Products of Variables in Structural Equation Models

Steven Boker The University of Virginia, Charlottesville

Timo von Oertzen University of the Bundeswehr, Munich

> Joshual Pritikin MakerDAO Immunefi Security

Michael D. Hunter, Timothy Brick The Pennsylvania State University

Andreas Brandmaier Max Planck Institute for Human Development, Berlin MSB Medical School Berlin, Berlin

> Michael Neale Virginia Commonwealth University

> > Draft September 2, 2022

Author Note

Funding for this work was provided in part by NIH Grant R01DA018673, the Max Planck Institute for Human Development, and a fellowship from the University of Zurich Research Priority Program on Dynamics of Healthy Aging. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Institutes of Health. The authors would also like to acknowledge Robert M. Kirkpatrick for his helpful advice on the Monte Carlo simulations reported in this manuscript and the OpenMx development team for their support in user interface design. Correspondence may be addressed to Steven M. Boker, Department of Psychology, The University of Virginia, PO Box 400400, Charlottesville, VA 22903, USA; email sent to boker@virginia.edu; or browsers pointed to https://psychology.as.virginia.edu/people/profile/smb3u

## Abstract

A general method is introduced in which variables that are products of other variables in the context of a structural equation model (SEM) can be decomposed into the sources of variance due to the multiplicands. The result is a new category of SEM which we call a Products of Variables Model (PoV). Some useful and practical features of PoV models include estimation of interactions between latent variables, latent variable moderators, manifest moderators with missing values, and manifest or latent squared terms. Expected means and covariances are analytically derived for a simple product of two variables and it is shown that the method reproduces previously published results for this special case. It is shown algebraically that using centered multiplicands results in an unidentified model, but if the multiplicands have non-zero means, the result is identified. The method has been implemented in OpenMx and  $\Omega$ nyx and is applied in five extensive simulations. 1

Products of Variables in Structural Equation Models

### Introduction

The expected covariance matrix of a Structural Equation Model (SEM) which is 2 composed of a linear combination of variables has long been known (Wright, 1921) as well 3 as its equivalence to a method of path coefficients (Wright, 1934). This equivalence was 4 formalized into what is called the Reticular Action Model (RAM) by McArdle and 5 McDonald (1984). Although linear combinations of variables are powerful and useful, 6 nonlinear models have become increasingly important as latent interaction models, latent 7 moderation models, and nonlinear dynamical systems models have become necessary in the 8 current data-rich research environment. 9

When a variable in an SEM is the outcome of the product of two variables and at 10 least one of the product terms is a manifest variable, there are a variety of methods that 11 can be used to estimate coefficients in the SEM (e.g., Mehta & Neale, 2005; Neale, 1998). 12 But when both of the product terms are latent, the problem becomes more difficult. 13 Methods have been proposed that require nonlinear constraints on measurement models 14 (the product indicant technique, Kenny & Judd, 1984), approximating non-normal 15 distributions using mixture distributions (the latent moderated structure technique, Klein 16 & Moosbrugger, 2000; Moosbrugger, Schermelleh-Engel, & Klein, 1997), multiple group 17 approaches and estimating latent scores so that they can be dealt with as manifest 18 variables (Schumacker, 2002), and variations on unconstrained approaches (unconstrained 19 and quasi-maximum likelihood techniques, Marsh, Wen, & Hau, 2004; Marsh et al., 2007), 20 the approach proposed and elaborated by Wall and Amemiya (Wall & Amemiya, 2001, 21 2003), and semiparametric approaches introduced by Bauer, (Bauer, 2005; Bauer & 22 Hussong, 2009). This is only a small sample of the large literature on nonlinear SEM, but 23 most of these methods apply only in limited cases or are difficult to implement in current 24 SEM software. 25

26

In the current article we take a different approach and present a novel method for the

decomposition and estimation of the variance and covariance of products of variables in 1 SEM models specified in RAM notation. The approach presented here makes the 2 assumption that the non-normality of a variable constructed as the product of other 3 variables can be completely accounted for if a model can account for all sources of variance 4 that are included in the product and that those sources of variance are multivariate 5 normally distributed. The method is currently implemented in OpenMx (Neale et al., 6 2016) and  $\Omega_{nyx}$  (Oertzen, Brandmaier, & Tsang, 2015) and requires only a one line change 7 in the R model script. 8

We begin our argument by providing an extension to RAM path diagrams (see Boker & McArdle, 2005, for an introduction to RAM matrices and path diagrams) in order to show how the specification of these models can be a minimal change from current user interfaces. By adding one new node to the pallette of RAM path diagrams, we can add products of variables into existing SEM models. We hope that authors of other SEM packages will find that this method is straightforward to implement and has many advantages for end users.

