

## CONSISTENT PARTIAL LEAST SQUARES PATH MODELING

Theo K. Dijkstra

Faculty of Economics and Business, University of Groningen, Nettelbosje 2,  
9747 AE Groningen THE NETHERLANDS {t.k.dijkstra@rug.nl}

Jörg Henseler

Faculty of Engineering Technology, University of Twente, Drienerlolaan 5,  
7522 NB Enschede THE NETHERLANDS {j.henseler@utwente.nl}  
NOVA IMS, Universidade Nova de Lisboa, 1070-312 Lisbon PORTUGAL {jhenseler@isegi.unl.pt}

### Appendix A

In variance-based SEM, the path coefficients can be determined based on the construct score correlation matrix  $\mathbf{R}$ , which for the model depicted in Figure 1 is of the following form (“cor” = correlation; “rel” = reliability):

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & \text{cor}(\xi_1, \xi_2) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\xi}_2)} & \text{cor}(\xi_1, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\eta})} \\ \text{cor}(\xi_1, \xi_2) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\xi}_2)} & 1 & \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_2) \cdot \text{rel}(\tilde{\eta})} \\ \text{cor}(\xi_1, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\eta})} & \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_2) \cdot \text{rel}(\tilde{\eta})} & 1 \end{bmatrix}$$

From  $\mathbf{R}$ , we can extract the submatrix  $\mathbf{R}_x = \begin{bmatrix} 1 & r_{12} \\ r_{12} & 1 \end{bmatrix}$

and the subvector  $\mathbf{R}_{xy} = \begin{bmatrix} r_{13} \\ r_{23} \end{bmatrix}$ .

Under the assumption that the construct scores are standardized, the regression equation equals

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \mathbf{R}_x^{-1} \mathbf{R}_{xy} = \begin{bmatrix} 1 & r_{12} \\ r_{12} & 1 \end{bmatrix}^{-1} \begin{bmatrix} r_{13} \\ r_{23} \end{bmatrix} = \frac{1}{\det(\mathbf{R}_x)} \begin{bmatrix} 1 & -r_{12} \\ -r_{12} & 1 \end{bmatrix} \begin{bmatrix} r_{13} \\ r_{23} \end{bmatrix} = \frac{1}{1-r_{12}^2} \begin{bmatrix} r_{13} - r_{12}r_{23} \\ r_{23} - r_{12}r_{13} \end{bmatrix}$$

Since the true  $\beta_1$  equals zero, we can exploit that  $\text{cor}(\xi_1, \eta) = \text{cor}(\xi_1, \xi_2) \cdot \text{cor}(\xi_2, \eta)$  and thus find the following result for the path  $\hat{\beta}_1$  coefficient as estimated by variance-based SEM:

$$\begin{aligned}
 \hat{\beta}_1 &= \frac{1}{1-r_{12}^2} \cdot (r_{13} - r_{12}r_{23}) \\
 &= \frac{\text{cor}(\xi_1, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\eta})} - \text{cor}(\xi_1, \xi_2) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\xi}_2)} \cdot \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_2) \cdot \text{rel}(\tilde{\eta})}}{1 - \text{cor}^2(\tilde{\xi}_1, \tilde{\xi}_2)} \\
 &= \frac{\text{cor}(\xi_1, \xi_2) \cdot \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\eta})} - \text{cor}(\xi_1, \xi_2) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\xi}_2)} \cdot \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_2) \cdot \text{rel}(\tilde{\eta})}}{1 - \text{cor}^2(\tilde{\xi}_1, \tilde{\xi}_2)} \\
 &= \frac{\text{cor}(\xi_1, \xi_2) \cdot \text{cor}(\xi_2, \eta) \cdot \sqrt{\text{rel}(\tilde{\xi}_1) \cdot \text{rel}(\tilde{\eta})} \cdot (1 - \text{rel}(\tilde{\xi}_2))}{1 - \text{cor}^2(\tilde{\xi}_1, \tilde{\xi}_2)} \\
 &= \frac{\text{cor}(\tilde{\xi}_1, \tilde{\xi}_2)}{1 - \text{cor}^2(\tilde{\xi}_1, \tilde{\xi}_2)} \cdot \text{cor}(\tilde{\xi}_2, \tilde{\eta}) \cdot \frac{1 - \text{rel}(\tilde{\xi}_2)}{\text{rel}(\tilde{\xi}_2)}
 \end{aligned}$$

