

Chapter 6

Revisiting and Extending PLS for Ordinal Measurement and Prediction



Tamara Schamberger, Gabriele Cantaluppi, and Florian Schubert

Abstract Traditionally, partial least squares (PLS) and consistent partial least squares (PLSc) assume the indicators to be continuous. To relax this restrictive assumption, ordinal partial least squares (OrdPLS) and ordinal consistent partial least squares have been developed. They are extensions of PLS and PLSc, respectively, that are able to take into account the nature of ordinal variables—both belonging to exogenous and endogenous constructs. In the PLS context, assessing the out-of-sample predictive power of models has increasingly gained interest. In contrast to PLS and PLSc, performing out-of-sample predictions is not a straightforward process for OrdPLS and OrdPLSc because the two assume that ordinal indicators are the outcome of categorized unobserved continuous variables, i.e., they rely on polychoric and polyserial correlations. In this chapter, we present OrdPLSPredict and OrdPLScPredict to perform out-of-sample predictions with models estimated by OrdPLS and OrdPLSc. A Monte Carlo simulation demonstrates the performance of

An earlier version of this chapter was published in the following Ph.D. thesis: Schamberger T. (2022) Methodological Advances in Composite-based Structural Equation Modeling. University of Würzburg/University of Twente, <https://doi.org/10.3990/1.9789036553759>.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-37772-3_6.

T. Schamberger (✉) · F. Schubert
Faculty of Engineering Technology, Department of Design, Production and Management,
University of Twente, Drienerlolaan 5, 7522 NB Enschede, The Netherlands
e-mail: t.s.schamberger@utwente.nl; tamara.schamberger@uni-wuerzburg.de

F. Schubert
e-mail: f.schubert@utwente.nl

T. Schamberger
Faculty of Business Management and Economics, University of Würzburg, Sanderring 2, 97070
Würzburg, Germany

G. Cantaluppi
Faculty of Economics, Department of Statistical Science, Università Cattolica del Sacro Cuore,
Largo A. Gemelli, 1, 20123 Milan, Italy
e-mail: gabriele.cantaluppi@unicatt.it

our proposed approach. Finally, we provide concise guidelines using the open source R package *cSEM* to enable researchers to apply *OrdPLSPredict* and *OrdPLScPredict* using an empirical example.

6.1 Introduction

In empirical research, scholars often encounter variables that are measured on an ordinal scale. Such ordinal variables are often the outcomes of surveys that measure constructs such as organizational identification using Likert scales (e.g., 1: strongly disagree to 5: strongly agree; Hwang & Takane, 2004). Similarly, classified variables such educational level, e.g., high school or less, university student, university graduate and postgraduate or more, and grouped age, e.g., below 18 years, between 18 and 35 years and above 35 years, are ordinal variables. In contrast to metric variables, the differences between the outcomes of ordinal variables are not interpretable. Specifically, the outcomes of ordinal variables are discrete and possess a natural order (Vogt, 1993). Therefore, not taking into account the ordinal scale of variables and treating them as metric can lead to distorted results as acknowledged in the partial least squares (PLS) literature. For instance, Lohmöller (1989, p.155) recognizes that “standard procedures cannot be used for the categorical and ordinal-scaled variables”. To address this issue, ordinal partial least squares (OrdPLS, Cantaluppi, 2012; Cantaluppi & Boari, 2016) and ordinal consistent partial least squares (OrdPLSc, Schubert et al., 2018; Schubert & Cantaluppi, 2017) were developed. They are similar to PLS (Wold, 1974, 1982) and consistent partial least squares (PLSc, Dijkstra & Henseler, 2015a, 2015b) and can deal with models containing exogenous and endogenous constructs associated with ordinal indicators. However, the originally used Pearson correlations in PLS and PLSc are replaced by polychoric and polyserial correlations which assume that an ordinal indicator is the outcome of a categorized unobserved standard normally distributed random variable (Poon & Lee, 1987). Consequently, researchers who want to apply PLS and PLSc and strive for consistent estimates in the case of ordinal indicators are advised to apply OrdPLS and OrdPLSc, respectively.

Over the last years, the causal-predictive nature of PLS was emphasized in which the assessment of a model’s out-of-sample predictive performance plays a crucial role (Chin et al., 2020; Sarstedt et al., 2023). Although Cantaluppi and Schubert (2019) proposed an approach that is based on OrdPLS and OrdPLSc to perform out-of-sample predictions and thus takes into account the nature of ordinal indicators, this approach is limited to situations where all indicators are on an ordinal scale. However, empirical studies applying PLS often deal with both metric and ordinal indicators. For instance, IT integration capability was modeled as a composite composed of ordinal variables each measured on a 5-point Likert scale, while the control variable firm performance was measured by a metric variable, namely the natural logarithm of the number of employees (Braojos et al., 2020). The necessity to make predictions of ordinal variables based on a mix of metric and ordinal variables is also observable

in other fields such as credit scoring. In this setting, final predictions need to be formulated on ordinal scales in order to take decisions, e.g., grant or do not grant a credit or assign a rating to a customer.

To address this issue, we propose `OrdPLScpredict` and `OrdPLSpredict` which are extensions of the approaches of Cantaluppi and Schuberth (2019) and Shmueli et al. (2016) to perform out-of-sample predictions based on models estimated by `OrdPLS` and `OrdPLSc` containing both continuous and ordinal indicators. The remainder of the chapter is organized as follows: Sect. 6.2 presents `OrdPLS` and `OrdPLSc`. Section 6.3 gives an overview of performing out-of-sample predictions using `PLSpredict` and `PLScpredict`, i.e., the approaches originally proposed to perform out-of-sample prediction based on models estimated by PLS and PLSc, respectively (Shmueli et al., 2016). In Sect. 6.4, we present `OrdPLSpredict` and `OrdPLScpredict`, our two proposed approaches to perform out-of-sample predictions using models estimated by `OrdPLS` and `OrdPLSc`, respectively. In Sect. 6.5, we conduct a Monte Carlo simulation to evaluate the performance of our two proposed approaches. Section 6.6 provides guidelines for the two approaches and shows how they can be applied in the open source R package `cSEM` (Rademaker & Schuberth, 2020) using an illustrative example. Our chapter closes with a discussion given in Sect. 6.7.

6.2 Ordinal (Consistent) Partial Least Squares Path Modeling

Wold (1966) originally developed PLS as an approach for principal component analysis and (generalized) canonical correlation analysis, which at the time were still known as nonlinear iterative least squares and nonlinear iterative partial least squares, respectively (Tenenhaus et al., 2005). A few years later, Wold proposed PLS as a computational efficient estimation method for structural models containing latent variables (Wold, 1974; 1982). In this case, weights are determined by the PLS algorithm to form proxies and subsequently these proxies are used to estimate the relationships between the latent variables. As various researchers emphasized, PLS estimates for this type of model are only consistent at large; i.e., only as both the number of observations and the number of indicators go to infinity, will PLS estimates converge in probability to the respective population parameters (e.g., Hui & Wold, 1982, Dijkstra, 1985). However, recently, various studies have shown that PLS produces consistent estimates for models containing interrelated emergent variables (Dijkstra 2017; Cho & Choi 2020; Henseler 2021; Schuberth 2021).¹ For an elaboration about emergent variables and their potential use, the interested reader is referred to Yu et al. (2021).

¹ In line with recent literature (e.g., Benitez et al., 2020; Yu et al., 2021; Schamberger et al., 2023), we use the term ‘emergent variable’ to emphasize that the variable is not only a composite, i.e., a weighted linear combination of variables, but also a composite that conveys all the information between its indicators and other variables in the model.

In its most modern appearance known as PLSc, it produces consistent parameter estimates for structural models containing latent and emergent variables (Dijkstra & Henseler, 2015b). Similar to PLS, PLSc relies on the PLS algorithm to determine the weights to build proxies for the constructs. In cases that constructs are modeled as latent variables, it applies a correction for attenuation to correlations. In this way, it is ensured that the construct correlation matrix is consistently estimated, and thus, consistent path coefficient estimates can be obtained. Moreover, in contrast to PLS which relies on ordinary least squares (OLS) to estimate the model parameters, PLSc applies two-stage least squares (2SLS) in the case of non-recursive structural models (Dijkstra & Henseler, 2015a). Finally, a recent development allows PLSc to deal with correlated random measurement errors within a block of indicators measuring a latent variable (Rademaker et al., 2019). For an overview on latent variable models that can be estimated by PLSc, we refer to the study of Schuberth et al. (2023b).

PLS including PLSc assumes that the observed variables are measured on a metric scale. However, in empirical research observed variables are often measured on an ordinal scale. For instance, a model about customer satisfaction (Sarstedt et al., 2011), and a model about young consumer's adoption intentions (Miltgen et al., 2016) rely on observed variables that are measured on an ordinal scale. To overcome the limitation of PLS considering the scale of the observed variables, various modifications of PLS have been developed to cope with non-metric variables such as partial maximum-likelihood partial least squares (Jakobowicz & Derquenne, 2007) and non-metric partial least squares (Russolillo, 2012). A further approach that was developed to deal with ordinal indicators in a classic psychometric way is OrdPLS (Cantaluppi, 2012). OrdPLS is similar to PLS, but applies polychoric and polyserial correlations instead of Pearson correlations as input for the PLS algorithm to take the nature of ordinal indicators into account. Consequently, the original PLS algorithm remains untouched. In the same way as PLS was extended by PLSc, OrdPLS was extended by OrdPLSc to consistently estimate structural models containing latent variables and ordinal indicators (Schuberth & Cantaluppi, 2017; Schuberth et al., 2018). Figure 6.1 illustrates the four steps of OrdPLSc: (i) calculating the polychoric/polyserial correlations, (ii) performing the PLS algorithm, (iii) correcting for attenuation if some constructs are modeled as latent variables, and (iv) estimating the path coefficients by OLS and 2SLS, respectively. Obviously, OrdPLSc only differs from traditional PLSc in step (i). The three other steps are the same for PLSc and OrdPLSc. In the following, we elaborate each of the four steps.²

6.2.1 *Calculating Polychoric/Polyserial Correlations*

Following Pearson's idea of a polytomous variable, we assume an ordinal indicator x to be the result of a categorized unobservable standard normally distributed random

² Note that the following subsections contain large parts adapted from Schuberth et al. (2018). Which is released under a Creative Commons Attribution 4.0 International License (CC BY 4.0).