To start, consider the construction of a linear combination: estimated coefficients are constants that are multiplied by variables and then all of these products are summed together and assigned to an outcome variable. For instance, a simple bivariate regression using mean centered multivariate normal variables can be written as,

$$z_i = b_1 x_i + b_2 y_i + e_i \tag{1}$$

where  $x_i$  and  $y_i$  are the predictor variables and  $e_i$  is the residual at each row index *i*, and  $b_1, b_2$ , and 1 are constants. If the variance of *x*, *y*, and *e* are  $V_x, V_y$ , and  $V_e$ , respectively and the covariance between *x* and *y* is  $C_{xy}$ , a RAM path diagram isomorphic to this bivariate regression can be drawn as in Figure 1-a. Compare Equation 1 with the path diagram and note that two things are happening when the three one-headed arrows attach to the box representing the manifest variable *z*. In the equation,  $b_1x_i, b_2y_i$ , and  $e_i$  are first summed. Second, the result of that sum is assigned to the variable  $z_i$ . In the equation, summation is represented by one symbol, "+", and assignment is represented by another
symbol, "=". The path diagram fuses together these two symbols from the equation. This
is not a problem as long as *only* linear combinations are used. In that case, summation and
assignment are always used together. However, if one wishes to use an n-ary<sup>1</sup> operator
other than summation, the assignment operator and the n-ary operator must be
disambiguated in the path diagram.



Figure 1. Two equivalent path diagrams of the bivariate regression in Equation 1. (a) Standard RAM path diagram of a bivariate regression. (b) Equivalent path diagram where the n-ary operator addition, represented by  $\oplus$ , is disambiguated from the assignment operator.

6

In Figure 1-b, a new n-ary operator path diagram node is proposed: the summation
operator is represented by a plus sign within a circle. If this summation operator is treated
just like a latent variable with zero unique variance, then the algorithmic application of the
RAM path tracing rules (Boker, McArdle, & Neale, 2002) will correctly produce the
components of covariance of the expected covariance matrix of any SEM model that is

<sup>&</sup>lt;sup>1</sup> Operators often take one or more elements from a set back onto that set. A unary operator maps one member of a set onto the set; whereas a binary operator maps two elements of a set onto that set, and a n-ary operator maps n elements onto that set.

<sup>1</sup> composed solely of linear combinations of variables. Note that in Figure 1-b, the

<sup>2</sup> summation operator could have been mixed together with standard

addition-and-assignment while the standard RAM tracing rules would continue to produce
the correct components of covariance. In many published SEM path diagrams, the plus
sign in a circle is drawn without the plus sign and is called a "dummy variable" whose only
purpose is to sum its input arrows and send the sum to any output arrows. We suggest
that "dummy variables" are actually summation operator nodes since they do not have any
residual variance.

Now that the assignment operator has been distinguished from the summation operator, it is possible to propose a new n-ary operation node: multiplication. Consider the simplest case, a product of two predictor variables,

$$z_i = b_1 x_i y_i + e_i av{2} av{2}$$

<sup>9</sup> where the variables x and y are multivariate normal with mean zero and variances  $V_x$ ,  $V_y$ <sup>10</sup> and covariance  $C_{xy}$ , and the residual e is independent normally distributed with mean zero <sup>11</sup> and variance  $V_e$ .

