., 2019), and it was favored in earlier methods for obtaining full posterior distributions in Bayesian SEM as it enabled the use of Gibbs sampling (Lee, 2007). With the development of algorithms such as Hamiltonian Monte Carlo (HMC; Neal, 2012; Betancourt, 2018), it was no longer necessary to use conditional distributions that could be sampled from and contemporary Bayesian SEM software has moved to use marginal likelihoods due to the increased sampling efficiency gained by not having the latent variables as additional parameters (Merkle et al., 2021). As mentioned above, however, latent moderators cannot be handled in full generality when using marginal likelihoods, which is why we adopt a conditional likelihood approach in this paper.

We implement our framework in the probabilistic programming language Stan, which provides an expressive syntax and powerful algorithms to specify and fit open-ended Bayesian models (Stan Development Team, 2023). To avoid users having to interact with Stan directly, we extended the R package `brms`, designed to simplify the process of fitting Bayesian regression models in Stan while still providing access to advanced regression techniques that can

be combined in a modular fashion (R Core Team, 2023; Børkner, 2017). We realize latent variable models with moderators in `brms` by utilizing its functionality for model-based imputation and distributional regression models, models predicting distributional parameters beyond the mean, for example, also variances or standard deviations (see Fahrmeir et al., 2021, Chapter 10; Børkner, 2018, 2021). By representing latent variables as missing observations and placing them as predictors of distributional parameters, we obtain an MFA-like procedure that admits latent moderators with more flexibility than methods based on marginal likelihoods.

In the remainder of this article we describe and evaluate our conditional likelihood approach for continuous latent moderators on both latent means and variances. In Section 2, we formally introduce the model and establish the corresponding notation. In Section 3, we present a large-scale simulation study to evaluate our approach, with results showing good convergence and parameter recovery in all investigated models. We demonstrate an application to a substantive hypothesis from personality psychology in Section 4. Finally, in Section 5, we discuss limitations and future directions.

2. Model description

Below, we formally introduce the developed SEM framework. Going forward, we will make an important simplifying assumption: all the variables in the model are conditionally normally distributed. This is not an inherent limitation of the approach we present, as it allows the specification of any continuous distribution for the latent variables, with moderation on other parameters beyond the mean. However, we find that this simplified setting already involves enough complexity for a rich discussion and practical relevance, so we omit a more general treatment in order to keep a reasonable scope for this paper. Additionally, we will omit structural manifest variables and fixed covariates from the following exposition as their inclusion is straightforward and our interest here is to discuss latent-to-latent regressions.

2.1. Model likelihood

We begin by describing the general structure of the model. Let I be some set indexing individual observations over the relevant units of analysis (e.g. institutions, individuals, time points, etc.). For each $i \in I$, we have a vector $\mathbf{f}_i = (f_{i1}, \dots, f_{iJ}, \dots, f_{iL})'$ of latent variables

and for each f_{ij} , the vector $\mathbf{y}_{ij} = (y_{ij1}, \dots, y_{ijM_i}, \dots, y_{ijM})'$ holds the corresponding manifest indicator variables. Here, M denotes the number of manifest variables of the i th factor, with M being allowed to vary over I . Then, the distribution of the variables can be written as

$$\mathbf{f}_{ij} | \mathbf{f}_i \sim \text{Normal}(\mu_{ij}, \Sigma_{ij}) \quad (1)$$

where μ_{ij} is the intercept of the manifest variable and Σ_{ij} is its factor loading. Both the mean μ_{ij} and standard deviation Σ_{ij} of each latent variable can depend on other latent variables. We consider dependencies of the form given by a generalized additive predictor

$$g_{h,i}(\mathbf{f}_i) = \sum_{k \in K_h} b_{kh_i} f_{kh_i} \quad (2)$$

where h stands for the likelihood parameter of interest (μ or Σ) and b_{kh_i} are the coefficients for each continuous (possibly non-linear) transformation f_{kh_i} of the latent variables. For clarity, we point out that if one wishes to include an intercept in the model, this can be done by setting $f_{0h_i} = 1$ and incorporating fixed covariates more generally is a matter of putting their values as a constant part in the f_{kh_i} . With this notation, each parameter is related to its predictor *via* the appropriate link function:

$$\mu_{ij} = g_{h,i}(\mathbf{f}_i) \quad (3)$$

Let us set $\mathbf{h}_F = (b_{K1}, \dots, b_{K1}, b_{1r_1}, \dots, b_{K1r_L})'$ to denote the vector of structural parameters and $\mathbf{h}_Y = (b_{11}, \dots, b_{LM}, k_{11}, \dots, k_{LM}, s_{11}, \dots, s_{LM})'$ to denote the vector of measurement model parameters. Then, the full likelihood can be written as

$$p(\mathbf{y} | \mathbf{f}, \mathbf{h}_F, \mathbf{h}_Y) \quad (4)$$

Following terminology from the latent variable model literature (e.g. Merkle et al., 2019), one could obtain the *marginal likelihood* by integrating out the latent variables:

$$p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{f}, \mathbf{h}_F, \mathbf{h}_Y) p(\mathbf{f}) d\mathbf{f} \quad (5)$$

The use of marginal likelihoods is a necessity in frequentist settings as latent variables play the role of incidental parameters, which results in inconsistent estimates if they are included in the estimation process (Neyman & Scott, 1948; also see discussion at the end of Hessen & Dolan, 2009). On the other hand, there are no formal impediments for performing Bayesian inference while including latent variables as part of the model parameters. The form of the

likelihood in which latent variables are explicitly included is called the *conditional likelihood*, as one can decompose Equation (4) into a likelihood for the indicator variables conditioned on the latent variables, and a likelihood for the latent variables themselves:

$$p(y|f)h_f, h_y \propto p(y|f)h_f, h_y p(f) \quad (6)$$

As discussed in Section 1, we use conditional likelihoods in this paper because marginalization would not produce a closed-form expression in the presence of latent predictors for latent variances.

Because we use the Bayesian framework for inference (see the Estimation section for more information), a complete specification must also include priors for the parameters. We can write the resulting joint posterior as

$$p(y|f, h_f, h_y) \propto p(y|f, h_y) p(f|f_{-j}) p(h_f) p(h_y) \prod_{j=1}^J \frac{p(y_j|f_j, h_{y_j})}{\text{PA}(f_j|f_{-j}, h_{y_j})} \quad (7)$$

where $\text{PA}(f_j|f_{-j})$ denotes the *parents* of latent variable f_j among the set of all other latent variables f_{-j} ; that is, all latent variables that contribute to the additive predictors l_j or r_j .

2.2. Identification

The model as given above is underidentified. Unless otherwise noted, identification for models in this paper is obtained by setting the expectation of all latent variables to 0, which identifies their mean, and the loading factor of one item to 1, which identifies their scale (see (Bollen, 1989, p. 238) for an introduction to identification in SEM). Our model additionally introduces coefficients for the latent variance linear predictor, so it is valid to ask whether these parameters are identified too. Below, we provide a formal argument demonstrating that a link from observed data to parameter values can be drawn without the need for any new constraints, thus showing identification.