## Appendix B<sup>1</sup>

A starting point for PLS-analyses is the so-called basic design, in essence a factor model. It is assumed that we have  $N$  i.i.d.<sup>2</sup> column vectors of observed scores  $y_1, y_2, y_3, \dots, y_N$  (i.e., we have  $N$  observations). We assume the usual standardization in PLS of each indicator, which means that each vector  $y_n$  for  $n = 1, 2, 3, \dots, N$  has zero mean and its elements have unit variance. The vectors can be partitioned into  $M$  subvectors,  $M \geq 2$ , each with at least two components. For the  $i^{\text{th}}$  subvector  $y_{in}$  of  $y_n$ , we have

$$y_{in} = \lambda_i \cdot \eta_{in} + \varepsilon_{in} \tag{1}$$

where the loading vector  $\lambda_i$  and the vector of measurement errors  $\varepsilon_{in}$  have the same dimensions as  $y_{in}$ , and the unobservable latent variable  $\eta_{in}$  is real-valued.<sup>3</sup> For convenience, we make the sufficient but by no means necessary assumption that all components of all error vectors are mutually independent, and independent of all latent variables.

The latent variables have also zero mean and unit variance. The correlation between  $\eta_{in}$  and  $\eta_{jn}$  will be denoted by  $\rho_{ij}$ . A particular set of easy implications is that the covariance matrix  $\Sigma_{ii}$  of  $y_{in}$  can be written as

$$\Sigma_{ij} := E y_{ij} y_{jn}^T = \rho_{ij} \lambda_i \lambda_j^T \tag{2}$$

where the covariance matrix of the measurement errors  $\Theta_i$  is diagonal with non-negative diagonal elements, and we have for the covariance between  $y_{in}$  and  $y_{jn}$

<sup>1</sup>To keep the paper self-contained, we collect here the main algebraic results underlying PLSc (for a more elaborate explication, see Dijkstra 2010). Note that there are many ways to correct for inconsistency and the one we use here is quite possibly the simplest one. See Dijkstra (2013) for alternatives and extensions. Any combination of PLS with one of the alternatives could legitimately be called PLSc.

<sup>2</sup>The use of i.i.d. (independent and identically distributed) signifies that the data are seen as a random sample from a population. This is for conceptual convenience only. Together with the classical laws of large numbers and the central limit theorem (Cramér 1946), this entails the consistency and asymptotic normality properties we need. We can make much more general assumptions, but the statistical subtleties would be a distraction and would not lead to different results (CAN-estimators).

<sup>3</sup>We follow exactly the specification of Joreskog (1969). Please note that variables are arranged per rows, and observations per columns. Moreover, note that, as opposed to our custom in the main body of the paper, it pays here to use a subscript for the blocks.

$$\Sigma_{ij} := E y_{in} y_{jn}^T = \rho_{ij} \lambda_i \lambda_j^T \tag{3}$$

The covariance matrix of each  $y_n$  will be denoted by  $\Sigma$ , and the sample covariance matrix by  $S$ . The sample counterparts of  $\Sigma_{ii}$  and  $\Sigma_{ij}$  are denoted by  $S_{ii}$  and  $S_{ij}$ , respectively. Since we assume that the sample data are standardized before being analyzed, we have  $S_{ij} = \frac{1}{N} \sum_{n=1}^N y_{in} y_{jn}^T$ . Note that the assumptions made so far entail that the sample counterparts are consistent and jointly asymptotically normal estimators of the theoretical variance and covariance matrices (Cramér 1946, Chapter 28).