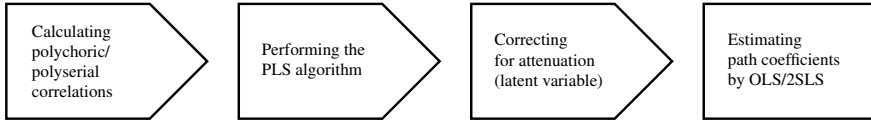


Fig. 6.1 Ordinal consistent partial least squares (adopted from Schuberth et al. (2018) which is released under a Creative Commons Attribution 4.0 International License (CC BY 4.0))

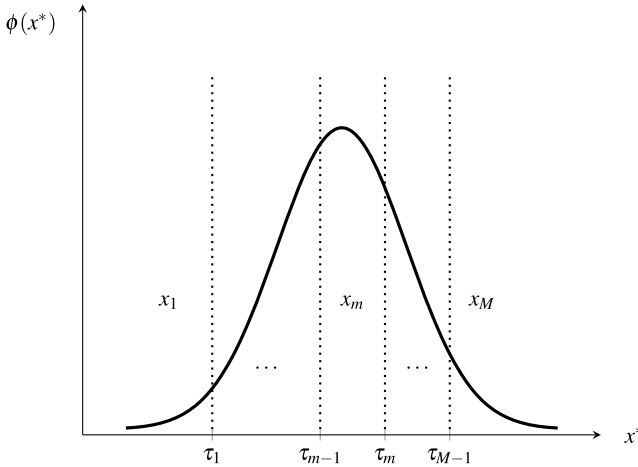


Fig. 6.2 Pearson’s idea of an ordinal variable (taken from Schuberth et al. (2018) which is released under a Creative Commons Attribution 4.0 International License (CC BY 4.0))

variable x^* (Pearson, 1900, 1913)

$$x = x_m \text{ if } \tau_{m-1} \leq x^* < \tau_m \quad m = 1, \dots, M \tag{6.1}$$

where the threshold parameters τ_0, \dots, τ_M determine the observed categories. The first and last thresholds are fixed: $\tau_0 = -\infty$ and $\tau_M = \infty$. Moreover, we assume the thresholds to be strictly increasing: $\tau_0 < \tau_1 < \dots < \tau_M$.³ Figure 6.2 depicts the idea of an underlying continuous variable, i.e., for indicator x , category x_m is observed if the realization of the underlying continuous variable x^* falls between the thresholds τ_{m-1} and τ_m .

Since we assume an ordinal variable to be determined by an underlying continuous variable, it is more appropriate to consider the correlation between these underlying continuous variables for evaluating the linear relationship of interest. This is achieved by using the polychoric correlation (Drasgow, 1986). In cases where the correlation between an ordinal variable and a continuous variable is calculated the polyserial correlation can be used (Lee & Poon, 1986) which considers the correlation between the

³ In empirical work two consecutive threshold parameters can be equal, $\tau_{m-1} = \tau_m$, if the corresponding category x_m is not observed.

continuous variable and the ordinal variable's underlying continuous variable. Various estimators have been proposed to obtain these correlation coefficients (Drasgow, 1986). In the following, we apply the two-step estimator which is computationally efficient.

As input for OrdPLS and OrdPLSc, the indicators' sample correlation matrix S is required and thus needs to be calculated. Since OrdPLS and OrdPLSc allow for both ordinal and continuous indicators, the sample correlation matrix S can comprise polychoric, polyserial and Pearson correlation coefficients. In specific, the sample correlation between two ordinal indicators equals their estimated polychoric correlation, the sample correlation between an ordinal and a continuous indicator equals their estimated polyserial correlation and the sample correlation between two continuous indicators equals their estimated Pearson correlation.

6.2.2 Performing the PLS Algorithm

The second step of OrdPLS and OrdPLSc involves applying the PLS algorithm to the sample correlation matrix S calculated in the previous step. The PLS algorithm remains untouched by OrdPLS and OrdPLSc and thus is the same as for their traditional counterparts. For simplicity, the K_j indicators belonging to one construct η_j , i.e., a latent variable or an emergent variable, are grouped to form the block j with $j = 1, \dots, J$ and where $\sum_{j=1}^J K_j = K$, i.e., each indicator belongs exactly to one block.

The PLS algorithm is an iterative algorithm which starts with initial arbitrary weights $\hat{\mathbf{w}}_j^{(0)}$ ($K_j \times 1$). The initial weights are chosen in such a way that they satisfy the following condition: $\hat{\mathbf{w}}_j^{(0)'} S_{jj} \hat{\mathbf{w}}_j^{(0)} = 1$ for each block j where the $(K_j \times K_j)$ matrix S_{jj} contains the sample correlations of the indicators of block j . This condition holds for all weights in each iteration i and can be achieved by using a scaling factor $(\hat{\mathbf{w}}_j^{(i)'} S_{jj} \hat{\mathbf{w}}_j^{(i)})^{-\frac{1}{2}}$ for the weights $\hat{\mathbf{w}}_j^{(i)}$.

The PLS algorithm aims to determine weights to build proxies for the J constructs. This can be done in three ways, identified as *Mode A*, *Mode B*, and *Mode C*. In the case of *Mode A*, the weights, also known as correlation weights, are determined as follows:

$$\hat{\mathbf{w}}_j^{(i+1)} \propto \sum_{l=1}^J S_{jl} \hat{\mathbf{w}}_l^{(i)} e_{jl}^{(i)} \quad \text{with} \quad \hat{\mathbf{w}}_j^{(i+1)'} S_{jj} \hat{\mathbf{w}}_j^{(i+1)} = 1. \quad (6.2)$$

In the case of *Mode B*, the weights, also known as regression weights, are calculated as follows:

$$\hat{\mathbf{w}}_j^{(i+1)} \propto S_{jj}^{-1} \sum_{l=1}^J S_{jl} \hat{\mathbf{w}}_l^{(i)} e_{jl}^{(i)} \quad \text{with} \quad \hat{\mathbf{w}}_j^{(i+1)'} S_{jj} \hat{\mathbf{w}}_j^{(i+1)} = 1. \quad (6.3)$$

Mode C, also known as *MIMIC mode*, is a mixture of mode A and B, which we do not consider here. The inner weights e_{jl} can be obtained in three different ways, following the *centroid* (Wold, 1982), *factorial* (Lohmöller, 1989), and *path* weighting scheme. All inner weighting schemes produce essentially the same results (Noonan & Wold, 1982), hence, we consider the path weighting scheme here.⁴ For the path weighting scheme, the inner weight jl is chosen as follows:

$$e_{jl}^{(i)} = \begin{cases} \hat{\mathbf{w}}_j^{(i)'} \mathbf{S}_{jl} \hat{\mathbf{w}}_l^{(i)} & \text{if } \eta_l \text{ is a consequence of } \eta_j \\ \hat{\beta}_l & \text{if } \eta_l \text{ is an antecedent of } \eta_j \\ 0 & \text{otherwise} \end{cases} \quad (6.4)$$

As Eq. (6.4) shows, the inner weight e_{jl} equals the covariance between the proxies of the constructs η_j and η_l if construct η_l is a consequence of the construct η_j . In contrast, if the construct η_l is an antecedent of the construct η_j , the inner weight e_{jl} is equal to the regression coefficient $\hat{\beta}_l$ of a multiple regression of the construct η_j on its antecedents. Otherwise, if the two constructs are not connected via the structural model, the inner weight is set to 0.

Since the PLS algorithm has no single criterion to be optimized, the new weights $\hat{\mathbf{w}}_j^{(i+1)}$ are checked for significant changes compared with the weights $\hat{\mathbf{w}}_j^{(i)}$ in the previous iteration step. When the change in weights exceeds a certain limit, the algorithm starts again. Otherwise, the final weights $\hat{\mathbf{w}}_j$ equal the stable weights determined in the last iteration. Finally, the standardized loading estimates, which are in PLS equal to the estimated correlations between a proxy and its indicators, are calculated as:

$$\hat{\boldsymbol{\lambda}}_j = \mathbf{S}_{jj} \hat{\mathbf{w}}_j \quad (6.5)$$

If we apply OrdPLS, the standardized loading estimates are calculated in the same way. The only difference is that the polychoric/polyserial correlation matrix is taken into account and therefore the calculation considers correlations with the underlying continuous variables that correspond to ordinal indicators of a latent variable. Note, as for PLS, the loading estimates of OrdPLS are not consistent for latent variable models.

6.2.3 Correcting for Attenuation if Constructs Are Modeled as Latent Variables

PLS creates composites as proxies for constructs. Consequently, its estimates are not consistent if the constructs are modeled as latent variables. To overcome this issue,

⁴ Note that the choice of inner weighting scheme can substantially affect the estimates in the case of models containing second-order constructs (Becker et al., 2012; Schuberth et al., 2020). For more details on the other inner weighting schemes, see Tenenhaus et al. (2005).

Dijkstra and Henseler (2015a, 2015b) proposed PLSc which applies a correction to obtain consistent parameter estimates. OrdPLSc applies the same correction to obtain consistent estimates for models containing latent variables. Consequently, this step is the same for PLSc and OrdPLSc.