To start with, consider the simplified scenario where we assume the latent values to be observed directly. Let f_0 be the variable whose variance we are interested in predicting based on the values of f_1, \dots, f_K so that

$$\text{Var } f_0 \propto \exp \left(\sum_{k=1}^K b_k f_k \right), \quad (8)$$

which means that the coefficients are related to the ratio of variances given a unit increase in one of the predictors. Without loss of generality, consider

increasing f_1 by one. Then, we find

$$\frac{\text{Var } f_0 | f_1, f_2, \dots, f_K}{\text{Var } f_0 | f_1, f_2, \dots, f_K} \propto \frac{\exp \left(\sum_{k=1}^K b_k f_k \right)}{\exp \left(\sum_{k=1}^K b_k f_k \right)} \quad (9)$$

thus showing identification of the coefficients b_1, \dots, b_K if there is variation in the corresponding latent variables. Estimation of the coefficients in such a model is a well-studied topic (e.g. Harvey, 1976) and can be regarded as a particular case of the more general *distributional regression* framework (e.g. see Chapter 10 of Fahrmeir et al., 2021).

Coming back to our non-simplified model, we do not actually observe the latent variables but rather noisy measurements as defined in Equation (1), and we want to show whether it's possible to infer changes in the latent variance from those available observations. For this purpose, a more helpful way of writing the measurement model is

$$y_{lm} = k_{lm} f_l + e_{lm}, \quad e_{lm} \sim \text{Normal}(0, s_{lm}^2) \quad (10)$$

which can be combined with Equation (8) to untangle the latent variance from the error variance. Let us examine the observed variance of some measurement y_{0m} of f_0 ; conditional on observed measurements y_{f1}, \dots, y_{Kgm} for the latent predictors. For simplicity, and without loss of generality, we take a single measurement per latent variable so the m subscript is dropped. This gives us

$$\text{Var } y_{0j} | y_{f1}, \dots, y_{Kgm}, e_{f1}, \dots, e_{Kg} \propto \frac{\text{Var } f_0 \text{Var } e_0}{k_0^2 \text{Var } f_0 + s_0^2} \quad (11)$$

$$\propto \frac{k_0^2 \exp \left(\sum_{k=1}^K b_k f_k \right) s_0^2}{k_0^2 \exp \left(\sum_{k=1}^K b_k f_k \right) + s_0^2} \quad (12)$$

The last expression still contains unobserved variables in the form of the error terms e_{f1}, \dots, e_{Kg} ; but we can apply the law of iterated expectations to deal with them. We use Equations (11) and (12) to work in terms of the latent variance we are interested in:

$$\begin{aligned} \text{Var } f_0 | y_{f1}, \dots, y_{Kg} &\propto \text{E} \left[\text{Var } f_0 | y_{f1}, \dots, y_{Kg}, e_{f1}, \dots, e_{Kg} \right] \\ &\propto \text{E} \left[\frac{\exp \left(\sum_{k=1}^K b_k f_k \right) s_0^2}{k_0^2 \exp \left(\sum_{k=1}^K b_k f_k \right) + s_0^2} \right] \\ &\propto \text{E} \left[\frac{\exp \left(\sum_{k=1}^K b_k f_k \right) s_0^2}{k_0^2 \exp \left(\sum_{k=1}^K b_k f_k \right) + s_0^2} \right] \end{aligned} \quad (13)$$

where the last step uses the fact that the expectation of the exp-sum of error terms is equivalent to the product of expectations of independent log-normal random variables.

With Equations (11) and (13), we can obtain an expression that is analogous to Equation (9), but fully expressed in terms of observable measurements (common factors are omitted):

[illegible]

Hence, observable changes in the variance of the measurement y_0 across measurements y_k of the latent variance predictors provide sufficient information to estimate the coefficients and no additional identification constraints are required because s_0^2 and k_1 are already identified by the usual constraints on the measurement model.

2.3. Estimation

We use Bayesian inference for model fitting. At a high level, the process consists of first specifying a *prior distribution* (further discussed in the next subsection), which describes our state of knowledge before seeing the data, and combining it with the data-informed model likelihood to obtain a *posterior distribution*, which represents our updated state of knowledge about the parameters' values. An accessible introduction to Bayesian inference can be found in Johnson et al. (2022).

Calculating the posterior distribution is the main challenge during inference as the expression involves a high-dimensional integral which will not have a closed-form beyond a few special cases; hence, it becomes necessary to resort to numerical methods. We use Markov chain Monte Carlo (MCMC), specifically adaptive Hamiltonian Monte Carlo as implemented in the Stan probabilistic programming language (Hoffman et al., 2014; Stan Development Team, 2023). Adaptive HMC is a class of efficient algorithms that can accurately sample complicated parameter spaces and Stan is a well-tested project that is freely available for all major operating systems. All MCMC algorithms produce sequences of samples (known as *chains*) from the target distribution as its output, which we can then directly use to obtain estimates of parameter means, credibility intervals, transformations, and other quantities of interest (Gelman et al., 2014).

2.4. Prior specification

In the ideal Bayesian workflow, all model parameters are given priors that represent some state of knowledge which will be updated through the likelihood as new data arrives. The purpose of this paper, however, is to investigate the set of conditions under which our approach can produce useful results. Hence, we adopt the *minimalist* position (Gelman et al., 2017) for all simulations, i.e. we attempt to identify the weakest priors for each model that will still produce reliable inferences. The criteria we use to assess reliability are described in Section 3.1 and the specific priors are introduced along with their corresponding models in Section 3.3. Readers looking for practical advice on how to set priors for SEM can find an excellent resource in van Erp (2020) and Winter and Depaoli (2023).

3. Simulations

We investigated the viability of our approach through statistical simulation. Specifically, we tested four structural models that are likely to be relevant for practitioners: a simple two-factor model, a model with mediators, a model with interactions and a model with a sequential structure. The metrics used for assessment are introduced next, followed by a description of the computational setup, and then each model is presented together with the respective results.

3.1. Model diagnostics

3.1.1. Convergence

While MCMC methods can work well in practice, convergence to a target distribution is an asymptotic property, so it is always necessary to verify convergence empirically (Bürkner et al., 2023). This can be achieved by running multiple chains with randomized initial values and then examining whether they exhibit similar distributions; one commonly recommended convergence diagnostic is the *potential scale reduction factor* \hat{R} (often just called “Rhat”). Briefly, it compares the variance between and within chains as a proxy for convergence and returns a value in $[1/2, 1]$ where values closer to 1 indicate the chains have more similar distributions. A detailed treatment can be found in Vehtari et al. (2021), where they also provide the recommendation to consider $\hat{R} \leq 1.01$ as a reliable indicator of convergence. For the purpose of our simulation study, we relaxed the threshold to 1.05, as we have access to the ground truth values and

Figure 1. Simulation-based rank histograms (top) and corresponding empirical cumulative distribution function (ECDF) difference plots (bottom) for three hypothetical quantities of interest. The blue areas in the ECDF difference plots indicate 95%-confidence intervals under the assumptions of uniformity and thus allow for a null-hypothesis significance test of self-consistent calibration. Left: A well-calibrated quantity. Center: A miscalibrated quantity with too many lower ranks indicating a positive bias in the estimated posteriors. Right: A miscalibrated quantity with too many extreme ranks indicating overconfident posteriors (i.e. variance underestimated).

therefore were able to verify that posterior estimates retained acceptable quality up to that point.

3.1.2. Calibration

Convergence alone does not tell us whether our MCMC draws provide a good approximation to the true posterior distribution. However, we can use the draws themselves to diagnose the quality of our approximation if we also have knowledge of the true data-generating distribution; this is the key idea behind *Simulation-Based Calibration* (SBC; Talts et al., 2020). For this method, one samples parameters from the prior which are passed to the likelihood for data generation, the model is then fit over the resulting datasets, and the sum of ranks of the posterior draws relative to the true value is calculated; when a uniform distribution of rank sums is recovered, our posterior is said to be *calibrated* (explained below). To assess uniformity, we used the graphical tests proposed by Sallin et al. (2022) (see Figure 1).