A path model for Equation 2 is displayed in Figure 2-a, where the n-ary operator 12 multiplication is represented by an asterisk within a circle. Although the summation 13 operator could be treated as if it were a latent variable with zero residual variance, the 14 multiplication operator cannot, as will be demonstrated below. Thus, we will need to 15 specify that an n-ary operator is a fourth type of node in a path diagram: i) squares 16 represent manifest variables; ii) circles represent latent variables; iii) triangles represent 17 constants for means and intercepts; and iv) n-ary operators are represented by a small 18 circle surrounding the selected operator symbol. This distinction is necessary for two 19 reasons that will be described in more detail later in the article: different n-ary operators 20 have different identity operations and also trigger different path tracing rules. 21

Diagramming the problem in this way has numerous advantages. Since scalar multiplication is commutative, neither the vector **x** nor **y** has a privileged place in the



*Figure 2*. Path diagrams of a product of two variables. The asterisk surrounded by a circle represents the product of the variables connected to it by incoming arrows. (a) Mean-centered variables with a single product term. (b) Equivalent path diagram with variances sources isolated to be independent normally distributed variables with mean zero and variance one.

multiplication, just as one would expect algebraically. This is not explicit in some 1 diagrammatic representations for moderation where a moderating variable is singled out. 2 This might be perceived as just a philosophical difference. However, placing both 3 multiplicands of a commutative binary operation on equal visual footing clarifies that 4 moderation is just interaction without a direct effect of the moderating variable and 5 emphasizes that the SEM network of variables can be perturbed at either multiplicand 6 with an equivalent downstream effect. In addition, by diagramming the model in this way, 7 it sets up how the user interface for specifying SEM models that include products of 8 variables can be implemented as a minimal change from current SEM conventions. 9

1

## Variance of the Product of Two Variables

Given two vectors  $\mathbf{x}$  and  $\mathbf{y}$ , each with N elements, Goodman (1960) first derived the

<sup>3</sup> variance of the product of two variables and showed that this was an unbiased estimator. If <sup>4</sup>  $\Delta \mathbf{x} = \mathbf{x} - \mu(\mathbf{x})$  and  $\Delta \mathbf{y} = \mathbf{x} - \mu(\mathbf{x})$ , the exact variance is given by

$$\operatorname{var}(\mathbf{xy}) = \operatorname{cov}(\Delta \mathbf{x}^{2}, \Delta \mathbf{y}^{2}) + \operatorname{var}(\mathbf{x})\operatorname{var}(\mathbf{y}) + \mu(\mathbf{x})^{2}\operatorname{var}(\mathbf{y}) + \mu(\mathbf{y})^{2}\operatorname{var}(\mathbf{x}) + (\mu(\mathbf{x})\mu(\mathbf{y}))^{2} + 2\mu(\mathbf{x})\mu(\mathbf{y})\operatorname{cov}(\mathbf{x}, \mathbf{y}) + 2\mu(\mathbf{x})\operatorname{cov}(\Delta \mathbf{x}, \Delta \mathbf{y}^{2}) + 2\mu(\mathbf{y})\operatorname{cov}(\Delta \mathbf{x}^{2}, \Delta \mathbf{y}) - \operatorname{cov}(\mathbf{x}, \mathbf{y})^{2}.$$
(3)

<sup>5</sup> Bohrnstedt and Goldberger (1969) applied this result to the bivariate normal case and

<sup>6</sup> showed that the variance of the product, var(xy), can be estimated as (Bohrnstedt &

7 Marwell, 1978, Equation 15)

$$\operatorname{var}(\mathbf{x}\mathbf{y}) = \mu(\mathbf{x})^{2}\operatorname{var}(\mathbf{y}) + \mu(\mathbf{y})^{2}\operatorname{var}(\mathbf{x}) + 2\mu(\mathbf{x})\mu(\mathbf{y})\operatorname{cov}(\mathbf{x},\mathbf{y}) + \operatorname{var}(\mathbf{x})\operatorname{var}(\mathbf{y}) + \operatorname{cov}(\mathbf{x},\mathbf{y})^{2}.$$
(4)

Bohrnstedt and Marwell (1978) also derived reliability for this estimator when the observed
product variable included a normally distributed error term.

<sup>10</sup> Now consider the regression equation

$$z_i = b_1 x_i y_i + e_i \tag{5}$$

<sup>11</sup> where **x** and **y** are mean centered, multivariate normal, and the vector of residuals, **e**, is <sup>12</sup> mean centered, normally distributed, and independent of the predictor variables **x** and **y**. <sup>13</sup> This may be drawn as a path diagram as shown in Figure 2-a. A new set of tracing rules <sup>14</sup> will be required when this product symbol is used. However, note that the total variance at <sup>15</sup> the product symbol can be estimated from Equation 6. The downstream contribution of <sup>16</sup> this variance to **z** will be scaled by  $b_1^2$ . We can thus linearly decompose the variance of **z** <sup>17</sup> into terms due to the contribution of the error term **y** and the product of **x** and **y** such that