The correction exploits the linearity between standardized population factor loadings and the population weights, $\lambda_j = c_j \mathbf{w}_j$ and requires that each latent variable be measured by at least two indicators. The estimated correction factor for block j satisfies the following condition

$$\text{plim } \hat{c}_j = \sqrt{\lambda_j' \Sigma_{jj} \lambda_j}, \quad (6.6)$$

where λ_j is a column vector of length K_j containing the population loadings of latent variable η_j and Σ_{jj} is the $(K_j \times K_j)$ population correlation matrix of the indicators of block j .⁵ The correction factor \hat{c}_j can be obtained by

$$\hat{c}_j^2 = \frac{\hat{\mathbf{w}}_j' (\mathbf{S}_{jj} - \text{diag}(\mathbf{S}_{jj})) \hat{\mathbf{w}}_j}{\hat{\mathbf{w}}_j' (\hat{\mathbf{w}}_j \hat{\mathbf{w}}_j' - \text{diag}(\hat{\mathbf{w}}_j \hat{\mathbf{w}}_j')) \hat{\mathbf{w}}_j}. \quad (6.7)$$

It is chosen in such a way that the Euclidean distance between

$$\mathbf{S}_{jj} - \text{diag}(\mathbf{S}_{jj}) \quad \text{and} \quad (c_j \hat{\mathbf{w}}_j)(c_j \hat{\mathbf{w}}_j)' - \text{diag}((c_j \hat{\mathbf{w}}_j)(c_j \hat{\mathbf{w}}_j')) \quad (6.8)$$

is minimized (Dijkstra & Henseler, 2015a). For other ways to obtain correction factors, the interested reader is referred to Dijkstra (2013). Finally, the standardized factor loadings of block j can be consistently estimated as

$$\hat{\lambda}_j = \hat{c}_j \hat{\mathbf{w}}_j. \quad (6.9)$$

6.2.4 Estimating Path Coefficients by OLS/2SLS

In the last step, we estimate the path coefficients based on the proxies' correlation matrix, i.e., $\hat{\mathbf{W}}' \mathbf{S} \hat{\mathbf{W}}$, where the matrix $\hat{\mathbf{W}}$ of dimension $K \times J$ contains all the weight estimates. In PLS and OrdPLS, this matrix is directly applied to estimate the parameters of the structural model by OLS. In contrast, if constructs are modeled as latent variables, PLSc and OrdPLSc apply a correction for attenuation to the proxies' correlation matrix before calculating the path coefficients. The correlation between the two latent variables η_j and η_l where $j \neq l$ can be consistently estimated by:

$$\widehat{\text{cor}}(\eta_j, \eta_l) = \frac{\hat{\mathbf{w}}_j' \mathbf{S}_{jl} \hat{\mathbf{w}}_l}{\hat{c}_j^2 \hat{\mathbf{w}}_j' \hat{\mathbf{w}}_j \hat{c}_l^2 \hat{\mathbf{w}}_l' \hat{\mathbf{w}}_l} \quad (6.10)$$

⁵ Here we do not consider the use of *Mode B* for latent variables. For a consistent version of PLS using *Mode B*, the interested reader is referred to Dijkstra (2011).

Similarly, if construct η_j is modeled as a latent variable and construct η_l as an emergent variable, the consistently estimated correlation is obtained by

$$\widehat{\text{cor}}(\eta_j, \eta_l) = \frac{\hat{\mathbf{w}}_j' \mathbf{S}_{jl} \hat{\mathbf{w}}_l}{\hat{c}_j^2 \hat{\mathbf{w}}_j' \hat{\mathbf{w}}_j}. \quad (6.11)$$

In the case of both constructs being modeled as emergent variables, no correction of the correlation is required because we assume that the correlation between two emergent variables is not affected by attenuation. Finally, in OrdPLSc the path coefficients are estimated by OLS or 2SLS depending on the structure of the underlying structural model.

6.3 Model-Based Predictions Using PLS and PLSc (PLSpredict and PLScpredict)

In the PLS context, out-of-sample predictions have increasingly gained attention (Evermann & Tate, 2014; Carrión et al., 2016; Shmueli et al., 2016, 2019; Sarstedt & Danks, 2022). To perform such out-of-sample predictions, a procedure called PLSpredict was introduced (Shmueli et al., 2016). In PLSpredict, values of variables are predicted based on a model estimated by PLS. In cases where the model parameters are estimated by PLSc, we label the procedure PLScpredict. In the following exposition we present the steps of PLSpredict and PLScpredict.

We begin by splitting a sample into two datasets, i.e., the train dataset $\mathbf{X}_{\text{train}}$ and the test dataset \mathbf{X}_{test} . The train dataset contains observations for all indicators and is used to estimate the model parameters by PLS or PLSc, i.e., the weights $\hat{\mathbf{w}}_j$, the loadings $\hat{\boldsymbol{\lambda}}_j$, and the path coefficients of the exogenous and endogenous constructs, which are captured in the matrices $\hat{\mathbf{\Gamma}}$ and $\hat{\mathbf{B}}$, respectively. Subsequently, out-of-sample predictions can be performed based on the estimated model and the observations given in the test dataset. The test dataset comprises N observations for at least all indicators connected to exogenous constructs, i.e., constructs that are not explained by other constructs in the structural model. Importantly, the observations of the test dataset are not used during the model estimation.

In the context of PLSpredict, we can distinguish different types of predictions (Shmueli et al., 2016; Lohmöller, 1989): (i) *valid predictions* in which predictions for scores of exogenous constructs are obtained by observations of their associated indicators, (ii) *structural predictions* in which predictions for scores of endogenous constructs are obtained by exogenous construct scores, (iii) *communal predictions* in which predictions for values of indicators associated with endogenous constructs are obtained by scores of their associated constructs, (iv) *redundant predictions* in which predictions for values of the indicators associated with endogenous constructs are obtained by exogenous construct scores and the estimated structural model, (v) *latent predictions* in which predictions for scores of endogenous constructs are obtained by

observations of the indicators associated with exogenous constructs and the estimated structural model, and (vi) *operative predictions* in which predictions for values of indicators associated with endogenous constructs are obtained by observations for the indicators associated with exogenous constructs and the estimated structural model.

Obviously, operative predictions are the most general case in that they involve all the steps of the other types of predictions. Additionally, predictions can only be evaluated if they are performed on item level, i.e., if values of the indicators are predicted. Against this background, we will now focus on operative predictions. Note that other types of predictions can be obtained by starting or stopping the approach we describe below at a later or earlier stage.

To obtain operative predictions, valid predictions have to be performed first. To do this, we standardize the N observations of the test dataset \mathbf{X}_{test} for the indicators associated with the exogenous constructs using the corresponding moments estimated on the basis of the train dataset (Shmueli et al., 2016). Subsequently, for all J_{exo} exogenous constructs, we predict scores as the weighted sum of their associated indicators using the observations from the test dataset. Consequently, the predicted scores of the exogenous constructs are obtained as follows:

$$\hat{\eta}_{j,\text{exo}} = \mathbf{X}_{j,\text{test}} \hat{\mathbf{w}}_j \quad j = 1, \dots, J_{\text{exo}} \quad (6.12)$$

In a next step, we use the predicted scores of the exogenous constructs to predict the scores of the J_{end} endogenous constructs in accordance with the structural model, i.e., we perform structural predictions:

$$\hat{\eta}_{\text{end}} = \hat{\eta}_{\text{exo}} \hat{\mathbf{\Gamma}}' (\mathbf{I} - \hat{\mathbf{B}}')^{-1}, \quad (6.13)$$

where $\hat{\eta}_{\text{end}}$ is a matrix of dimension $N \times J_{\text{end}}$ that contains the predictions for the scores of the endogenous constructs in its columns.

Finally, in the last step, we use the scores of the endogenous constructs to predict values of the indicators connected to endogenous constructs, i.e., we perform communal predictions:

$$\hat{\mathbf{X}}_{\text{end}} = \hat{\eta}_{\text{end}} \hat{\mathbf{\Lambda}}'_{\text{end}} \quad (6.14)$$

where the matrix $\hat{\mathbf{\Lambda}}_{\text{end}}$ contains the estimated loadings of the indicators connected to endogenous constructs in its columns. To obtain the final predictions for continuous indicators, the values in $\hat{\mathbf{X}}_{\text{end}}$ are brought back to their original scale using the mean and standard deviation of the train dataset (see Shmueli et al., 2016). In cases where ordinal indicators are associated with endogenous constructs, Cantaluppi and Schubert (2019) proposed rounding the predicted values to an integer. Thereby, we obtain predictions that are in line with the domain of the ordinal indicators. However, the scale of the ordinal indicators was not taken into account during parameter estimation.

To evaluate the model's predictive power, the test dataset must contain observations for all indicators. In such a case, the observed values of the indicators can be compared to their predicted counterparts (Shmueli, 2010). As predictive performance measures, the mean absolute error (MAE), and the root mean squared error (RMSE) can be used to evaluate the predictive power of the model (Evermann & Tate, 2014). The MAE is the average absolute deviation of the predicted value of an indicator from its observed counterpart, $\frac{1}{N} \sum_{i=1}^N |\hat{x}_i - x_i|$, where N is the sample size of the test dataset. Similarly, the RMSE is the square root of the average squared deviation of the predicted value from its observed counterpart, $\sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{x}_i - x_i)^2}$.

6.4 Model-Based Predictions Using OrdPLS and OrdPLSc

In this section, we present an approach to perform predictions based on a model estimated by OrdPLS and OrdPLSc, which we label OrdPLSpredict and OrdPLScpredict, respectively. These approaches allow for both ordinal and continuous indicators. In fact, our approaches generalize the idea of Cantaluppi and Schubert (2019) to perform predictions based on a model estimated by OrdPLS or OrdPLSc in cases where all indicators are ordinal.

6.4.1 Relationship Between the Ordinal Indicators and Their Underlying Latent Variables in the Test Dataset

We rely on the idea presented in Sect. 6.2.1 that an ordinal indicator x is the outcome of a polytomized standard normally distributed unobservable random variable x^* , see Eq. (6.1). In cases with more than one ordinal indicator, we assume that the ordinal indicators \mathbf{x} are the outcome of categorized underlying multivariate standard normally distributed latent random variables \mathbf{x}^* . Consequently, the observations of the ordinal indicators \mathbf{x}_j belonging to construct j are the outcome of columnwise transformations (as expressed by Eq. (6.1)) of the observations of the underlying multivariate normally distributed random variables that are stacked in the matrix \mathbf{X}_j^* , expressed as:

$$\mathbf{X}_j^* \rightarrow \mathbf{X}_j \tag{6.15}$$

As shown in Sect. 6.2, OrdPLS and OrdPLSc can deal with both ordinal and continuous indicators. Note that the transformation is only performed for the observations of the ordinal indicators and not for those of continuous indicators.

As in PLSpredict, the first step is to estimate the model parameters. In the context of OrdPLSpredict and OrdPLScpredict this is done by OrdPLS and OrdPLSc, respectively, based on the train dataset which contains observations for at least one

ordinal indicator. Otherwise, if the train dataset contains no ordinal indicator, there is no need to apply OrdPLS or OrdPLSc.