Calibration in the context of SBC is, strictly speaking, a statement about the expected coverage of posterior intervals over the joint distribution of data and parameters. In practice, this means that it can readily detect posterior approximations that consistently under-/overestimate the location or uncertainty that the true posterior would output for a given parameter; however, it can miss less obvious mismatches and hence does not provide a global guarantee of correctness (Modrak et al., 2023). Fortunately, the procedure can be augmented with data-dependent quantities to provide a more stringent test; in particular, we also test the model likelihood, which greatly increases the

sensitivity of the test as also demonstrated by Modrak et al.

3.1.3. Effective sample size and efficiency

Even if the model is calibrated and has converged, we only have a finite sample of MCMC draws from the posterior, so we must ensure that the estimation error is small enough to give us reliable inference. It is also necessary to account for the autocorrelation that is often present in the chains as this further reduces their information content; Effective Sample Size (ESS) is a diagnostic that addresses this by estimating the number of independent draws that the information in our chains is equivalent to. Section 11.5, Gelman et al. (2014). We consider an ESS of at least 100 per independent MCMC chain to be sufficient for reliable estimation and separately report bulk ESS and tail ESS, as suggested by Vehtari et al. (2021).

As having a high enough ESS is a prerequisite for accurate inference, a question of practical importance is how long one has to run a model for in order to achieve the desired precision. We calculate ESS per second (ESS/s), as it provides a simple measure of sampling efficiency for each model. However, this will vary considerably depending on the priors used, the data at hand, and the computer one uses to fit the model; the intent here is only to determine whether the models can run in a reasonable time.

3.1.4. Parameter recovery

To evaluate parameter recovery, we use bias and the Root Mean Squared Error (RMSE). Given a set $h^{(i)}/s_i$ of S posterior draws and a true value h ; we have

$$\text{Bias} = \frac{1}{S} \sum_{s=1}^S \bar{h}_{i,j}^{(s)} - h_{i,j}, \quad (15)$$

$$\text{RMSE} = \sqrt{\frac{1}{S} \sum_{s=1}^S (\bar{h}_{i,j}^{(s)} - h_{i,j})^2} \quad (16)$$

Posterior means will almost surely (in the formal sense) have non-zero bias whenever proper priors are used. However, they are also consistent estimators and we show that, for our models, and this leads to the bias being negligible. As bias does not account for posterior uncertainty, we also report RMSE because it provides an overall indication of estimation error by incorporating both posterior bias and variance into a single measure. This relation can be shown explicitly by rearranging Equation (16):

$$\text{RMSE}^2 = \frac{1}{S} \sum_{s=1}^S (\bar{h}_{i,j}^{(s)} - h_{i,j})^2 = \frac{1}{S} \sum_{s=1}^S (\bar{h}_{i,j}^{(s)} - h_{i,j} + h_{i,j} - h_{i,j})^2 = \frac{1}{S} \sum_{s=1}^S (\bar{h}_{i,j}^{(s)} - h_{i,j})^2 + \frac{1}{S} \sum_{s=1}^S (h_{i,j} - h_{i,j})^2 + 2 \frac{1}{S} \sum_{s=1}^S (\bar{h}_{i,j}^{(s)} - h_{i,j})(h_{i,j} - h_{i,j}) \quad (17)$$

Variance

where $\bar{h}_{i,j} = \frac{1}{S} \sum_{s=1}^S \bar{h}_{i,j}^{(s)}$ is the posterior mean.

3.2. Computational setup

The workflow of our simulation study can be summarized in four key steps:

1. For each model, we found a relatively tight *generative* prior distribution, such that parameter vectors drawn from the joint prior could be used to simulate datasets without degeneracies (e.g. without values close to zero for variances or factor loadings) with high probability.
2. We drew 250 parameter vectors from the generative prior and with each of these, we subsequently generated a dataset of 500 observations from the model likelihood.
3. We fitted the model twice, using two different priors, for each of the generated datasets: first with the generative prior itself and second with a much wider, weakly informative prior.
4. Finally, we processed the resulting posterior samples to obtain model diagnostics in the following manner:

Convergence was examined in both sets of models. Below, we only report convergence for the models with the weakly informative prior, as convergence for the models with the generative prior was always superior (see Appendix A).

Calibration was assessed on the models that used the generative prior, because only there SBC is valid.

Metrics for parameter recovery and sampling efficiency were calculated from the models with the weakly informative prior. We excluded models with any $\hat{R} > 1.05$ from this calculation to avoid artifacts caused by clear non-convergence.

Our simulations were fully implemented using the R programming language (R Core Team, 2023). We specified our models in the *brms* package (Bürkner, 2017), which provides a user-friendly interface for generation of Stan code, and wraps the *cmdstanr* and *posterior* packages, which respectively provide functions for interfacing with Stan itself and for extracting model diagnostics (Gabry et al., 2025; Bürkner et al., 2023a). Functions needed for dataset generation, linking true parameter values to specific fits and plotting calibration diagnostics were provided by the *SBC* package (Kim et al., 2023). To facilitate reproducibility, the simulation pipeline itself was built to run *via* the *targets* package (Landau, 2021). The full code is available at an online repository.¹ The complete pipeline was run on a MacBook Pro with M2 chip, where it took approximately 42 h to complete.

3.3. Results

To aid visualization of the models that follow, we introduce a novel graphical representation to complement the previously established notation. As our models allow the mean (μ) and standard deviation (σ) of a latent variable distribution to be independently influenced by other latent variables, we extend the usual path diagram notation to show both of these parameters explicitly.

The diagram for a single latent variable with a five-item measurement model is shown in Figure 2. This corresponds to the following statistical model:

$$\begin{aligned} f_1 & \sim \text{Normal}(\mu_1, \sigma_1^2) \\ y_{1m} & \sim \text{Normal}(\mu_m, \sigma_m^2) \end{aligned} \quad (18)$$

with $m \in \{1, \dots, 5\}$ and $k_{11} = 1$ for identification. All the models in this section use five items per latent variable and a unit factor loading identification constraint, as in the example above. For brevity, we omit the measurement models in the descriptions that follow.

¹See the folder simulation-study at <https://github.com/bdlvm-project/gdsem-paper>

Figure 2. Extended path diagram notation which explicitly shows the parameters (l and r) that determine the latent variable's distribution.

Figure 3. Two-factor model. f_1 influences both the mean and standard deviation of f_2 :

Table 1. Prior specifications for the two-factor model.

Parameter type	Notation	Generative prior	Weakly informative prior
Latent mean			
Slope	b_{11_2}	Normal(1, 0.3)	Normal(0, 2.5)
Latent std. dev.			
Initial	r_1	Gamma(0.71, 11.1)	Gamma(5, 5)
Intercept	b_{0r_2}	Exp-Gamma(1.11, 1)	Exp-Gamma(5, 5)
Slope	b_{1r_2}	Normal(0.15, 0.05)	Normal(0, 0.5)
Item parameters			
Factor loadings	k	Normal(1, 0.3)	Normal(0, 2.5)
Error std. dev.	s	Normal(0.31, 0.15)	Gamma(2.5, 5)

3.3.1. Two-factor model

We start with a two-factor model as this is the simplest setup where we can have a latent variance that is conditional on the value of another latent variable. The mathematical notation for the model is

$$\begin{aligned}
 f_1 & \sim \text{Normal}(0, r_1) \\
 f_2 & \sim \text{Normal}(l_2, r_2) \\
 l_2 & \sim b_{11_2} f_1 \\
 r_2 & \sim b_{0r_2} + b_{1r_2} f_1
 \end{aligned} \tag{19}$$

The path diagram representation is shown in Figure 3. The priors used for generation and fitting are described in Table 1.