$$\operatorname{var}(\mathbf{z}) = b_1^2[\mu(\mathbf{x})^2 \operatorname{var}(\mathbf{y}) + \mu(\mathbf{y})^2 \operatorname{var}(\mathbf{x}) + 2\mu(\mathbf{x})\mu(\mathbf{y})\operatorname{cov}(\mathbf{x},\mathbf{y}) + \operatorname{var}(\mathbf{x})\operatorname{var}(\mathbf{y}) + \operatorname{cov}(\mathbf{x},\mathbf{y})^2] + \operatorname{var}(\mathbf{e}) .$$
(6)

<sup>1</sup> When  $\mu(\mathbf{x}) = \mu(\mathbf{y}) = 0$ , Equation 6 simplifies to

$$\operatorname{var}(\mathbf{z}) = b_1^2[\operatorname{var}(\mathbf{x})\operatorname{var}(\mathbf{y}) + \operatorname{cov}(\mathbf{x}, \mathbf{y})^2] + \operatorname{var}(\mathbf{e}) .$$
(7)

We next derive a general estimator for the variance decomposition of products of variables and apply it to the example problem in Equation 7 in order to demonstrate how the known estimator for the variance of the product of two variables is a special case result of the general estimator. This method relies on the model in the equivalent path diagram shown in Figure 2-b.

#### 7

## Expectations by Method of Moments for Products of Variables

<sup>8</sup> Suppose we are given an SEM as a path diagram in which all nodes with incoming <sup>9</sup> edges are additionally labeled as sum or product of their incoming edges. For the time <sup>10</sup> being, we will assume that the graph is acyclic (recursive). It may be that the results can <sup>11</sup> be extended to cases with cycles as well, provided that for all node values, an infinite series <sup>12</sup> calculated along each cycle converges.

We are interested in all moments of the vector of all observed variables. The following describes how those can be computed analytically assuming a fixed set of parameters. The main idea is to represent the variance sources of the SEM as independent

standard-normally distributed sources and then represent all variables as polynomial over these sources. Then, all moments become expectations of polynomials, which in turn are sums of monomials. In this way, the computation of the moments is reduced to computing the expectation of a monomial of independently standard-normally distributed variables—a problem that has a known computational solution.

21

At the highest level, the algorithm proceeds in three steps,

1. All variables in the SEM are represented by a linear combination of some

independently normally distributed variables  $w_1, ..., w_n$  with known variances such

that the covariance matrix of all variables is the symmetrical matrix  $\mathbf{S}$  from the RAM

<sup>25</sup> matrix formulation.

- 2. Progressing top-down in the asymmetrical graph of the path diagram, polynomial representations of all variables in the  $w_1, ..., w_n$  are computed.
- 3. Polynomial representations of all requested moments are computed and evaluated
   4 into numbers.

Suppose that we have an SEM represented in standard RAM notation such that the, the model-expected covariance matrix of the observed variables,  $\mathbf{C}_{xx}$  is calculated as

$$\mathbf{C}_{xx} = \mathbf{F}(\mathbf{I} - \mathbf{A})^{-1} \mathbf{S}((\mathbf{I} - \mathbf{A})^{-1})^T \mathbf{F}^T$$
(8)

<sup>5</sup> where for all variables, both latent and manifest, **A** is the matrix of regression coefficients, <sup>6</sup> **S** is the matrix of variances and covariances, and **I** is the identity matrix. The matrix **F** <sup>7</sup> filters out the latent variables so that  $C_{xx}$  contains only the model-expected covariance <sup>8</sup> matrix of the observed variables.

In order to transform the model into a model with only independent variables, we will operate on  $\mathbf{S}$ , the matrix of model variances and covariances both latent and manifest. We first compute the Eigenvalue decomposition of  $\mathbf{S}$ ,

$$\mathbf{S} = \mathbf{Q}\mathbf{D}\mathbf{Q}^T \tag{9}$$

Let  $W = w_1, ..., w_n$  be independently normally distributed variables with zero mean and variances given by the diagonal entries of **D**, the Eigenvalues of **S**, where *n* is the number of total variables in the SEM. Then **Q***W* is an *n*-dimensional random variable with zero mean and covariance

$$\mathbb{V}(\mathbf{Q}W) = \mathbb{E}(\mathbf{Q}\sqrt{\mathbf{D}}\sqrt{\mathbf{D}}\mathbf{Q}^T) = \mathbf{S}$$
(10)

<sup>9</sup> So each variable can be expressed by a linear combination of the W variables with the <sup>10</sup> corresponding row of Q as weights, plus a constant term that gives the mean of that <sup>11</sup> variable.