Next, the estimated model and the observations of the test dataset are used to perform out-of-sample predictions. For this purpose, the test dataset must at least contain observations for the indicators associated with the exogenous constructs which are stored in the matrix $\mathbf{X}_{\text{test, exo}}$.⁶ We assume the observations of ordinal indicators of the test dataset to be the columnwise transformations of a multivariate *truncated* normally distributed dataset, as stated by Eq. (6.1):

$$\mathbf{X}_{\text{test, exo}}^{\text{Trunc,*}} \rightarrow \mathbf{X}_{\text{test, exo}} \quad (6.16)$$

The observations from the multivariate truncated normal distribution $\mathbf{X}_{\text{test, exo}}^{\text{Trunc,*}}$ are standardized and have the same correlation matrix as the polychoric correlation between the indicators connected to the train data's exogenous constructs. For each subject in the test dataset, given the expressed categories, the domain of $\mathbf{X}_{\text{test, exo}}^{\text{Trunc,*}}$ is defined by the corresponding pairs of thresholds in the set of thresholds (τ_{j-1}, τ_j) . These are obtained from the polychoric correlation matrix used for model parameter estimation, i.e., the one based on the train dataset $\mathbf{X}_{\text{train}}$. If the test dataset contains additional observations for the indicators associated with endogenous constructs, the model's predictive performance can be evaluated by comparing the indicators' observed values to their predicted counterparts.

6.4.2 *OrdPLSpredict and OrdPLScpredict*

In the following explication, we present the steps OrdPLSpredict and OrdPLScpredict take to perform out-of-sample predictions. Similar to PLSpredict and PLScpredict, the only difference between OrdPLSpredict and OrdPLScpredict is that the former uses OrdPLS estimates, while the latter employs OrdPLSc estimates.

1. Standardize the test dataset $\mathbf{X}_{\text{test, exo}}$ using the means and standard deviations of the train dataset. Note that only the continuous indicators are standardized, i.e., the ordinal indicators comprised in the test dataset remain untouched.
2. Predict the scores of the exogenous constructs, i.e., valid predictions. In PLSpredict, scores of construct j are obtained as linear combinations of the observed indicators x_j and the corresponding weight estimates, regardless of the indicators' measurement scale. In contrast, in OrdPLSpredict and OrdPLScpredict the nature of ordinal indicators is explicitly taken into account. Since the number of ordinal indicators associated with exogenous constructs can differ, three cases have to be distinguished, namely ones in which (i) all indicators are continuous, (ii) all indicators are ordinal, and (iii) there is a mixture of continuous and ordinal indicators.

Considering cases in which all indicators associated with exogenous constructs

⁶ In cases where only values of a subset of the indicators associated with endogenous constructs are predicted, a subset of the indicators associated with the exogenous constructs might be sufficient.

are continuous, construct scores are obtained as in PLSpredict:

$$\hat{\eta}_{j,\text{exo}} = \mathbf{X}_{j,\text{test,exo}} \hat{\mathbf{w}}_j, \quad j = 1, \dots, J_{\text{exo}} \quad (6.17)$$

where $\hat{\mathbf{w}}_j$ are the weight estimates obtained by OrdPLS/OrdPLSc based on the train data.

Considering cases in which all indicators associated with exogenous constructs are ordinal, the unknown values of the unobservable variables underlying these indicators (see Eq. (6.1)) need to be aggregated (Cantaluppi and Schubert, 2019). Specifically, the exogenous constructs' scores can be calculated as linear combinations of multivariate truncated normally distributed random variables $\mathbf{x}_{j,\text{test,exo}}^{\text{Trunc},*}$ which are continuous. They also have the domain (τ_{j-1}, τ_j) defined by the threshold parameters of the polychoric correlations based on the train dataset $\mathbf{X}_{\text{train}}$, conditional on categories that characterize the test dataset regarding ordinal indicators. Consequently, we obtain the construct scores as follows:

$$\hat{\eta}_{j,\text{exo}} = \mathbf{X}_{j,\text{test,exo}}^{\text{Trunc},*} \hat{\mathbf{w}}_j \quad j = 1, \dots, J_{\text{exo}} \quad (6.18)$$

As Eq. (6.18) shows, the distribution of the construct scores is a linear combination of multivariate truncated normally distributed random variables with OrdPLS/OrdPLSc weight estimates based on the train data. The distribution of the construct scores has no simple form but can be approximated by simulation. To simulate this distribution for each subject, we generate $n_{\text{pred}} = 100$ drawings from a multivariate truncated normal distribution with a variance-covariance matrix that equals the polychoric correlation matrix of the train dataset and truncation limits that equal the threshold parameter estimates of this polychoric correlation matrix. As a consequence, we obtain n_{pred} draws in total for the unobservable variables underlying the ordinal indicators associated with exogenous constructs $\mathbf{X}_{j,\text{test,exo}}^{\text{Trunc},*,p}$ for $p = 1, \dots, n_{\text{pred}}$, and thus, n_{pred} sets of predicted scores for each exogenous construct:

$$\hat{\eta}_{j,\text{exo}}^p = \mathbf{X}_{j,\text{test,exo}}^{\text{Trunc},*,p} \hat{\mathbf{w}}_j \quad j = 1, \dots, J_{\text{exo}}, \quad p = 1, \dots, n_{\text{pred}} \quad (6.19)$$

Considering a case in which there is a mixture of continuous and ordinal indicators associated with exogenous constructs, we generate n_{pred} drawings from a multivariate truncated normal distribution for both the ordinal and the continuous indicators to obtain construct scores. The variance-covariance matrix of the multivariate truncated normal distribution equals the estimated correlation matrix of the indicators based on the train dataset, which can contain polychoric, polyserial and Pearson correlations. We take the continuous indicators into account during the simulation to preserve the correlation structure. However, their generated values in $\mathbf{X}_{j,\text{test,exo}}^{\text{Trunc},*,p}$ are replaced by the corresponding observations from the test data. For the ordinal indicators, the truncation limits are appropriately chosen, conditional on categories that characterize the test data set by using the

threshold estimates obtained by the polychoric/polyserial correlations based on the train data. In contrast, for the continuous indicators we use arbitrary lower and upper truncation limits, e.g., -10 and 10 . Consequently, we obtain n_{pred} datasets for the indicators connected to exogenous constructs where the observations of the continuous indicators equal the observations from the test data, while for the ordinal indicators we use the generated dataset of the multivariate truncated normal distribution. Based on the resulting samples, we calculate the n_{pred} scores for each exogenous construct as follows:

$$\hat{\eta}_{j,\text{exo}}^p = \mathbf{X}_{j,\text{test,exo}}^{\text{Trunc},*,p} \hat{\mathbf{w}}_j \quad j = 1, \dots, J_{\text{exo}}, \quad p = 1, \dots, n_{\text{pred}} \quad (6.20)$$

3. Predict the endogenous constructs' scores using the exogenous constructs' scores in accordance with the structural model, i.e., structural predictions. Using the n_{pred} predicted scores of the exogenous constructs, n_{pred} scores for the endogenous constructs can be predicted via the structural model:

$$\hat{\eta}_{\text{end}}^p = \hat{\eta}_{\text{exo}}^p \hat{\mathbf{\Gamma}}' (\mathbf{I} - \hat{\mathbf{B}}')^{-1} \quad p = 1, \dots, n_{\text{pred}} \quad (6.21)$$

For the case in which only continuous indicators are connected to exogenous constructs, the matrix with the predicted construct scores $\hat{\eta}_{\text{exo}}^p$ is replaced by $\hat{\eta}_{\text{exo}}$ from Eq. (6.17). Consequently, we do not obtain n_{pred} matrices containing predicted scores for the endogenous constructs, but only one matrix $\hat{\eta}_{\text{end}}$.

4. Predict the values of the indicators belonging to the endogenous constructs, i.e., communal predictions. Here, two cases need to be distinguished, namely, ones in which (i) an indicator belonging to an endogenous construct is continuous, and (ii) an indicator belonging to an endogenous construct is ordinal.

If the indicator \mathbf{x}_k belonging to the j -th endogenous construct is continuous, first n_{pred} predictions are obtained by multiplying the construct scores with the estimated loading from the train dataset:

$$\hat{\mathbf{x}}_{k,\text{end}}^p = \hat{\eta}_{j,\text{end}}^p \hat{\lambda}_{j,k,\text{end}} \quad j = 1, \dots, J_{\text{end}} \quad k = 1, \dots, K_j \quad p = 1, \dots, n_{\text{pred}} \quad (6.22)$$

In contrast, if the indicator associated with an endogenous construct is ordinal, predictions for the continuous unobservable variables underlying the ordinal indicator have to be obtained first. As we have predicted n_{pred} scores for an endogenous construct, we also obtain n_{pred} predictions for the unobservable variable underlying the ordinal indicator. This we do by multiplying the endogenous construct's scores with the estimated loading corresponding to the k -th indicator \mathbf{x}_k of the j -th endogenous construct η_j :

$$\hat{\mathbf{x}}_{k,\text{end}}^{*,p} = \hat{\eta}_{j,\text{end}}^p \hat{\lambda}_{j,k,\text{end}} \quad j = 1, \dots, J_{\text{end}} \quad k = 1, \dots, K_j \quad p = 1, \dots, n_{\text{pred}} \quad (6.23)$$

Obviously, the only difference between the procedure for continuous and ordinal variables is that for the latter the values of the unobservable variable underlying the ordinal variable are predicted.

Finally, to obtain one prediction for each observation of the test dataset, the location of the distribution of the n_{pred} predictions for the indicators of the endogenous constructs need to be obtained. For this purpose, Cantaluppi and Boari (2016) proposed the *mean*, the *median*, and the *mode* approach. In the case of the mean approach, the i -th value of a continuous indicator is predicted as the mean of the n_{pred} draws, expressed as:

$$\hat{x}_{k,i,\text{end}}^* = \frac{1}{n_{\text{pred}}} \sum_{p=1}^{n_{\text{pred}}} \hat{x}_{k,i,\text{end}}^{*,p} \quad i = 1, \dots, N \quad (6.24)$$

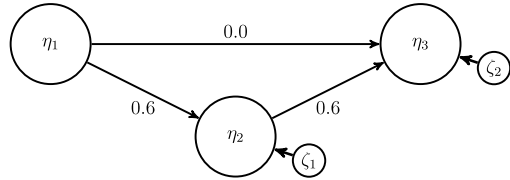
The median approach works similar to the mean approach; however, instead of using the mean to determine the location of the distribution of the n_{pred} predictions, the median is used. While for continuous indicators the final predictions equal these location values, for ordinal variables, the location values are transformed into ordinal values according to Eq. (6.1) using the estimated thresholds based on the train data.