PLS features a number of iterative *fixed-point* algorithms, of which we select one, the so-called Mode A algorithm. This is, in general, numerically the most stable algorithm.<sup>4</sup> As a rule, it converges from arbitrary starting vectors, and it is usually very fast. The outcome is an estimated weight vector  $\hat{w}$ , with typical subvector  $\hat{w}_i$  of the same dimensions as  $y_{in}$ . With these weights, *sample proxies* are defined for the latent variables:  $\hat{\eta}_{in} := \hat{w}_i^T y_{in}$  for  $\eta_{in}$ , with the customary normalization of a unit sampling variance, so  $\frac{1}{N} \sum_{n=1}^N (\hat{w}_i^T y_{in})^2 = \hat{w}_i^T S_{ii} \hat{w}_i = 1$ . In Wold's PLS approach the  $\hat{\eta}_{in}$ 's replace the unobserved latent variables, and loadings and structural parameters are estimated using ordinary least squares.<sup>5</sup> For Mode A, we have for each  $i$  ( $\propto$  stands for "proportional to")

$$\hat{w}_i \propto \sum_{j \in C(i)} sign_{ij} \cdot S_{ij} \hat{w}_j \tag{4}$$

Here,  $sign_{ij}$  is the sign of the sample correlation between  $\hat{\eta}_i$  and  $\hat{\eta}_j$ , and  $C(i)$  is a set of indices of latent variables. Traditionally,  $C(i)$  contains the indices of latent variables adjacent to  $\hat{\eta}_i$  (i.e., the indices of latent variables that appear on the other side of the structural or path equations in which  $\hat{\eta}_i$  appears). Clearly,  $\hat{w}_i$  is obtained by a regression of the indicators  $y_{in}$  on the sign-weighted sum:  $\sum_{j \in C(i)} sign_{ij} \cdot \hat{\eta}_{jn}$ . There are other versions (with correlation weights, for example); this is one of the very simplest, and it is the original one (see Wold, 1982). There is little motivation in the PLS literature for the coefficients of  $S_{ij} \hat{w}_j$ ,<sup>6</sup> but the particular choice can be shown to be irrelevant for the probability limits of the estimators. The algorithm takes an arbitrary starting vector, and then basically follows the sequence of regressions for each  $i$ , each time inserting updates when available (or after each full round; the precise implementation is not important).

Dijkstra (1981, 2010) has shown that the PLS modes converge with a probability tending to one when the sample size tends to infinity, for essentially arbitrary starting vectors. Moreover, the weight vectors that satisfy the fixed-point equations, are locally continuously differentiable functions of the sample covariance matrix of  $y$ . They, as well as other estimators that depend smoothly on the weight vectors and  $S$ , are therefore jointly asymptotically normal.

Let us denote the probability limit of  $\hat{w}_i$ ,  $\text{plim } \hat{w}_i$ , by  $\bar{w}_i$ . We can get it from the equation for  $\hat{w}_i$  by substitution of  $\Sigma$  for  $S$ .

$$\bar{w}_i \propto \sum_{j \in C(i)} sign_{ij} \cdot \sum_j \bar{w}_j = \sum_{j \in C(i)} sign_{ij} \cdot \rho_{ij} \lambda_i \lambda_j^T \bar{w}_j = \lambda_i \sum_{j \in C(i)} sign_{ij} \cdot \rho_{ij} \lambda_j^T \bar{w}_j \tag{5}$$

Since  $\lambda_i^T \bar{w}_j$  is a scalar, and all terms in the sum have the vector  $\lambda_i$  in common, it is clear that  $\bar{w}_i \propto \lambda_i$ . We must have  $\bar{w}_i^T \sum_{ii} \bar{w}_i = 1$ , so

$$\bar{w}_i = \frac{\lambda_i}{\sqrt{\lambda_i^T \sum_{ii} \lambda_i}} \tag{6}$$

<sup>4</sup>This is because Mode A ignores collinearity between the observed variable predictors of the proxy, while Mode B takes account of that collinearity and thus must find the inverse of the covariance matrix of the predictions, which itself is sometimes unstable.

<sup>5</sup>Other estimators like 2SLS are also possible.