Using only product and sum operators, every variable in the model (both observed and latent)  $x_i$  can now be represented as a polynomial of the original observed variables. However, the mathematics of expectations are greatly simplified by using multivariate polynomials of *independent* (i.e., uncorrelated Gaussian) variables. Thus, we express each variable as a multivariate polynomial,  $f_i$ , in the independently normally distributed variables  $w_1, ..., w_n$ . If there are m variables in total, the  $k = (k_1, ..., k_m)$ -th moment of the joint distribution of the vector  $X = (x_1, ..., x_m)$  is

$$M_k(X) = \mathbb{E}(\prod_{i=1}^m f_i^{k_i}) \tag{11}$$

where the expectation is taken with respect to the roots  $w_i$ . In particular,  $M_k(X)$  is the expected value of a polynomial in  $w_i$ , where the coefficients are a combination of the regression weights in the SEM. Let this polynomial be  $g = \prod_{i=1}^{m} f_i^{k_i}$ . This polynomial is a sum of monomials in  $w_i$ . Since the expectation is linear, we can separate the computation of the expectation to the single monomials, and thus reduce our problem to computing the expectation of a monomial of independently normally distributed variables, i.e., an expectation of the form

$$\mathbb{E}\left(w_1^{e_1}\cdots w_n^{e_n}\right)\tag{12}$$

which is given by the product of the expectations of the single variables with their exponents,

$$\mathbb{E}\left(w_1^{e_1}\cdots w_n^{e_n}\right) = \prod_{i=1}^n \mathbb{E}(w_i^{e_i}) \tag{13}$$

For completeness sake, we give a quick proof:

Theorem 1.

1

$$\mathbb{E}\left(W_1^{e_1}\cdots W_n^{e_n}\right) = \prod_{i=1}^n \mathbb{E}(W_i^{e_i})$$
(14)

1

$$\mathbb{E} \left( W_1^{e_1} \cdots W_n^{e_n} \right) = \int_{-\infty}^{\infty} w_1^{e_1} \cdots w_n^{e_n} p df(w_1) \cdots p df(w_n) dw_1 \dots dw_n$$

$$= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} w_1^{e_1} p df(w_1) dw \right) w_2^{e_2} \cdots w_n^{e_n} p df(w_2) \cdots p df(w_n) dw_2 \dots dw_n$$

$$= \mathbb{E}(W_1^{e_1}) \int_{-\infty}^{\infty} w_2^{e_2} \cdots w_n^{e_n} p df(w_2) \cdots p df(w_n) dw_2 \dots dw_n$$

$$\vdots$$

$$= \prod_{i=1}^{n} \mathbb{E}(W_i^{e_i})$$

Thus, we are left with computing the higher-order moments of independent standard-normally distributed variables with zero mean and unit variance,

$$\mathbb{E}(W^e) , \qquad (15)$$

where e is a positive integer. These moments are known (e.g., Papoulis & Pillai, 2002) to be

$$\mathbb{E}(W^e) = \left\{ \begin{array}{ll} 0 & \text{if } e \text{ is odd} \\ \mathbb{V}(W)^{e/2} \prod_{i=1}^{\frac{e-2}{2}} (2i+1) & \text{if } e \text{ is even} \end{array} \right\}$$
(16)

## <sup>2</sup> Example 1: Bivariate Product of Variables Regression

As an example, we will transform the SEM model in Figure 2-a into the equivalent 3 SEM model shown in Figure 2-b such that all sources of variance and covariances are 4 independent, normally distributed variables with mean zero and unit variance. The first 5 step is to remove all covariances between variables by replacing them with unit variance 6 sources. Thus, we add the latent variable  $w_2$  and replace the covariance path with the 7 value  $C_{xy}$  in Figure 2-a with regression paths from  $w_2$  to x and y with values  $C_{xy}^{\frac{1}{2}}$ . Thus, 8 the covariance between x and y can be calculated as  $C_{xy}^{\frac{1}{2}} \cdot 1 \cdot C_{xy}^{\frac{1}{2}} = C_{xy}$ . However, now the 9 variances for x and y become residual variances that must be reduced by the total effect of 10  $w_2$  which is  $C_{xy}^{\frac{1}{2}} \cdot 1 \cdot C_{xy}^{\frac{1}{2}} = C_{xy}$ . So, the residual variance of x is  $V_x - C_{xy}$  and the residual 11 variance of y is  $V_y - C_{xy}$ . 12