As a third approach to summarizing the n_{pred} predictions, we can use the mode approach. It uses the maximum of the predicted unobservable variable's empirical density on the intervals defined by the thresholds. Consequently, this approach cannot be used for continuous indicators.

Finally, the continuous indicators' predicted values are brought back to their original scale using the mean and standard deviation of the train data.

6.4.3 Evaluating the Predictive Performance of OrdPLSPredict and OrdPLScpredict

To evaluate the predictive performance of OrdPLSPredict and OrdPLScpredict, the RMSE or MAE as proposed for PLSpredict can be used. However, in the context of OrdPLSPredict and OrdPLScpredict, for ordinal indicators they are interpreted as penalties of 0 in the presence of exact predictions, as penalties of 1 if a category $h - 1$ or $h + 1$ is predicted for the observed category h , and as penalties of 2 (MAE) or 4 (RMSE) if a category $h - 2$ or $h + 2$ is predicted for category h , and so on. Moreover, the misclassification error rate (MER) can be computed as the fraction of incorrect classifications (James et al., 2021).

Fig. 6.3 Structural model

6.5 Monte Carlo Simulation

To assess the performance of OrdPLSpredict and OrdPLScpredict, we conducted a Monte Carlo simulation. Specifically, we compared the accuracy of predictions for continuous and ordinal indicators obtained by OrdPLScpredict, OrdPLSpredict, PLSpredict, and PLSpredict. For OrdPLSpredict and OrdPLScpredict, we used the *mean* and the *median* approach to obtain the final predictions of the indicators. To ensure a fair comparison of the different methods in terms of their ability to predict ordinal indicators, we rounded the original continuous predictions of PLSpredict and PLSpredict to integer values to obtain predicted categories.

6.5.1 Simulation Design

To compare the various approaches' performance, we considered a population model consisting of one exogenous latent variable η_1 and two endogenous latent variables η_2 and η_3 . We assumed all latent variables to be standardized and related via a structural model, as follows:

$$\eta_2 = 0.6 \cdot \eta_1 + \zeta_1 \quad (6.25)$$

$$\eta_3 = 0.0 \cdot \eta_1 + 0.6 \cdot \eta_2 + \zeta_2 \quad (6.26)$$

Note, the exogenous latent variable η_1 is assumed to be uncorrelated with the structural disturbance terms ζ_1 and ζ_2 . The structural model is displayed in Fig. 6.3.

Additionally, we measured each of the three latent variables by three indicators; therefore, x_{11} , x_{12} , and x_{13} loaded on η_1 with factor loadings of 0.8, 0.7, and 0.6, respectively; x_{21} , x_{22} , and x_{23} loaded on η_2 each with a factor loading of 0.7; and x_{31} , x_{32} , and x_{33} loaded on η_3 with factor loadings of 0.5, 0.7, and 0.9, respectively. Similar to the latent variables, the indicators were assumed to be standardized. Moreover, all structural disturbance terms and random measurement errors were assumed to be uncorrelated. Similarly, the latent variables were assumed to be uncorrelated with the random measurement errors. Consequently, we could give the population correlation matrix of the indicators as follows:

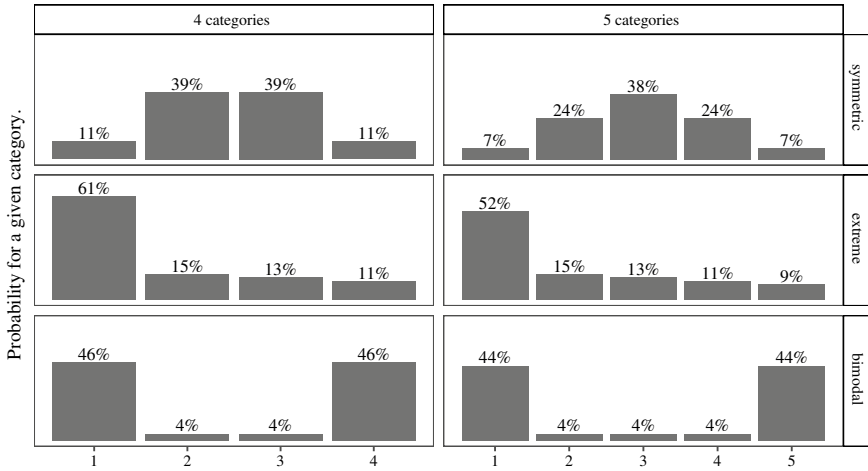


Fig. 6.4 Probability distribution of the ordinal indicators

$$\Sigma = \begin{pmatrix}
 x_{11} & x_{12} & x_{13} & x_{21} & x_{22} & x_{23} & x_{31} & x_{32} & x_{33} \\
 1.000 & & & & & & & & \\
 0.560 & 1.000 & & & & & & & \\
 0.480 & 0.420 & 1.000 & & & & & & \\
 0.336 & 0.294 & 0.252 & 1.000 & & & & & \\
 0.336 & 0.294 & 0.252 & 0.490 & 1.000 & & & & \\
 0.336 & 0.294 & 0.252 & 0.490 & 0.490 & 1.000 & & & \\
 0.144 & 0.126 & 0.108 & 0.210 & 0.210 & 0.210 & 1.000 & & \\
 0.202 & 0.176 & 0.151 & 0.294 & 0.294 & 0.294 & 0.350 & 1.000 & \\
 0.259 & 0.227 & 0.194 & 0.378 & 0.378 & 0.378 & 0.450 & 0.630 & 1.000
 \end{pmatrix} \quad (27)$$

To examine the approaches' performance in predicting ordinal indicators' values, the values of indicators x_{11} , x_{13} , x_{21} , x_{23} , and x_{33} were transformed as described in Sect. 6.2. In doing so, we varied the number of categories between four and five, and we considered three sets of threshold parameters. In the case of symmetrically distributed threshold parameters, all five indicators were categorized using the following threshold parameters: $-\infty, -1.25, 0, 1.25, \infty$ and $-\infty, -1.5, -0.5, 0.5, 1.5, \infty$, respectively. Similarly, for the extreme asymmetric threshold parameter distribution, we set the thresholds to $-\infty, 0.28, 0.71, 1.23, \infty$ in the case of four categories, and to $-\infty, 0.05, 0.44, 0.84, 1.34, \infty$ in the case of five categories. Moreover, we considered a bimodal distribution of the ordinal variables. In this case, the thresholds were set to $-\infty, -0.1, 0, 0.1, \infty$ for four categories and to $-\infty, -0.15, -0.05, 0.05, 0.15, \infty$ for five categories. Figure 6.4 shows the corresponding probability distributions for the categories of the ordinal indicators.

The complete Monte Carlo simulation was carried out in the statistical programming environment R (R Core Team, 2021). To assess the influence of the train dataset's sample size on the approaches' predictive performance, we varied the sample sizes of the train dataset from 200, 500, and 1,000 observations. Hence, in total, we had 36 conditions: three different sample sizes of the train dataset (200, 500, and 1,000 observations) \times two different numbers of categories for the ordinal indicators (four and five categories) \times three sets of threshold parameters (symmetric, extreme asymmetric and bimodal) \times two ways to obtain the final predictions for OrdPLScpredict and OrdPLSpredict (mean and median approach). To assess the approaches' predictive performance, we considered test datasets containing $N = 100$ observations. Additionally, we focused on the following three predictive performance measures: (i) MAE, (ii) RMSE, and (iii) MER. Small values of these measures indicate accurate predictions. For each condition, we conducted 500 simulation runs. In each run, we drew a dataset from the multivariate standard normal distribution with a mean vector of $\mathbf{0}$ and the correlation matrix shown in Eq. (27) using the `mvrnorm()` function of the MASS package (Venables & Ripley, 2002). The number of draws equaled the train dataset's sample size from the corresponding condition plus the 100 observations of the test dataset. Subsequently, we categorized the observations for the variables x_{11} , x_{13} , x_{21} , x_{23} , and x_{33} to obtain ordinal variables using threshold parameters from the corresponding condition. To estimate the model by PLS, PLSc, OrdPLS, and OrdPLSc, we used the `csem()` function of the R package cSEM (Rademaker & Schuberth, 2020). In doing so, the path weighting scheme was used for inner weighting and Mode A was used to calculate the weights to form the proxies for the latent variables. Additionally, we replaced inadmissible estimations, i.e., each condition was based on 500 valid estimations. An inadmissible estimation suffers from at least one of the following problems: (i) the PLS algorithm has not converged, (ii) at least one reliability estimate is larger than 1, (iii) at least one absolute factor loading estimate is larger than 1, (iv) the model-implied construct correlation matrix is not positive semi-definite, and/or (v) the model-implied indicator correlation matrix is not positive semi-definite. Next, to apply OrdPLScpredict, OrdPLSpredict, PLScpredict, and PLSpredict, we used the `predict()` function of the R package cSEM to obtain the predictions for the indicators associated with endogenous constructs.