<sup>6</sup>In favor of the sign-weights, one could argue that in sufficiently large samples  $sign_{ij} \cdot S_{ij} \hat{w}_j$  is approximately equal to  $\lambda_i \cdot |\rho_{ij}| \cdot (\lambda_j^T \bar{w}_j)$ , where the term in brackets measures the (positive) correlation between  $\eta_j$  and its proxy; see below for results that help justify this claim. So the tighter the connection between  $\eta_j$ , and the better  $\eta_j$  can be measured, the more important  $\hat{\eta}_j$  is in determining  $\hat{w}_i$ .

We conclude that PLS, Mode A, produces estimated weight vectors that tend to vectors proportional to the true loadings. One would like to have a simple estimate for the proportionality factor. We propose here as in Dijkstra (1981, 2010, 2011) to define  $\hat{\lambda}_i := \hat{c}_i \cdot \hat{w}_i$ , where the scalar  $\hat{c}_i$  is such that the *off-diagonal* elements of  $S_{ii}$  are reproduced as well as possible in a least squares sense. So we minimize the Euclidean distance between<sup>7</sup>

$$[S_{ii} - \text{diag}(S_{ii})] \text{ and } [(c_i \cdot \hat{w}_i)(c_i \cdot \hat{w}_i)^T - \text{diag}((c_i \cdot \hat{w}_i)(c_i \cdot \hat{w}_i)^T)] \tag{7}$$

as a function of  $c_i$  and obtain

$$\hat{c}_i := \left[ \frac{\hat{w}_i^T (S_{ii} - \text{diag}(S_{ii})) \hat{w}_i}{\hat{w}_i^T (\hat{w}_i \hat{w}_i^T - \text{diag}(\hat{w}_i \hat{w}_i^T)) \hat{w}_i} \right]^{\frac{1}{2}} \tag{8}$$

The use of “diag” means that only the off-diagonal elements of the matrices are taken into account. So the numerator within brackets is just  $\sum_{a \neq b} \hat{w}_{ia} \hat{w}_{ib} S_{ii,ab}$  and similarly for the denominator. In sufficiently large samples,  $\hat{c}_i$  will be well-defined, real, and positive. (In all samples in this paper and those in other studies,  $\hat{c}_i$  attained proper values.) Its calculation does not require additional numerical optimization. It is straightforward to verify, by replacing  $S_{ii}$  by  $\Sigma_{ii}$  and  $\hat{w}_i$  by  $\bar{w}_i$ , that the correction does its job: the matrix in the denominator equals the matrix in the numerator, apart from a factor  $\frac{1}{\lambda_i^T \sum_{ii} \lambda_i}$ , so

$$\bar{c}_i := \text{plim } \hat{c}_i = \sqrt{\lambda_i^T \sum_{ii} \lambda_i} \tag{9}$$

Now, in particular

$$\text{plim } \hat{\lambda}_i = \text{plim } (\hat{c}_i \cdot \hat{w}_i) = \bar{c}_i \cdot \bar{w}_i = \lambda_i \tag{10}$$

It will be useful to define a population proxy  $\bar{\eta}_{in}$  by  $\bar{\eta}_{in} := \bar{w}_i^T y_{in}$ . Clearly, the squared correlation between a population proxy and its corresponding latent variable is

$$R^2(\eta_i, \bar{\eta}_i) = (\bar{w}_i^T \lambda_i)^2 \tag{11}$$

which equals

$$(\lambda_i^T \lambda_i)^2 \div \lambda_i \sum_{ii} \lambda_i = \frac{(\lambda_i^T \lambda_i)^2}{(\lambda_i^T \lambda_i)^2 + \lambda_i^T \Theta_i \lambda_i} \tag{12}$$