The calibration plots in Figure 4(a) show that all model parameters and the log-likelihood are well-calibrated.

The heatmap in Figure 4(b) shows that about two-thirds of the simulations produced estimates with \hat{R} between 1.00 and 1.01. Factor loadings (k) appear to be the parameters most likely to present suboptimal convergence but major problems are rare, with only two simulations having $\hat{R} > 1.05$. We did not find any feature of the simulations that distinctly explained the variation in \hat{R} which suggests it is primarily caused by variations in the random initializations of parameter values. Such occasional convergence issues are known to occur when fitting Bayesian SEMs in general so attempting to refit with a different seed is an advisable first step.²

Finally, our parameter recovery plots (Figure 4(c)) show that bias is negligible across all model parameters. The values for RMSE need to be interpreted in the context of each parameter's relevant scale: for the slope on the standard deviation b_{1r_2} , the parameter we are most interested in, an RMSE of 0.05 is low enough to expect that the model can be usefully employed to obtain directional estimates. This result is encouraging given that it comes from datasets with only 500 observations; studies with a larger sample size would be able to produce even more precise inferences.

3.3.2. Mediation model

One common use of SEM is mediation analysis, which allows researchers to quantify direct and indirect effects for a given variable. We show that our approach allows investigating mediation for effects on both means and standard deviations by fitting the following model:

$$\begin{aligned}
 f_1 & \sim \text{Normal}(0, r_1) \\
 f_2 & \sim \text{Normal}(l_2, r_2) \\
 l_2 & \sim b_{11_2} f_1 \\
 r_2 & \sim b_{0r_2} + b_{1r_2} f_1 \\
 f_3 & \sim \text{Normal}(l_3, r_3) \\
 l_3 & \sim b_{11_3} f_1 \\
 r_3 & \sim b_{0r_3} + b_{1r_3} f_1 \\
 f_4 & \sim \text{Normal}(l_4, r_4) \\
 l_4 & \sim b_{11_4} f_1 + b_{21_4} f_2 \\
 r_4 & \sim b_{0r_4} + b_{1r_4} f_1 + b_{2r_4} f_2
 \end{aligned} \tag{20}$$

Here, f_1 has a direct effect on the l and r for f_2 , f_3 , and f_4 . Additionally, it has an indirect effect on l

²For instance, this is also recommended in the documentation for Bayesian SEM package blavaan.

Figure 4. Simulation diagnostics for the two-factor model. (a) ECDF difference plots. Curves are overlaid when there are multiple parameters of the same type. (b) Heatmap showing the average \hat{R} for each parameter type in each simulation. Simulations are arranged in ascending order across the x-axis according to their overall mean \hat{R} . (c) Box plots of the error distribution for average bias and average RMSE per simulation and parameter type. Simulations with convergence issues (any parameter with $\hat{R} > 1.05$) were excluded.

Table 2. Prior specifications for the mediation model.

Parameter type	Notation	Generative prior	Weakly informative prior
Latent mean			
Slope	b_{11}	Normal(1, 0.3)	Normal(0, 2.5)
Latent std. dev.			
Initial	Γ_1	Gamma(0.71, 1/11)	Gamma(5, 5)
Intercept	$b_{0r_2}, b_{0r_3}, b_{0r_4}$	Exp-Gamma(1, 1/11)	Exp-Gamma(5, 5)
Slope	$b_{1r_2}, b_{1r_3}, b_{2r_4}$	Normal(-0.15, 0.05)	Normal(0, 0.5)
	b_{1r_4}	Normal(0.15, 0.05)	
Item parameters			
Factor loadings	k	Normal(1, 0.3)	Normal(0, 2.5)
Error std. dev.	s	Normal(0.31, 0.5)	Gamma(2.5, 5)

Figure 5. Mediation model. f_1 has direct and indirect effects on both of f_4 's parameters.

of f_4 through f_2 and an indirect effect on r of f_4 through f_3 : The path diagram is shown in Figure 5. Priors for this model are shown in Table 2.

The calibration plots in Figure 6(a) show good calibration across all test quantities. Convergence in Figure 6(b) again appears to be good in general; there are 17 simulations with $\hat{R} > 1.05$, but we traced these cases back to generated datasets that produced very small variances in f_3 ; which led to unstable estimates for factor loadings and error variances. We verified that refitting these cases with a different initialization was sufficient to resolve the issue (results not shown).

Results for parameter recovery in Figure 6(c) show no significant bias and RMSE is only slightly increased for all slopes (on both l and r) in the model compared to the previous two-factor model. This is an expected consequence of including mediators in the model, as having multiple paths for a given effect widens the range of coefficient values that are compatible with the data. In general, assessing mediation imposes increases in sample size and methodological complexity (Rohrer et al., 2022; Montoya, 2023).

3.3.3. Interaction model

We have mentioned in the Introduction that the ability to study moderation in general is a feature of interest for users of SEM frameworks. Our approach was created with the intent of providing more flexible models for moderation on latent variances, but we

Figure 6. Simulation diagnostics for the mediation model. (a) ECDF difference plots. Curves are overlaid when there are multiple parameters of the same type. (b) Heatmap showing the average \hat{R} for each parameter type in each simulation. Simulations are arranged in ascending order across the x-axis according to their overall mean \hat{R} : (c) Box plots of the error distribution for average bias and average RMSE per simulation and parameter type. Simulations with convergence issues (any parameter with $\hat{R} > 1.05$) were excluded.

Table 3. Prior specifications for the interaction model.

Parameter type	Notation	Generative prior	Weakly informative prior
Latent mean			
Slope	b_{11_4} b_{21_4}	Normal(1, 0.3) Normal(0.5, 0.3)	Normal(0, 2.5)
Latent std. dev.			
Initial	$\Gamma_1, \Gamma_2, \Gamma_3$	Gamma(0.71, 111)	Gamma(5, 5)
Intercept	b_{0r_4}	Exp-Gamma(1, 11)	Exp-Gamma(5, 5)
Slope	b_{1r_4} b_{2r_4}	Normal(0.1, 0.05) Normal(0.05, 0.05)	Normal(0, 0.5)
Item parameters			
Factor loadings	k	Normal(1, 0.3)	Normal(0, 2.5)
Error std. dev.	s	Normal(0.31, 0.50:15)	Gamma(2.5, 5)

Figure 7. Interaction model. Dashed lines represent deterministic transformations; in this case, taking the product of two latent variables.

found that it can also accommodate moderation on structural paths with ease, since it can be represented through interactions between variables on the same linear predictor. We show this by fitting the moderation model depicted in Figure 7, mathematically expressed as:

$$\begin{aligned}
 f_1 & \sim \text{Normal}(\mu_{f_1}, \sigma_{f_1}^2) \\
 f_2 & \sim \text{Normal}(\mu_{f_2}, \sigma_{f_2}^2) \\
 f_3 & \sim \text{Normal}(\mu_{f_3}, \sigma_{f_3}^2) \\
 f_4 & \sim \text{Normal}(\mu_{f_4}, \sigma_{f_4}^2) \\
 \log r_4 & \sim \text{Normal}(\mu_{r_4}, \sigma_{r_4}^2) \\
 \log r_4 & \sim \text{Normal}(\mu_{r_4}, \sigma_{r_4}^2)
 \end{aligned} \tag{21}$$

The priors used to fit this model are given in Table 3.

The results in Figure 8 show that the model is well-calibrated and has good convergence overall. Parameter recovery is in line with the previous models, showing no evidence of bias as well as low RMSE.