6.5.2 Simulation Results

In this section, we present the results of our Monte Carlo simulation. Since the results for the ordinal indicators and the continuous indicators, respectively, are very similar, we only present the results for the ordinal indicator x_{23} and for the continuous indicator x_{31} . Further, the results for four and five categories are very similar. Therefore, we only report the results for four categories. Furthermore, the results are only slightly affected by the train dataset's sample size. Hence, we only report the results for 500 observations. Finally, the results for the mean and median approaches used to obtain the predictions with OrdPLScpredict and OrdPLSpredict hardly differ. Therefore,

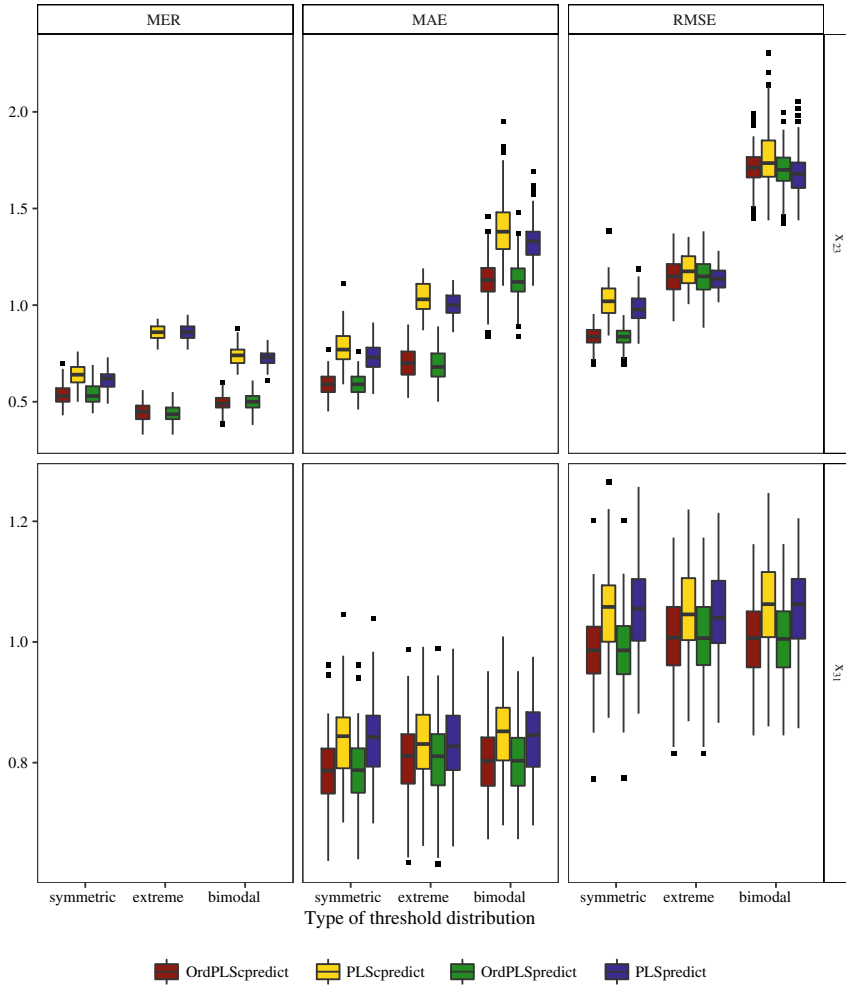


Fig. 6.5 Box plots of the performance measures

we report only the results for the mean approach. The complete results are given in the Supplementary Material.

Figure 6.5 shows the boxplots for the three performance measures, namely, the MER, the MAE, and the RMSE over the 500 simulation runs. In specific, the line in the boxes illustrate the median of the respective performance measure over all simulation runs. While the expanse of the boxes and the length of the whiskers illustrate the variation of the performance measures over all simulation runs, the dots depict outliers. Since the MER is not consistent with the nature of continuous variables for which exact matches are not expected, we only consider the MER for the ordinal indicator x_{23} .

Considering the ordinal indicator x_{23} , OrdPLScpredict and OrdPLSpredict produce very similar results. The same is observed for PLScpredict and PLSpredict. Also, the four approaches produce very similar results in terms of medians of the RMSE for an extremely asymmetric threshold distribution and a bimodal distribution of the ordinal variables. However, there are also differences between the approaches. Considering the median of the MAE and the MER, OrdPLScpredict and OrdPLSpredict outperform PLScpredict and PLSpredict. This is particularly obvious in the case of the extreme asymmetric threshold parameter distribution and a bimodal distribution of the ordinal variables. Considering the variation of predictions per condition, the approaches perform similar for all the performance measures. Considering the continuous indicator x_{31} , OrdPLSpredict and OrdPLScpredict outperform their traditional counterparts considering the median performance measures.

6.5.3 *Simulation Insights*

The results of our Monte Carlo simulation show that correcting for attenuation in case of latent variables does not increase the prediction accuracy. Moreover, they show that OrdPLSpredict and OrdPLScpredict outperform PLSpredict and PLScpredict for most of the simulation conditions in terms of performance measures' median. This is particularly the case for ordinal indicators with an extremely asymmetric or bimodal distribution of the categories. Finally, the approach to determine the location of the distribution of the predictions of the latent variables underlying the ordinal variables, i.e., mean and median, does not influence the predictive performance. However, the results are less clear for the RMSE in combination with extremely skewed threshold distribution or bimodally distributed ordinal variables.

6.6 Guidelines on Performing Predictions Using the R Package cSEM

To illustrate how researchers can apply OrdPLScpredict, OrdPLSpredict, PLScpredict, and PLSpredict, we provide guidelines for the open source R package cSEM. In doing so, we focus on a model that Hwang and Takane (2004) studied. We display their model in Fig. 6.6. To preserve clarity, we have omitted the random measurement error terms and the structural error terms. For a motivation of the model, the interested reader is referred to the article of Hwang and Takane (2004). As Fig. 6.6 shows, the model consists of the following four latent variables: organizational prestige (OrgPres), organizational identification (OrgIden), affective commitment (Love) (Afflove), and affective commitment (Joy) (AffJoy). Bergami and Bagozzi (2000) give an elaboration of the constructs. The considered dataset is part of the survey data used in Bergami and Bagozzi's (2000) study. It consists of 305 observations for the 21 indicators. Each indicator is measured on a 5-point scale ranging from 1

(=strongly disagree) to 5 (=strongly agree), i.e., all indicators are ordinal. A detailed description of the indicators can be found in Henseler (2021, Table 6.1).

As a first step, we need to estimate the model parameters. For this purpose, we can use the `csem()` function of the `cSEM` R package. In general, the `csem()` function requires a dataset and a model as input.

To specify models in `cSEM`, `lavaan` syntax (Rosseel, 2012) is used. Specifically, the `'=~'` operator is used to specify the relationship between indicators and latent variables, the `'<~'` operator is used to specify indicators forming an emergent variable, and the `'~'` operator is used to specify the structural model. The specification for the model illustrated in Fig. 6.6 is given as follows:

```
.model <-"
#Measurement models
OrgPres =~ cei1 + cei2 + cei3 + cei4 + cei5 + cei6 + cei7 + cei8
OrgIden =~ ma1 + ma2 + ma3 + ma4 + ma5 + ma6
AffJoy =~ orgcmt1 + orgcmt2 + orgcmt3 + orgcmt7
AffLove =~ orgcmt5 + orgcmt6 + orgcmt8

# Structural model
OrgIden ~ OrgPres
AffLove ~ OrgIden
AffJoy ~ OrgIden
"
```

The dataset we use here is publicly available and also provided in the `cSEM` R package. However, as it is provided in the `cSEM` package all indicators are labeled as *numeric*. In this case, the `csem()` function uses the Pearson correlations to estimate

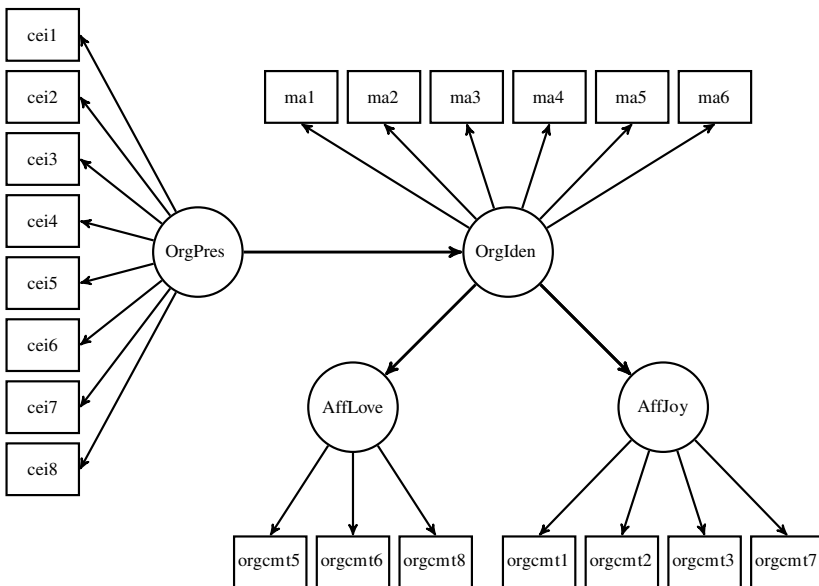


Fig. 6.6 Model from Hwang and Takane (2004)

the model parameters, i.e., PLS or PLSc is employed. To use the polychoric/polyserial correlations, and thus to apply OrdPLS or OrdPLSc, the ordinal indicators need to be labeled as *ordered factors*, as shown in the following:

```
library(cSEM)

# Load the data from the cSEM package
data(BergamiBagozzi2000)

# Transform the numerical indicators into factors
data_new <- as.data.frame(lapply(BergamiBagozzi2000, as.ordered))
```

Finally, to estimate the model parameters, the dataset and the specified model are provided as input to the `csem()` function as the following shows:

```
res <- csem(.model = .model, .data = data_new[1:250,], .resample_method = "bootstrap")
```

To evaluate our model's predictive performance, we use only the first 250 observations of the dataset for the estimation. Note that cSEM applies a correction for attenuation by default if latent variables are included in the model, i.e., PLSc or OrdPLSc is used. Further, by default, cSEM uses the path weighting scheme to calculate the inner weights. In our case, the specified model was estimated by OrdPLSc since at least one indicator is labeled as factor and the model comprises at least one latent variable. If the user aims for statistical inference about the parameter estimates, the argument `'resample_method'` has to be set to either `'bootstrap'` or `'jackknife'`, otherwise no standard errors will be estimated. In our case, we used bootstrap for statistical inference. By default, 499 bootstrap runs are conducted. A summary of the estimated model can be obtained via the `summarize()` function.

To assess the estimated model's predictive performance, the `predict()` function is used. Evaluating the predictive performance of a model requires benchmark predictions. For that purpose, the `'benchmark'` argument of the `predict()` function can be used to determine how benchmark predictions are obtained. In cases where the original model was estimated by OrdPLSc or OrdPLS, the benchmark predictions are rounded for the ordinal indicators if PLS or PLSc were used to estimate the benchmark model. If predictions based on OrdPLSpredict were to be used as benchmark, the `'benchmark'` argument must be set to `'PLS-PM'`, the argument `'treat_as_continuous'` must be set to `'FALSE'`, and the argument `'disattenuate'` has to be set to `'FALSE'` to prevent a correction for attenuation. In the case of OrdPLSpredict and OrdPLSpredict, by default $n_{\text{pred}} = 100$ draws are performed from the multivariate truncated normally distributed unobservable variables underlying the ordinal indicators associated with exogenous constructs. To determine the location of the distribution of the n_{pred} predictions for the indicators of the endogenous constructs, the `'mean'`, `'median'`, or `'mode'` approach can be used. The approach can be chosen separately for the target predictions and the benchmark predictions using the arguments `'approach_score_target'` and `'approach_score_benchmark'`, respectively.