We can now replace the variance terms in Figure 2-a with the independent normally 1 distributed unit variance variables  $w_1$ ,  $w_3$ , and  $w_4$  and regression paths to x, y, and e 2 respectively. The regression weights for these variables become the square root of the 3 residual variances for x, y, and e as shown in Figure 2-b. The *i*th variable in the path 4 diagram in Figure 2-b can now be represented as a polynomial  $f_i$  of what we will refer to as 5 root nodes, i.e., the independent, normally distributed variables  $w_1, \ldots, w_n$  with zero mean 6 and unit variance. This transformation of an SEM model into polynomials of root nodes 7 can be applied to any SEM that can be represented as a RAM model, including models 8 that have n-ary operators such as introduced here. We will call models that conform to 9 RAM conventions along with addition and product n-ary operators Products of Variables 10 (PoV) models. 11

Returning to the example in Figure 2-b, we can compute the moments, including all joint moments of any variable or pair of variables. Thus, we can decompose the variance of the product variable  $\mathbf{z} = b_1 \mathbf{x} \mathbf{y} + \mathbf{e}$  into components of variance and covariance of the other variables in the SEM as follows:

$$\mathbb{V}(z) = \mathbb{E}(zz) \tag{17}$$

$$= \mathbb{E}((b_1xy + e)(b_1xy + e)) \tag{18}$$

$$= \mathbb{E}(b_1^2 x^2 y^2 + 2b_1 x y e + e^2) .$$
(19)

<sup>16</sup> We next transform the variables x, y, and e into their monomial equivalents,

$$x = w_1 (V_x - C_{xy})^{\frac{1}{2}} + w_2 (C_{xy})^{\frac{1}{2}}$$
(20)

$$y = w_3 (V_y - C_{xy})^{\frac{1}{2}} + w_2 (C_{xy})^{\frac{1}{2}}$$
(21)

$$e = w_4 (V_e)^{\frac{1}{2}} . (22)$$

Next we substitute Equations 20, 21 and 22 into Equation 19 and then expand, collect, and
cancel terms to find that

$$\mathbb{V}(z) = b_1^2 (V_x V_y + C_{xy}^2) + V_e , \qquad (23)$$

<sup>1</sup> which is the result in Equation 7 derived from the results of Goodman (Goodman, 1960)

and Bohrnstedt and colleagues (Bohrnstedt & Marwell, 1978). The complete derivation of
this result is lengthy but is included in the supplemental material for this article.

Following the same logic, we can derive the expected covariance matrix of the model  $z = b_1 xy + e$  as

$$\mathbb{E}(\Sigma) = \begin{bmatrix} V_x & C_{xy} & 0 \\ C_{xy} & V_y & 0 \\ 0 & 0 & b_1^2 (V_x V_y + C_{xy}^2) + V_e \end{bmatrix}$$
(24)

<sup>6</sup> What becomes apparent here is that this model is unidentified, that is to say,  $b_1$  and  $V_e$ <sup>7</sup> only appear in one cell of the matrix and in that cell they form two parts of a sum. Thus, a <sup>8</sup> smaller  $b_1$  can be compensated by a larger  $V_e$  and vice versa. In addition, since  $b_1$  only <sup>9</sup> appears as a squared term, it is unidentified with respect to its sign. Thus the maximum <sup>10</sup> likelihood solution to this problem is not a single point, but lies on a line where all values <sup>11</sup> on that line are equally likely.

When pre-multiplying two variables and then adding them into a regression equation 12 as a way of estimating an interaction, one is taught to always subtract the mean of each 13 variable prior to multiplying them together. This preprocessing is done in order to remove 14 spurious covariances between the multiplicands and the outcome. But in the case of the 15 current PoV method of estimating the effects of products, we can account for these 16 covariances and they provide a way to resolve the underidentification of product coefficient 17 and error term. If we do not remove the means prior to entering the multiplicands into the 18 model and then estimate the means as part of a full information maximum likelihood 19 solution, the model becomes fully identified. 20

When the expected means and covariances are derived for the same model  $z_2 \quad z = b_1 x y + e$  we find 1

2

Again, the full derivation of this result can be found in the supplemental materials.