3.3.4. Sequential model

Model structures that contain longer sequences of latent variables can be relevant in studies that involve measurements over time (Asparouhov et al., 2018) or in those which seek to study detailed causal structures (Zugna et al., 2022). A thorough investigation of models in that space is well beyond the scope of this paper, but we considered it relevant to at least explore whether issues with our approach could become

Figure 8. Simulation diagnostics for the interaction model. (a) ECDF difference plots. Curves are overlaid when there are multiple parameters of the same type. (b) Heatmap showing the average \hat{R} for each parameter type in each simulation. Simulations are arranged in ascending order across the x-axis according to their overall mean \hat{R} . (c) Box plots of the error distribution for average bias and average RMSE per simulation and parameter type. Simulations with convergence issues (any parameter with $\hat{R} > 1.05$) were excluded.

Figure 9. Sequential model.

apparent only when fitting a longer sequence of dependencies between latent variables. The model we constructed additionally shows that transformations of latent variables are also supported in this approach (Figure 9). Using mathematical notation:

$$\begin{aligned} f_1 & \sim \text{Normal}(0, 1) \\ f_j & \sim \text{Normal}(b_{0j}, r_j^2) \\ & \quad | f_{j-1} \sim \text{Normal}(b_{1j}, r_{j-1}^2) \\ & \quad \log r_j \sim \text{Normal}(b_{0rj}, b_{1rj} f_{j-1}^2) \end{aligned} \quad (22)$$

where $j \in \{2, 3, 4, 5\}$ for this case. Model priors are shown in Table 4.

Readers familiar with time series may note that setting an equality constraint on the coefficients would give the form of an autoregressive model, but testing this for lengths of time series representative of real applications would involve computational challenges that we do not aim to tackle here.

The sequential structure of this model meant that, for certain parameter draws from the generative prior, there was a runaway increase in latent variance.

Table 4. Prior specifications for the sequential model.

Parameter type	Notation	Generative prior	Weakly informative prior
Latent mean			
Slope	b_{1j}	$\text{Normal}(0, 0.2)$	$\text{Normal}(0, 2.5)$
Latent std. dev.			
Intercept	b_{0r-1}	$\text{Gamma}(0.71, 11.1)$	$\text{Gamma}(5, 5)$
Intercept	b_{0r-2}	$\text{Exp-Gamma}(1, 11)$	$\text{Exp-Gamma}(5, 5)$
Slope	b_{1r}	$\text{Normal}(0, 0.05)$	$\text{Normal}(0, 0.5)$
Item parameters			
Factor loadings	k	$\text{Normal}(1, 0.3)$	$\text{Normal}(0, 2.5)$
Error std. dev.	s	$\text{Normal}(0.31, 0.50)$	$\text{Gamma}(2.5, 5)$

Therefore, we had to drop one out of the 250 datasets due to overflowing variances which caused numerical errors. Similarly to the mediation model, there were also 18 datasets with unrealistic values for the variance parameters, which is reflected as a larger fraction of models with convergence issues in Figure 10. Nonetheless, calibration for this model was good, the vast majority of simulations converged well, and estimates don't show any major sign of bias. We did note a longer tail of high-RMSE estimates for the b_{1r}

Figure 10. Simulation diagnostics for the sequential model. (a) ECDF difference plots. Curves are overlaid when there are multiple parameters of the same type. (b) Heatmap showing the average \hat{R} for each parameter type in each simulation. Simulations are arranged in ascending order across the x-axis according to their overall mean \hat{R} . (c) Box plots of the error distribution for average bias and average RMSE per simulation and parameter type. Simulations with convergence issues (any parameter with $\hat{R} > 1.05$) were excluded.

coefficients, but the bulk of simulations remained at levels comparable to the previously tested models.

3.3.5. Comparing ESS

All the diagnostics discussed so far are based on a finite amount of samples drawn from the posterior and so, in principle, are subject to estimation error. The left side of [Figure 11](#) shows that across all models and parameters, the resulting ESS was well above the recommended threshold of 100 samples per chain, both for bulk and tail estimates, which establishes that our sample-based diagnostics are providing reliable information.

On the top right side of [Figure 11](#), we depict the ESS per second, which is a measure of the sampling efficiency of our models. As expected, increasing the number of latent variables in the model generally leads to decreases in sampling efficiency. Available alternatives for SEM based on marginal likelihoods do not allow latent moderation to be modeled with the same degree of flexibility as our approach, so a direct comparison cannot be made, but we consider that our conditional likelihood models are sufficiently fast for practical everyday use: the lower right corner of [Figure 11](#) shows that the vast majority of models took less than a minute to fit.

4. Case study

We sought to address a substantively relevant research question using dataset from a real study in order to demonstrate the usefulness and flexibility of our model. We reached out to the authors of Mader et al. (2023), as their paper investigates the association between neuroticism and intra-person variation in negative affect, both of which can be conceptualized as latent variables.

A key finding of Mader et al. is that floor effects must be accounted for to reliably detect the neuroticism-emotional variability association. They used a hierarchical distributional regression model, where each individual's mean negative affect score was treated as a censored outcome and both its mean and variance were regressed against their mean neuroticism score. This is already quite an advanced model, but it has some shortcomings as it ignores the uncertainty in the score means and assumes every item should be weighted equally. However, methods available at the time would not have allowed to properly model the latent nature of neuroticism and negative affect while simultaneously estimating how one affects the variability of the other.

One of the authors kindly provided us with access to a subset of the *Goettingen Ovulatory Cycle Diaries 2*

Figure 11. Top two plots show the distribution of average effective sample size (ESS) and average ESS per second overall simulations, after excluding those with convergence issues (any parameter with $\hat{R} > 1.05$). Lower right corner is the distribution of wall time of those simulations.

dataset.³ The recruitment and data collection process is described in Arslan et al. (2021). Briefly, women filled out a form upon recruitment, which included a personality questionnaire, and were subsequently invited to fill out an online diary every day for 70 days, which included items on loneliness, irritability, self-esteem, stress and mood. We used this set of items as measurements for the emotional affect construct and the personality items in the initial questionnaire for the neuroticism construct. The study used a planned missingness design, which allows us to drop all incomplete observations without risk of bias. This leaves a total of 1039 women, each observed between 1 and 11 times (mean: 2.5).

We aim to follow the model given in Mader et al. (2023) as closely as possible. Therefore, we consider item responses at the extremes of the scale as censored and use a hierarchical structure to account for the repeated within-person observations. Abbreviating neuroticism as *Ne* and emotional affect as *Em*, our

model can be notationally expressed as

$$\begin{aligned}
 f_{Ne,i} & \sim \text{Normal}(\mu_{Ne,i}, \sigma_{Ne,i}^2) \\
 f_{Em,i,j} & \sim \text{Normal}(\mu_{Em,i,j}, \sigma_{Em,i,j}^2) \\
 l_{Em,i} & \sim b_{0Em} + b_{1Em} f_{Ne,i} + c_{1Em,i} \\
 \log \sigma_{Em,i}^2 & \sim b_{0Em} + b_{1Em} f_{Ne,i} + c_{1Em,i} \\
 c_{1Em,i} & \sim \text{Normal}(0, \sigma_{c1Em,i}^2) \\
 c_{1Em,i} & \sim \text{Normal}(0, \sigma_{c1Em,i}^2) \\
 y_{Ne,n,i} & \sim \text{Normal}(\mu_{Ne,n,i}, \sigma_{Ne,n,i}^2) \\
 & \quad \text{if } y_{Ne,n,i} \geq 1 \text{ or } y_{Ne,n,i} \leq -1 \\
 & \quad \text{if } y_{Ne,n,i} > 5 \text{ or } y_{Ne,n,i} < -5 \\
 y_{Em,m,i,j} & \sim \text{Normal}(\mu_{Em,m,i,j}, \sigma_{Em,m,i,j}^2) \\
 & \quad \text{if } y_{Em,m,i,j} \geq 0 \text{ or } y_{Em,m,i,j} \leq -4 \\
 & \quad \text{if } y_{Em,m,i,j} > 4 \text{ or } y_{Em,m,i,j} < -4
 \end{aligned} \tag{23}$$

where *i* indexes individual participants, *j* indexes their responses over the study duration, $n \in \{1, \dots, 8\}$

³Codebook available at <https://rubenarslan.github.io/gocd2/>

Figure 12. For subject i , the baseline measurement of neuroticism is used to predict the mean and variance of the emotional affect measurement taken at each time j .