The `predict()` function allows the user to provide a test dataset via the `'test_data'` argument. If no test dataset is provided, k fold cross-validation is applied, i.e., the dataset from the original estimation is randomly split into k (approximately)

equal parts. Subsequently, the values of each part are predicted based on a model estimated on the basis of the remaining parts. To adjust the number of cross-validation folds the `'cv_folds'` argument is used. By default this argument is set to 10. To minimize the effect of random splitting in k fold cross-validation, the k fold cross-validation is repeated several times (Shmueli et al., 2019). In the `predict()` function, the number of repetitions is adjusted via the argument `'r'`. If a test dataset is provided, no k fold cross-validation is conducted and predictions are performed based on the observations of the test dataset.

For the considered empirical example, we used `PLScpredict` as benchmark, and provided the last 55 observations of the original dataset as test dataset, i.e., only the observations of the test dataset were predicted. Additionally, we used the `'median'` approach to obtain predictions in `OrdPLScpredict`. The results of the `predict()` function are as follows.

```
pred <- predict (.object = res, .benchmark = "PLS-PM", .test_data = data_new[( 251): 305,],
               .treat_as_continuous = TRUE, .approach_score_target = "median")

pred
```

The output contains some general information in the top. Moreover, the user can choose the performance measures to assess the accuracy of the predictions for the indicators associated with the endogenous constructs. For our example, we report the MAE and the MER:

```
print(pred, .metrics = c("MAE", "MER"))

## -----
##                               Overview
## -----
## Number of obs. training      = 250
## Number of obs. test        = 55
## Number of cv folds          = NA
## Number of repetitions        = 1
## Handle inadmissibles        = stop
## Estimator target            = 'OrdPLS'
## Estimator benchmark         = 'PLS-PM'
## Disattenuation target       = 'TRUE'
## Disattenuation benchmark    = 'TRUE'
## Approach to predict         = 'earliest'
## -----
##                               Prediction metrics -----
## -----
## Name      MAE target MAE benchmark MER target MER benchmark
## ma1      0.5091     1.3455     0.4727     0.8909
## ma2      0.4545     1.3818     0.4364     0.9455
## ma3      0.5636     1.0000     0.5091     0.7273
## ma4      0.6000     1.6909     0.5273     0.9455
## ma5      0.6909     1.6545     0.5636     0.9273
## ma6      0.5636     1.2182     0.5091     0.8364
## orgcmt5  0.4000     0.9818     0.3636     0.8182
## orgcmt6  0.4727     0.6545     0.4545     0.5636
## orgcmt8  0.7455     0.9091     0.6182     0.6909
## orgcmt1  0.6727     1.1818     0.6000     0.8182
## orgcmt2  0.5818     1.1273     0.5455     0.8182
## orgcmt3  0.5455     1.1273     0.5091     0.8000
## orgcmt7  0.5636     0.8364     0.5091     0.7091
## -----
```

Considering our example's MAE, the results show that `OrdPLScpredict` outperforms `PLScpredict`. Moreover, it shows that the MER is smaller for all indicators in the case of `OrdPLScpredict`, which indicates more accurate predictions than those obtained by `PLScpredict`. These results are also in line with the findings of our Monte Carlo simulation.

In general, the cSEM R package provides users with a lot of flexibility. For more details about the package, we refer the interested reader to the manual. Also, additional tutorials using the cSEM package can be found in Henseler (2021).

6.7 Discussion

The past decade has seen increased scholarly attention to evaluating the predictive power of models estimated by PLS (e.g., Shmueli et al., 2016; Carrión et al., 2016; Shmueli et al., 2019). This is mainly due to the *causal-predictive* nature of PLS (Chin et al., 2020). However, as Schubert et al. (2023a) has emphasized, if PLS is applied in the context of explanatory modeling, i.e., in theory testing, researchers should not rely solely on predictive metrics for model evaluation, but should also consider all possible means known from explanatory modeling for model assessment, including overall model fit assessment.

In this chapter, we have focused on the predictive power of models estimated by OrdPLS and OrdPLSc. Specifically, we presented OrdPLScpredict and OrdPLSpredict. The two approaches are similar to those known from PLS and PLSc to perform predictions, namely PLSpredict and PLScpredict. In contrast to PLSpredict and PLScpredict, our two proposed approaches take the nature of ordinal indicators into account. Additionally, our approaches resemble those Cantaluppi and Schubert (2019) proposed. However, our two proposed approaches are not limited to models containing only ordinal indicators.

The results of our Monte Carlo simulation to evaluate OrdPLScpredict's performance provides several interesting insights. First, OrdPLScpredict and OrdPLSpredict outperform PLScpredict and PLSpredict in cases where values of continuous indicators are predicted. Second, considering the MER and the MAE evaluation metric, OrdPLScpredict and OrdPLSpredict outperform the other approaches in cases where values of ordinal indicators are predicted. Third, the approach to determine the location of the distribution of the predictions of the latent variables underlying the ordinal indicators, i.e., mean or median approach, does not influence the predictive performance of OrdPLSpredict and OrdPLScpredict. Finally, comparing the performance of OrdPLScpredict and OrdPLSpredict to the performance of PLScpredict and PLSpredict, the results show that not correcting for attenuation, even if the parameter estimates are not consistent, does not lead to a worse predictive performance.

A crucial point in predictive research is the principle that estimation should be based solely on the train dataset, while evaluating predictions should be based solely on the test dataset (James et al., 2021). In OrdPLSpredict and OrdPLScpredict, we simulate values for the indicators connected to exogenous constructs from a multivariate truncated normal distribution if ordinal variables are present. Specifically, we use a variance-covariance matrix that equals the estimated correlation matrix and truncation limits that equal the estimated thresholds of the train dataset. Although using estimated threshold parameters based on the train dataset is the only feasible solution in various situations, e.g., in cases of small test datasets, future research

should evaluate ways of rendering the current prediction method more robust to situations in which the test data set's correlation structure slightly differs from the one observed on the train data set. For instance, research could consider the effect on a model's predictive performance if the test dataset's correlation matrix instead the train dataset's is used for simulating the scores of exogenous constructs. Moreover, the polychoric and polyserial correlation used in OrdPLS/OrdPLSc show some limitations. For instance, they assume that the continuous variables underlying the ordinal indicators are multivariate normally distributed and require a sufficient sample size to avoid empty cells, i.e., a combination of categories is not present in the dataset. To address a violation of the normality assumption, a robust version of the polychoric correlation has recently been proposed (Lyhagen & Ornstein, 2023). Future research should investigate how this robust version can be used in OrdPLSpredict and OrdPLScpredict. To obtain the final predictions of OrdPLSpredict and OrdPLScpredict, we use the mean, median and mode approach. However, the linear combination of truncated multivariate normally distributed variables does not need to be unimodally distributed, which might negatively impact the performance of the mean, median and mode approach. Consequently, future research should study the distribution of the sets of the predicted exogenous construct scores in more detail. Furthermore, although our results were hardly affected by the performance measure used, i.e., the RMSE, the MAE and the MER, future research should consider further metrics to evaluate the predictive performance in case of ordinal indicators and provide guidelines on which measure to use in which situation. Finally, simulation studies are limited regarding their design. Consequently, future research should evaluate the effect of our chosen simulation parameters, e.g., test data sample size and model complexity, and compare the performance of OrdPLScpredict to predictions based on other estimation methods such as non-metric partial least squares (Russolillo, 2012), generalized structured component analysis (Hwang & Takane, 2004), approaches to generalized canonical correlation analysis (Kettenring, 1971) or maximum-likelihood estimator (Jöreskog, 1970; Schubert, 2023).

References

- Becker, J. M., Klein, K., & Wetzels, M. (2012). Hierarchical latent variable models in PLS-SEM: Guidelines for using reflective-formative type models. *Long Range Planning*, 45(5), 359–394.
- Benitez, J., Henseler, J., Castillo, A., & Schubert, F. (2020). How to perform and report an impactful analysis using partial least squares: Guidelines for confirmatory and explanatory IS research. *Information & Management*, 57(2), 1–16.
- Bergami, M., & Bagozzi, R. P. (2000). Self-categorization, affective commitment and group self-esteem as distinct aspects of social identity in the organization. *British Journal of Social Psychology*, 39(4), 555–577.
- Braojos, J., Benitez, J.e., Llorens, J., & Ruiz, L. (2020). Impact of IT integration on the firm's knowledge absorption and desorption. *Information & Management*, 57(7), 103–290.
- Cantaluppi, G. (2012). A partial least squares algorithm handling ordinal variables also in presence of a small number of categories. arXiv preprint, [arXiv:1212.5049](https://arxiv.org/abs/1212.5049)