#### Simulations

We ran a series of simulations in order to test the performance of the implementation of the PoV method in OpenMx on five common use cases. The only new syntax required in scripting a PoV model is that one must add a line declaring **productVars=** with the names of any product nodes shown as a circle around an asterisk in the following path diagrams. Then one may use the standard **mxPath()** statements to specify the paths between named variables. All R scripts for each of the example models simulated and estimated below are included in the supplemental materials.

## <sup>10</sup> Simulation One: Manifest Variable Moderation

A moderation model with only observed variables as shown in Figure 3 was simulated where the parameters took on one of the values  $b_1 = \{-0.5, 0.5\}, b_2 = \{-0.5, 0.5\}, b_3 = \{-0.5, 0.5\}, \mu_x = \{-0.5, 0, 0.5\}, \mu_y = \{-0.5, 0, 0.5\}, C_{xy} = \{-1, 0, 1\}, V_e = \{0.2, 0.7, 1.2\},$  V<sub>x</sub> = V<sub>y</sub> = 3.0, resulting in 2 × 2 × 3 × 3 × 3 × 3 = 324 conditions. Each condition was
replicated 30 times, resulting in a total of 9,720 data sets. Prior to simulating the data for
N = 1000 simulated participants, a normally distributed random number (μ = 0, σ = 0.1)
was added to each parameter to better cover the parameter space.



*Figure* 3. Path diagram and simulation outcomes for a bivariate moderation model. For the scatter plots, the color of the plotted point corresponds to the error variance condition. In the plot for  $b_1$ , note that there are some estimated values that have the opposite sign than the generating values. This occurs when both the mean of x and y are very close to zero and the  $b_1$  parameter is underidentified.

The simulation resulted in 99.8% convergence. The estimated parameters from each 5 of the 9,720 data sets was compared to the generating coefficients for that data set. Mean 6 relative bias was calculated as the signed difference between the generating and estimated 7 coefficient divided by the generating coefficient. Mean relative bias for  $b_1$  was 0.046 and for 8  $b_2$  was < 0.001. Assuming an alpha level of 0.95, coverage as calculated by 1.96 time the 9 parameter standard error was 0.974 for  $b_1$  and was 0.958 for  $b_2$ . Thus, standard errors were 10 correct for the direct effect,  $b_2$ , but were conservative for the product effect,  $b_1$ . The model 11 was estimated on the same 9,720 data sets but with  $b_1$  fixed to zero and the likelihood ratio 12

test was calculated by subtracting the minus two log likelihood of the full model from the model with the parameter set to zero. Again assuming an alpha level of 0.95 and a  $\chi^2$  test with one degree of freedom, this resulted in a type I error rate of 0.0494. Thus, although the standard error for the product direct effect is conservative, the likelihood ratio test performs exactly as expected.

## <sup>6</sup> Simulation Two: Manifest Variable Moderation with Missing Values

One of the problems with current methods used to estimate moderation is that they are not able to account for missing data. In order to test whether the PoV method could account for missingness in the same way as full information maximum likelihood, we re-ran Simulation One where we substituted a random selection of 20% of the values of y and a separate random selection of the values of z with NA. Thus, this fits with the definition of values of y and z being missing completely at random. Figure 4 presents the results of this simulation.



Figure 4. Path diagram and simulation outcomes for a bivariate moderation model where 20% of y and z are missing completely at random.

This simulation resulted in 97.3% convergence. Mean relative bias for  $b_1$  was 0.015 and for  $b_2$  was < 0.001. Assuming an alpha level of 0.95, coverage as calculated by 1.96 <sup>1</sup> time the parameter standard error was 0.969 for  $b_1$  and was 0.950 for  $b_2$ . Clearly, full

<sup>2</sup> information maximum likelihood is working as expected when data are missing completely
<sup>3</sup> at random.