Table 5. Prior specification for the case study.

Parameter type	Notation	Prior
Latent mean		
Slope	$b_{1i_{Em}}$	Normal(0, 2)
Std. dev. for varying intercept	$\tau_{1i_{Em}}$	Half-normal $^{100:25i_{1/2}}$
Latent std. dev.		
Intercept	$b_{0r_{Em}}$	Normal(0, 0.25)
Slope	$b_{1r_{Em}}$	Normal(0, 0.25)
Std. dev. for varying intercept	$\tau_{r_{Em}}$	Half-normal $^{100:25i_{1/2}}$
Item parameters		
Intercepts (centered)		Student- $t(3, 0, 2.5)$
Factor loadings	k	Normal(0, 2)
Error std. dev.	s	Gamma(5, 5)

indexes the items measuring neuroticism, $m \in \{1, \dots, 5\}$ indexes the items measuring emotional affect, c represents person-specific random intercepts, and y represents the observed (censored) responses to the questionnaires. The corresponding path diagram is shown in Figure 12.

To obtain an identified model, we set $b_{0i_{Ne}} = 0$ and $\tau_{Ne} = k_{Em, 1} = 1$: We chose the variance identification constraints pragmatically by fitting each latent variable in a separate model under each possible choice of constraint and picking the one, which resulted in higher ESS for the remaining parameters, but in general one should also keep in mind the way constraints interact with priors (Graves & Merkle, 2022).

We fit the model using weakly informative priors for all parameters as described in Table 5. It was

implemented using the `brms` package and code is available online.⁴ Posterior means and credible intervals are shown in Figure 13. The key parameter of interest b_r was well-estimated with an R_{hat} of 1.00 and ESS above 1200 for bulk and tail, and had a posterior mean of 0.11 with a [0.05, 0.18] 95% credible interval; this is consistent with Mader et al. (2023), which pooled 13 studies to produce an estimate of 0.10 [0.07, 0.13]. Strictly speaking, the parameters cannot be directly compared as the free factor loadings in our model lead to items being weighed differently, but investigating measurement invariance is beyond the scope of this example (however, see Robitzsch & Lüdtke, 2023).

This case study illustrates how our approach opens the door for more truthful modeling of measurement processes. The results show that the latent parameters can be well-estimated even as part of a more complex model structure. Furthermore, the implementation cost of adding censoring and varying intercepts for this analysis was essentially negligible as these features were already available in the `brms` package. This extensibility arises from the underlying conditional likelihood formulation of the latent variable model and it should offer more flexibility for researchers who wish to build sophisticated models without the need to develop distinct, SEM-specific implementations.

⁴See the folder `case-study` at <https://github.com/bdvlm-project/gdsem-paper>

Figure 13. Posterior means and credible intervals (thick line: 50%, thin line: 95%) for model 23.

5. Discussion

In this article, we developed a Gaussian distributional SEM framework for the flexible estimation of latent variable models that include latent moderators of both latent means and latent variances. We achieved this by using a Bayesian framework, which confirms the suggestion put forward by de Kort et al. (2017) that Bayesian estimation should be a viable approach for handling latent heteroscedasticity within more complex model structures.

Our simulation study used four distinct model structures to test the reliability of the estimates obtained through the conditional likelihood approach. Although our results show that all model parameters were well-calibrated, we wish to emphasize that the SBC procedure only provides information for the parameter region covered by the generative prior, and the ones we used were relatively narrow. We did not systematically investigate wider generative priors because they produced datasets with unrealistic properties and, as a result, lead to convergence issues too often to be reliably fitted. However, we did employ weakly informative priors for the assessment of convergence and parameter recovery with favorable results. Therefore, we anticipate our framework to function well across a broad range of prior specifications. That said, it is highly recommended that users employ prior predictive checks to ensure the appropriateness of their choices in any particular analysis (Winter & Depaoli, 2023, provide a practical illustration of this technique in the specific context of SEM).

One important practical consideration that we did not address here is the impact of model misspecification on the resulting inferences. Results from de Kort et al. (2017) show that biased estimates will be obtained in situations where heteroscedasticity and nonlinearity are simultaneously present. Given that the model examined by de Kort et al. can be seen as a special case of the structures we have considered here, we expect the same caveats to carry over.

For future research, the analogy with distributional regression directly suggests the possibility of using conditional likelihood SEMs to explore non-Gaussian distributions for latent variables, including the specification of moderators on distributional parameters beyond the variance. As mentioned above, we also faced the challenge of specifying sensible generative priors while designing our simulation study. It has already been highlighted by Merkle et al. (2023) that the default approach of using non-informative priors implies data-generating processes that are incompatible with the patterns that would motivate using SEM in the first place. However, during model validation (e.g. when performing SBC), one would also like to cover as much of the parameter space as possible. Therefore, a relevant direction may be to move away from the usual approach of specifying priors on each individual model parameter and instead explore methods that use information expressed on more intuitive scales to construct the implied prior on the parameter scale (e.g. Aguilar & Burkner, 2023; Bockting et al., 2023). Another possibility that we recently became aware of is to keep non-informative

priors while simultaneously introducing imaginary data in order to produce an updated prior (for an overview of the approach see Ibrahim et al., 2015; we demonstrate an application to SEM in Fazio et al., 2024).

Finally, while we have shown that the performance of our Gaussian distributional SEMs is sufficient for practical everyday applications, there is certainly room for further optimization. For this paper we used the `brms`-generated Stan code as-is, but we are aware that applying a non-centered parametrization (Papaspiliopoulos et al., 2007) to the latent variables leads to noticeable performance gains, so it would be helpful to implement this as an option in `brms` itself. Alternatively, variational approximations can be used in place of MCMC for fast posterior estimation. Initial results for SEM estimation have been encouraging (Dang & Maestrini, 2022), but the statistical performance of these approximate methods still needs to be studied in a wider range of scenarios. Another promising set of approaches is those from the field of simulation-based inference, in particular, machine learning-based methods of posterior estimation and amortized inference (Cranmer et al., 2020; Radev et al., 2022; Zammit-Mangion et al., 2024). These techniques offer much faster inference-time results at the upfront cost of an initial training phase, but we have not found works that show their specific application to SEM estimation at this time.

Article information

Conflict of Interest Disclosures: Each author signed a form for disclosure of potential conflicts of interest. No authors reported any financial or other conflicts of interest in relation to the work described.

Ethical Principles: The authors affirm having followed professional ethical guidelines in preparing this work. These guidelines include obtaining informed consent from human participants, maintaining ethical treatment and respect for the rights of human or animal participants, and ensuring the privacy of participants and their data, such as ensuring that individual participants cannot be identified in reported results or from publicly available original or archival data.