- Cantaluppi, G., & Boari, G. (2016). A partial least squares algorithm handling ordinal variables. In H. Abdi, V. Esposito Vinzi, G. Russolillo, G. Saporta, & L. Trinchera (Eds.), *The multiple facets of partial least squares and related methods: PLS, Paris, France, 2014* (pp. 295–306). Switzerland: Springer International Publishing.
- Cantaluppi, G., & Schuberth, F. (2019). A prediction method for ordinal consistent partial least squares. In G. Arbia, S. Peluso, A. Pini, & G. Rivellini (Eds.), *Smart statistics for smart applications—Book of short papers SIS2019*. Milan.
- Carrión, G. C., Henseler, J., Ringle, C. M., & Roldán, J. L. (2016). Prediction-oriented modeling in business research by means of PLS path modeling: Introduction to a JBR special section. *Journal of Business Research*, 69(10), 4545–4551.
- Chin, W., Cheah, J. H., Liu, Y., Ting, H., Lim, X. J., & Cham, T. H. (2020). Demystifying the role of causal-predictive modeling using partial least squares structural equation modeling in information systems research. *Industrial Management & Data Systems*, 120(12), 2161–2209.
- Cho, G., & Choi, J. Y. (2020). An empirical comparison of generalized structured component analysis and partial least squares path modeling under variance-based structural equation models. *Behaviormetrika*, 47, 243–272.
- Dijkstra, T. K. (1985). *Latent variables in linear stochastic models: Reflections on “Maximum Likelihood” and “Partial Least Squares” methods* (Vol. 2). Amsterdam: Sociometric Research Foundation.
- Dijkstra, T. K. (2011). Consistent partial least squares estimators for linear and polynomial factor models. Technical Report. <https://doi.org/10.13140/RG.2.1.3997.0405>
- Dijkstra, T. K. (2013). A note on how to make PLS consistent. Technical Report. <https://doi.org/10.13140/RG.2.1.4547.5688>
- Dijkstra, T. K. (2017). A perfect match between a model and a mode. In H. Latan & R. Noonan (Eds.), *Partial least squares path modeling: Basic concepts, methodological issues and applications* (pp. 55–80). Cham: Springer.
- Dijkstra, T. K., & Henseler, J. (2015a). Consistent and asymptotically normal PLS estimators for linear structural equations. *Computational Statistics & Data Analysis*, 81, 10–23.
- Dijkstra, T. K., & Henseler, J. (2015b). Consistent partial least squares path modeling. *MIS Quarterly*, 39(2), 29–316.
- Dragso, F. (1986). Polychoric and polyserial correlations. In S. Kotz & N. Johnson (Eds.), *The encyclopedia of statistics* (Vol. 7, pp. 68–74). New York: John Wiley.
- Evermann, J., & Tate, M. (2014). Comparing out-of-sample predictive ability of PLS, covariance, and regression models. In *Proceedings of the 35th International Conference on Information Systems*. Association for Information Systems (AIS).
- Henseler, J. (2021). *Composite-based structural equation modeling: Analyzing latent and emergent variables*. New York, NY: Guilford Press.
- Hui, B. S., & Wold, H. (1982). Consistency and consistency at large of partial least squares estimates. In K. G. Jöreskog & H. Wold (Eds.), *Systems under indirect observation: Causality, structure, prediction Part II* (pp. 119–130). Amsterdam: North-Holland.
- Hwang, H., & Takane, Y. (2004). Generalized structured component analysis. *Psychometrika*, 69(1), 81–99.
- Jakobowicz, E., & Derquenne, C. (2007). A modified PLS path modeling algorithm handling reflective categorical variables and a new model building strategy. *Computational Statistics & Data Analysis*, 51(8), 3666–3678.
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). *An introduction to statistical learning*. New York: Springer.
- Jöreskog, K. G. (1970). A general method for estimating a linear structural equation system. *ETS Research Bulletin Series*, 1970(2), i–41.
- Kettenring, J. R. (1971). Canonical analysis of several sets of variables. *Biometrika*, 58(3), 433–451.
- Lee, S. Y., & Poon, W. Y. (1986). Maximum likelihood estimation of polyserial correlations. *Psychometrika*, 51(1), 113–121.

- Lohmöller, J. B. (1989). *Latent variable path modeling with partial least squares*. Heidelberg: Physica-Verlag.
- Lyhagen, J., & Ornstein, P. (2023). Robust polychoric correlation. *Communications in Statistics—Theory and Methods*, *52*(10), 3241–3261.
- Miltgen, C. L., Henseler, J., Gelhard, C., & Popovič, A. (2016). Introducing new products that affect consumer privacy: A mediation model. *Journal of Business Research*, *69*(10), 4659–4666.
- Noonan, R., & Wold, H. (1982). PLS path modeling with indirectly observed variables: A comparison of alternative estimates for the latent variable. In K. G. Jöreskog & H. Wold (Eds.), *Systems under indirect observation: Causality, structure, prediction part II* (pp. 75–94). Amsterdam: North-Holland.
- Pearson, K. (1900). Mathematical contributions to the theory of evolution. VII. on the correlation of characters not quantitatively measurable. *Philosophical Transactions of the Royal Society of London Series A (Containing Papers of a Mathematical or Physical Character)*, *195*, 1–47 & 405
- Pearson, K. (1913). On the measurement of the influence of “broad categories” on correlation. *Biometrika*, *9*(1/2), 116–139.
- Poon, W. Y., & Lee, S. Y. (1987). Maximum likelihood estimation of multivariate polyserial and polychoric correlation coefficients. *Psychometrika*, *52*(3), 409–430.
- R Core Team (2021). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>
- Rademaker, M. E., & Schubert, F. (2020). cSEM: Composite-based structural equation modeling. <https://m-e-rademaker.github.io/cSEM/> package version: 0.4.0.9000
- Rademaker, M. E., Schubert, F., & Dijkstra, T. K. (2019). Measurement error correlation within blocks of indicators in consistent partial least squares. *Internet Research*, *29*(3), 448–463.
- Rosseel, Y. (2012). Lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, *48*(2), 1–36.
- Russolillo, G. (2012). Non-metric partial least squares. *Electronic Journal of Statistics*, *6*, 1641–1669.
- Sarstedt, M., & Danks, N. P. (2022). Prediction in HRM research—A gap between rhetoric and reality. *Human Resource Management Journal*, *32*, 485–513.
- Sarstedt, M., Henseler, J., & Ringle, C. M. (2011). Multigroup analysis in partial least squares (PLS) path modeling: Alternative methods and empirical results. *Advances in Interational Marketing*, *22*, 195–218.
- Sarstedt, M., Hair, J. F., & Ringle, C. M. (2023). “PLS-SEM: Indeed a silver bullet”—Retrospective observations and recent advances. *Journal of Marketing Theory and Practice*, *31*(3), 261–275.
- Schamberger, T., Schubert, F., & Henseler, J. (2023). Confirmatory composite analysis in human development research. *International Journal of Behavioral Development*, *47*(1), 89–100.
- Schubert, F. (2021). Confirmatory composite analysis using partial least squares: Setting the record straight. *Review of Managerial Science*, *15*, 1311–1345.
- Schubert, F. (2023). The Henseler-Ogasawara specification of composites in structural equation modeling: A tutorial. *Psychological Methods*, *28*(4), 843–859.
- Schubert, F., & Cantaluppi, G. (2017). Ordinal consistent partial least squares. In L. Hengky & R. Noonan (Eds.), *Partial least squares path modeling* (pp. 109–150). Switzerland: Springer.
- Schubert, F., Henseler, J., & Dijkstra, T. K. (2018). Partial least squares path modeling using ordinal categorical indicators. *Quality & Quantity*, *52*(1), 9–35.
- Schubert, F., Rademaker, M. E., & Henseler, J. (2020). Estimating and assessing second-order constructs using PLS-PM: the case of composites of composites. *Industrial Management & Data Systems*, *120*(12), 2211–2241.
- Schubert, F., Rademaker, M. E., & Henseler, J. (2023a). Assessing the overall fit of composite models estimated by partial least squares path modeling. *European Journal of Marketing*, *57*(6), 1678–1702.

- Schuberth, F., Zaza, S., Henseler, J. (2023b). Partial least squares is an estimator for structural equation models: A comment on Evermann and Rönkkö. *Communications of the Association for Information Systems*, 52, 711–714.
- Shmueli, G. (2010). To explain or to predict? *Statistical Science*, 25(3), 289–310.
- Shmueli, G., Ray, S., Estrada, J. M. V., & Chatla, S. B. (2016). The elephant in the room: Predictive performance of PLS models. *Journal of Business Research*, 69(10), 4552–4564.
- Shmueli, G., Sarstedt, M., Hair, J. F., Cheah, J. H., Ting, H., Vaithilingam, S., & Ringle, C. M. (2019). Predictive model assessment in PLS-SEM: Guidelines for using PLSpredict. *European Journal of Marketing*, 53(11), 2322–2347.
- Tenenhaus, M., Vinzi, V. E., Chatelin, Y. M., & Lauro, C. (2005). PLS path modeling. *Computational Statistics & Data Analysis*, 48(1), 159–205.
- Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* (4th ed.). New York: Springer.
- Vogt, W. (1993). *Dictionary of statistics and methodology*. London: Sage.
- Wold, H. (1966). Estimation of principal components and related models by iterative least squares. In P. Krishnaiah (Ed.), *Multivariate Analysis* (pp. 391–420). New York: Academic Press.
- Wold, H. (1974). Causal flows with latent variables: Partings of the ways in the light of NIPALS modelling. *European Economic Review*, 5(1), 67–86.
- Wold, H. (1982). Soft modeling: The basic design and some extensions. In K. G. Jöreskog & H. Wold (Eds.), *Systems under indirect observation: Causality, structure, prediction Part II* (pp. 1–54). Amsterdam: North-Holland.
- Yu, X., Zaza, S., Schuberth, F., Henseler, J. (2021). Counterpoint: Representing forged concepts as emergent variables using composite-based structural equation modeling. *ACM SIGMIS Database: the DATABASE for Advances in Information Systems*, 52(S1), 114–130.

Tamara Schamberger is a Postdoctoral Researcher at the Chair of Econometrics in the Faculty of Business Management and Economics at the University of Würzburg, Germany. She obtained her Ph.D., in Econometrics in a joint program of the Faculty of Business Management and Economics, University of Würzburg (Germany) and the Faculty of Engineering Technology, University of Twente (the Netherlands). Her main research interests are on methodological advances in composite-based structural equation modeling.

Gabriele Cantaluppi is Associate Professor of Statistics. He graduated in Economics at Università Cattolica del Sacro Cuore and got his Ph.D., in Methodological Statistics from University of Trento. His main research interests are Total Quality Management, in particular focusing on the analysis and the measurement of Customer Satisfaction by means of Structural Equation Models with latent variables.

Florian Schuberth is Assistant Professor at the Chair of Product-Market Relations in the Faculty of Engineering Technology at the University of Twente, the Netherlands. He obtained his Ph.D., Degree in Econometrics at the Faculty of Business Management and Economics, University of Würzburg, Germany. His main research interests are structural equation modeling, in particular on composite-based estimators and their enhancement. He is also co-inventor of confirmatory composite analysis (CCA) and co-developer of cSEM, an R package for composite-based structural equation modeling. His research has been published in various international peer-reviewed journals, including *Psychological Methods*, *European Journal of Marketing*, *Information and Management*, and *International Journal of Information Management*.