With a large number of high quality indicators, this correlation can be close to one (“consistency at large” in PLS parlance). A trivially deduced but important algebraic relationship is

$$R^2(\bar{\eta}_i, \bar{\eta}_j) = (\bar{w}_i^T \Sigma_{ij} \bar{w}_j)^2 = \rho_{ij}^2 \cdot R^2(\eta_i, \bar{\eta}_i) \cdot R^2(\eta_j, \bar{\eta}_j) \tag{13}$$

indicating that the PLS proxies will tend to underestimate the squared correlations between the latent variables. In fact, one can show that this is true for multiple correlations as well (see Dijkstra 2010). Also note that

$$R^2(\eta_i, \bar{\eta}_j) = (\bar{w}_i^T \lambda_j)^2 = (\bar{w}_i^T \cdot (\bar{w}_i \cdot \bar{c}_i))^2 = (\bar{w}_i^T \bar{w}_i)^2 \cdot \bar{c}_i^2 \tag{14}$$

so that we can estimate the (squared) quality of the proxies consistently by

<sup>7</sup>This assumes that the measurement errors within a block are uncorrelated, the basic design. If we know that some errors are correlated or we have doubts about them, we can delete the items from the difference to be minimized.

$$\hat{R}^2(\eta_i, \bar{\eta}_i) := (\hat{w}_i^T \hat{w}_i)^2 \cdot \hat{c}_i^2 \quad (15)$$

Moreover, with

$$\hat{R}^2(\bar{\eta}_i, \bar{\eta}_j) := (\hat{w}_i^T S_{ij} \hat{w}_j)^2 \cdot \hat{c}_i^2 \quad (16)$$

we can estimate the correlations between the latent variables consistently (see Equation 13)

$$\hat{\rho}_{ij} := \sqrt{\frac{\hat{R}^2(\bar{\eta}_i, \bar{\eta}_j)}{\hat{R}^2(\eta_i, \bar{\eta}_i) \cdot \hat{R}^2(\eta_j, \bar{\eta}_j)}} \quad (17)$$

We close this appendix with four observations:

1. Standard PLS software for Mode A will produce all the necessary, simple ingredients for consistent estimation.
2. The approach can be and has been extended to Mode B (Dijkstra 1981, 2010, 2011), but since Mode A is numerically more stable and faster than Mode B, it has first priority.
3. In the main body of this paper, we used  $\rho(\eta_i)$  for  $R(\eta_i, \bar{\eta}_i)$ .
4. In the main body of this paper, we used (without subscript  $i$ )

$$\rho_{A,i} := (\hat{w}_i^T \hat{w}_i)^2 \cdot \frac{\hat{w}_i^T (S_{ii} - \text{diag}(S_{ii})) \hat{w}_i}{\hat{w}_i^T (\hat{w}_i, \hat{w}_i^T - \text{diag}(\hat{w}_i, \hat{w}_i^T)) \hat{w}_i} \quad (18)$$

This is just  $\hat{R}^2(\eta_i, \bar{\eta}_i)$ .

## References

- Cramér, H. 1946. *Mathematical Models of Statistics* (Volume 9), Princeton, NJ: Princeton University Press.
- Dijkstra, T. K. 1981. *Latent Variables in Linear Stochastic Models: Reflections on "Maximum Likelihood" and "Partial Least Squares" Methods*, Ph.D. Thesis, Groningen University. (A second edition was published in 1985 by Sociometric Research Foundation, Amsterdam.)
- Dijkstra, T. K. 2010. "Latent Variables and Indices: Herman Wold's Basic Design and Partial Least Squares," in *Handbook of Partial Least Squares: Concepts, Methods, and Applications*, V. E. Vinzi, W. W. Chin, J. Henseler, and H. Wang (eds.), New York: Springer, pp. 23-46.
- Dijkstra, T. K. 2011. "Consistent Partial Least Squares Estimators for Linear and Polynomial Factor Models," working paper, University of Groningen, Groningen, The Netherlands (<http://www.rug.nl/staff/t.k.dijkstra/research>).
- Dijkstra, T. K. 2013. "A Note on How to Make PLS Consistent," working paper, University of Groningen, Groningen, The Netherlands (<http://www.rug.nl/staff/t.k.dijkstra/how-to-make-pls-consistent.pdf>).
- Jöreskog, K. G. 1969. "A General Approach to Confirmatory Maximum Likelihood Factor Analysis," *Psychometrika* (34:2), pp. 183-202.
- Wold, H. O. A. 1982. "Soft Modeling: The Basic Design and Some Extensions," in *Systems under indirect observation*, K. G. Jöreskog and H. O. A. Wold (eds.), Amsterdam: North-Holland, pp. 1-54.