Funding: This work was supported by Project 497785967 from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation).

Role of the Funders/Sponsors: None of the funders or sponsors of this research had any role in the design and conduct of the study; collection, management, analysis, and interpretation of data; preparation, review, or approval of the manuscript; or decision to submit the manuscript for publication.

Acknowledgments: The authors would like to thank Timo Stenz for producing the path diagrams shown in this paper and the anonymous reviewers who provided us with valuable feedback. The ideas and opinions expressed herein are those of the authors alone, and endorsement by the authors'

institution or the funding agency is not intended and should not be inferred.

References

- Aguilar, J. E., & Bürkner, P.-C. (2023). Intuitive joint priors for Bayesian linear multilevel models: The R2D2M2 prior. *Electronic Journal of Statistics*, 17(1), 1711–1767. <https://doi.org/10.1214/23-EJS2136>
- Arslan, R. C., Reitz, A. K., Driebe, J. C., Gerlach, T. M., & Penke, L. (2021). Routinely randomize potential sources of measurement reactivity to estimate and adjust for biases in subjective reports. *Psychological Methods*, 26(2), 175–185. <https://doi.org/10.1037/met0000294>
- Asparouhov, T., Hamaker, E. L., & Muthén, B. (2018). Dynamic structural equation models. *Structural Equation Modeling: A Multidisciplinary Journal*, 25(3), 359–388. <https://doi.org/10.1080/10705511.2017.1406803>
- Basu, D. (1977). On the elimination of nuisance parameters. *Journal of the American Statistical Association*, 72(358), 355–366. <https://doi.org/10.1080/01621459.1977.10481002>
- Bauer, D. J. (2017). A more general model for testing measurement invariance and differential item functioning. *Psychological Methods*, 22(3), 507–526. <https://doi.org/10.1037/met0000077>
- Bauer, D. J., & Hussong, A. M. (2009). Psychometric approaches for developing commensurate measures across independent studies: Traditional and new models. *Psychological Methods*, 14(2), 101–125. <https://doi.org/10.1037/a0015583>
- Betancourt, M. (2018). A conceptual introduction to Hamiltonian Monte Carlo. *arXiv Preprint*, <https://doi.org/10.48550/arXiv.1701.02434>
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika*, 46(4), 443–459. <https://doi.org/10.1007/BF02293801>
- Bockting, F., Radev, S. T., & Bürkner, P.-C. (2023). Simulation-based prior knowledge elicitation for parametric Bayesian models. *Sci Rep*, 14, 17330. <https://doi.org/10.48550/arXiv.2308.11672>
- Bollen, K. A. (1989). *Structural equations with latent variables*. John Wiley & Sons.
- Bürkner, P.-C. (2017). Brms: An R package for Bayesian multilevel models using Stan. *Journal of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- Bürkner, P.-C. (2018). Advanced Bayesian multilevel modeling with the R Package brms. *The R Journal*, 10(1), 395. <https://doi.org/10.32614/RJ-2018-017>
- Bürkner, P.-C. (2021). Bayesian item response modeling in R with BRMS and Stan. *Journal of Statistical Software*, 100(5), 1–54. <https://doi.org/10.18637/jss.v100.i05>
- Bürkner, P.-C., Gabry, J., Kay, M., & Vehtari, A. (2023). *Posterior: Tools for working with posterior distributions*. R package version 1.6.0, <https://mc-stan.org/posterior/>.
- Bürkner, P.-C., Scholz, M., & Radev, S. T. (2023). Some models are useful, but how do we know which ones? Towards a unified Bayesian model taxonomy. *Statistics Surveys*, 17, 216–310. <https://doi.org/10.1214/23-SS145>
- Cranmer, K., Brehmer, J., & Louppe, G. (2020). The frontier of simulation-based inference. *Proceedings of the National*

- Academy of Sciences of the United States of America, 117(48), 30055–30062. <https://doi.org/10.1073/pnas.1912789117>
- Dang, K.-D., & Maestrini, L. (2022). Fitting structural equation models via variational approximations. *Structural Equation Modeling: A Multidisciplinary Journal*, 29(6), 839–853. <https://doi.org/10.1080/10705511.2022.2053857>
- de Kort, J. M., Dolan, C. V., Lubke, G. H., & Molenaar, D. (2017). Studying the strength of prediction using indirect mixture modeling: Nonlinear latent regression with heteroskedastic residuals. *Structural Equation Modeling: A Multidisciplinary Journal*, 24(2), 301–313. <https://doi.org/10.1080/10705511.2016.1250636>
- Fahrmeir, L., Kneib, T., Lang, S., & Marx, B. D. (2021). *Regression: Models, methods and applications*. Springer. <https://doi.org/10.1007/978-3-662-63882-8>
- Fazio, L., Scholz, M., & Bürkner, P.-C. (2024, August 12). *Generative Bayesian modeling with implicit priors*.
- Gabry, J., Cesnovar, R., Johnson, A., & Bröder S. (2025). *Cmdstanr: R interface to 'CmdStan'*. R package version 0.9.0, <https://mcstan.org/cmdstanr/>.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (Eds.). (2014). *Bayesian data analysis* (3 ed.). CRC Press, Taylor & Francis.
- Gelman, A., Simpson, D., & Betancourt, M. (2017). The prior can generally only be understood in the context of the likelihood. *Entropy*, 19(10), 555. <https://doi.org/10.3390/e19100555>
- Graves, B., & Merkle, E. C. (2022). A note on identification constraints and information criteria in Bayesian latent variable models. *Behavior Research Methods*, 54(2), 795–804. <https://doi.org/10.3758/s13428-021-01649-8>
- Harvey, A. C. (1976). Estimating regression models with multiplicative heteroscedasticity. *Econometrica*, 44(3), 461–465. <https://doi.org/10.2307/1913974>
- Hessen, D. J., & Dolan, C. V. (2009). Heteroscedastic one-factor models and marginal maximum likelihood estimation. *The British Journal of Mathematical and Statistical Psychology*, 62(Pt 1), 57–77. <https://doi.org/10.1348/000711007X248884>
- Hildebrandt, A., Lüdtke, O., Robitzsch, A., Sommer, C., & Wilhelm, O. (2016). Exploring factor model parameters across continuous variables with local structural equation models. *Multivariate Behavioral Research*, 51(2–3), 257–258. <https://doi.org/10.1080/00273171.2016.1142856>
- Hoffman, M. D., & Gelman, A., et al. (2014). The no-U-turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(1), 1593–1623.
- Ibrahim, J. G., Chen, M.-H., Gwon, Y., & Chen, F. (2015). The power prior: Theory and applications. *Statistics in Medicine*, 34(28), 3724–3749. <https://doi.org/10.1002/sim.6728>
- Johnson, A. A., Ott, M. Q., & Dogucu, M. (2022). *Bayes Rules!: An introduction to applied Bayesian modeling* (1st ed.). Chapman and Hall/CRC. <https://doi.org/10.1201/9780429288340>
- Jøreskog, K. G. (1971). Simultaneous factor analysis in several populations. *Psychometrika*, 36(4), 409–426. <https://doi.org/10.1007/BF02291366>
- Kim, S., Moon, H., Modrak, M., & Sailyoja, T. (2023). *SBC: Simulation based calibration for rstan/cmdstanr models*.
- Klein, A., & Moosbrugger, H. (2000). Maximum likelihood estimation of latent interaction effects with the LMS method. *Psychometrika*, 65(4), 457–474. <https://doi.org/10.1007/BF02296338>
- Kline, R. B. (2016). *Principles and practice of structural equation modeling*. (4th ed.). The Guilford Press.
- Landau, W. M. (2021). The targets R package: A dynamic Make-like function-oriented pipeline toolkit for reproducibility and high-performance computing. *Journal of Open Source Software*, 6(57), 2959. <https://doi.org/10.21105/joss.02959>
- Lee, S.-Y. (2007). *Structural equation modeling: A Bayesian approach*. Wiley.
- Leitgeb, H., Seddig, D., Asparouhov, T., Behr, D., Davidov, E., De Roover, K., Jak, S., Meitinger, K., Menold, N., Muthen, B., Rudnev, M., Schmidt, P., & van de Schoot, R. (2023). Measurement invariance in the social sciences: Historical development, methodological challenges, state of the art, and future perspectives. *Social Science Research*, 110, 102805. <https://doi.org/10.1016/j.ssresearch.2022.102805>
- Mader, N., Arslan, R. C., Schmukle, S. C., & Rohrer, J. M. (2023). Emotional (in)stability: Neuroticism is associated with increased variability in negative emotion after all. *Proceedings of the National Academy of Sciences of the United States of America*, 120(23), e2212154120. <https://doi.org/10.1073/pnas.2212154120>
- Martin, S. R., & Rast, P. (2022). The reliability factor: Modeling individual reliability with multiple items from a single assessment. *Psychometrika*, 87(4), 1318–1342. <https://doi.org/10.1007/s11336-022-09847-9>
- Merkle, E. C., Ariyo, O., Winter, S. D., & Garnier-Villareal, M. (2023). Opaque prior distributions in Bayesian latent variable models. *Methodology*, 19(3), 228–255.
- Merkle, E. C., Fitzsimmons, E., Uanhoro, J., & Goodrich, B. (2021). Efficient Bayesian structural equation modeling in Stan. *Journal of Statistical Software*, 100(6), 1–22. <https://doi.org/10.18637/jss.v100.i06>
- Merkle, E. C., Furr, D., & Rabe-Hesketh, S. (2019). Bayesian comparison of latent variable models: Conditional versus marginal likelihoods. *Psychometrika*, 84(3), 802–829. <https://doi.org/10.1007/s11336-019-09679-0>
- Modrak, M., Moon, A. H., Kim, S., Bürkner, P., Huurre, N., Faltejskova, K., Gelman, A., & Vehtari, A. (2023). Simulation-based calibration checking for Bayesian computation: The choice of test quantities shapes sensitivity. Bayesian Analysis, Advance publication.
- Molenaar, D. (2015). Heteroscedastic latent trait models for dichotomous data. *Psychometrika*, 80(3), 625–644. <https://doi.org/10.1007/s11336-014-9406-0>
- Molenaar, D., Dolan, C. V., & van der Maas, H. L. (2011). Modeling ability differentiation in the second-order factor model. *Structural Equation Modeling: A Multidisciplinary Journal*, 18(4), 578–594. <https://doi.org/10.1080/10705511.2011.607095>
- Molenaar, D., Dolan, C. V., & Verhelst, N. D. (2010). Testing and modelling non-normality within the one-factor model. *The British Journal of Mathematical and*

- Statistical Psychology*, 63(Pt 2), 293–317. <https://doi.org/10.1348/000711009X456935>
- Molenaar, D., Van Der Sluis, S., Boomsma, D. I., & Dolan, C. V. (2012). Detecting specific genotype by environment interactions using marginal maximum likelihood estimation in the classical twin design. *Behavior Genetics*, 42(3), 483–499. <https://doi.org/10.1007/s10519-011-9522-x>
- Montoya, A. K. (2023). Selecting a within- or between-subject design for mediation: Validity, causality, and statistical power. *Multivariate Behavioral Research*, 58(3), 616–636. <https://doi.org/10.1080/00273171.2022.2077287>
- Neal, R. M. (2012). MCMC using Hamiltonian dynamics. *arXiv Preprint*. <https://doi.org/10.48550/arXiv.1206.1901>
- Nestler, S. (2020). Modelling inter-individual differences in latent within-person variation: The confirmatory factor level variability model. *The British Journal of Mathematical and Statistical Psychology*, 73(3), 452–473. <https://doi.org/10.1111/bmsp.12196>
- Neyman, J., & Scott, E. L. (1948). Consistent estimates based on partially consistent observations. *Econometrica*, 16(1), 1–32. <https://doi.org/10.2307/1914288>
- Papaspiliopoulos, O., Roberts, G. O., & Skögl, M. (2007). A general framework for the parametrization of hierarchical models. *Statistical Science*, 22(1), 59–73. <https://doi.org/10.1214/088342307000000014>
- R Core Team (2023). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing.
- Radev, S. T., Mertens, U. K., Voss, A., Ardizzone, L., & Kothe, U. (2022). BayesFlow: Learning complex stochastic models with invertible neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 33(4), 1452–1466. <https://doi.org/10.1109/TNNLS.2020.3042395>
- Robitzsch, A., & Lüdtke, O. (2023). Why full, partial, or approximate measurement invariance are not a prerequisite for meaningful and valid group comparisons. *Structural Equation Modeling: A Multidisciplinary Journal*, 30(6), 859–870. <https://doi.org/10.1080/10705511.2023.2191292>
- Rohrer, J. M., Hünemund, P., Arslan, R. C., & Elson, M. (2022). That's a lot to process! Pitfalls of popular path models. *Advances in Methods and Practices in Psychological Science*, 5(2), 251524592210958. <https://doi.org/10.1177/25152459221095827>
- Säilynoja, T., Bürkner, P.-C., & Vehtari, A. (2022). Graphical test for discrete uniformity and its applications in goodness-of-fit evaluation and multiple sample comparison. *Statistics and Computing*, 32(2), 32. <https://doi.org/10.1007/s11222-022-10090-6>
- Stan Development Team. (2023). *Stan modeling language users guide and reference manual*, 2.33.
- Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. (2020). Validating Bayesian inference algorithms with simulation-based calibration. *arXiv Preprint*. <https://doi.org/10.48550/arXiv.1804.06788>
- Thomson, G. H., & Lederunn, W. (1939). The influence of multivariate selection on the factorial analysis of ability. *British Journal of Psychology*, 29(3), 288–306. <https://doi.org/10.1111/j.2044-8295.1939.tb00919.x>
- van Erp, S. (2020). Bayesian structural equation modeling: The power of the prior. Gildeprint.
- Vehtari, A., Gelman, A., Simpson, D., Carpenter, B., & Bürkner, P.-C. (2021). Rank-normalization, folding, and localization: An improved \hat{R} for assessing convergence of MCMC (with discussion). *Bayesian Analysis*, 16(2), 667–718. <https://doi.org/10.1214/20-BA1221>
- Winter, S. D., & Depaoli, S. (2023). Illustrating the value of prior predictive checking for Bayesian structural equation modeling. *Structural Equation Modeling: A Multidisciplinary Journal*, 30(6), 1000–1021. <https://doi.org/10.1080/10705511.2022.2164286>
- Zammit-Mangion, A., Sainsbury-Dale, M., & Huser, R. (2024, October 10). *Neural methods for amortized inference*.
- Zugna, D., Popovic, M., Fasanelli, F., Heude, B., Scelo, G., & Richiardi, L. (2022). Applied causal inference methods for sequential mediators. *BMC Medical Research Methodology*, 22(1), 301. <https://doi.org/10.1186/s12874-022-01764-w>

Appendix A. Additional diagnostics

Figure A1. Convergence and parameter recovery results for fits with the generative prior.