Missingness that can be accounted for by the model, the case of missing at random, 4 was also tested in simulation for this model. Here, missingness in y was set to be 5 conditional on the value of x. When x was more than one standard deviation above its 6 mean, y had a 50% chance of being missing. In addition, z had a 20% chance of being 7 missing, but not conditional on any part of the model. Thus y was missing at random and 8 z was missing completely at random. This simulation resulted in 97.4% convergence. Mean 9 relative bias for  $b_1$  was 0.049 and for  $b_2$  was 0.003. Assuming an alpha level of 0.95, 10 coverage as calculated by 1.96 time the parameter standard error was 0.811 for  $b_1$  and was 11 0.855 for  $b_2$ . Thus, the missing at random case increased bias by a very small amount, but 12 decreased coverage considerably in both the direct effect and the product effect. 13

## <sup>14</sup> Simulation Three: Latent Variable Moderation

A moderation model with a latent variable moderating a second latent variable and a 15 latent outcome as shown in Figure 5 was simulated where the parameters took on one of 16 the values  $b_1 = \{-0.5, 0, 0.5\}, b_6 = \{-0.5, 0, 0.5\}, \mu_x = \{-0.5, 0.5\}, \mu_y = \{-0.5,$ 17  $C_{xy} = \{-1, 0, 1\}$ , and  $V_e = \{0.1, 0.6, 1.1\}$ . The variance of the predictor latent variables was 18 set to  $V_x = V_y = 3.0$  and the residual variance of the outcome latent variable was set to 19  $V_z = V_e$ . The loadings for the latent variables were set to  $b_2 = b_{x2} = b_{y2} = \{-.05, 0.5\}$  and 20  $b_3 = b_{x3} = b_{y3} = \{-.05, 0.5\}$  and the unique variances were all set to the current value of 21  $V_e$ . This resulted in  $3 \times 3 \times 2 \times 2 \times 3 \times 3 \times 2 \times 2 = 1296$  conditions. Each condition was 22 replicated 8 times, resulting in a total of 10,368 data sets, each with N = 1,000 simulated 23 participants. Prior to simulating the data, a normally distributed random number ( $\mu = 0$ , 24  $\sigma = 0.1$ ) was added to each parameter to better cover the parameter space. 25

The simulation resulted in 99.9% convergence. Mean relative bias for  $b_1$  was -0.092



*Figure 5*. Path diagram and simulation outcomes for a bivariate latent moderation model.

and for  $b_6$  was < 0.001. Assuming an alpha level of 0.95, coverage as calculated by 1.96 1 time the parameter standard error was 0.975 for  $b_1$  and was 0.958 for  $b_6$ . Thus, standard 2 errors were correct for the direct effect,  $b_6$ , but were conservative for the product effect,  $b_1$ . 3 The model was estimated on the same 9,720 data sets but with  $b_1$  fixed to zero and the 4 likelihood ratio test was calculated by subtracting the minus two log likelihood of the full 5 model from the model with the parameter set to zero. Again assuming an alpha level of 6 0.95 and a  $\chi^2$  test with one degree of freedom, this resulted in a type I error rate of 0.0221. 7 Thus, both the standard error for the product direct effect and the likelihood ratio test 8 were conservative. This was likely due to the fact that in the latent simulation we allowed 9 the true product effect  $b_1$  to take on values very close to zero. 10

## <sup>11</sup> Simulation Four: Latent Variable Interaction

A model with a product of two latent variables and direct effects on a latent outcome as shown in Figure 6 was simulated where the parameters took on one of the values  $b_1 = \{-0.5, 0.5\}, b_6 = \{-0.5, 0.5\}, b_7 = \{-0.5, 0.5\}, \mu_x = \{-0.5, 0.5\}, \mu_y = \{-0.5, 0.5\}, \mu_y = \{-1, 0, 1\}, \text{ and } V_e = \{0.1, 0.6, 1.1\}.$  The variance of the predictor latent variables was set to  $V_x = V_y = 3.0$  and the residual variance of the outcome latent variable was set to  $V_z = V_e$ . The loadings for the latent variables were set to  $b_2 = b_{x2} = b_{y2} = \{-.05, 0.5\}$  and

b<sub>3</sub> = b<sub>x3</sub> = b<sub>y3</sub> = {-.05, 0.5} and the unique variances were all set to the current value of
V<sub>e</sub>. This resulted in 2 × 2 × 2 × 2 × 2 × 3 × 3 × 2 × 2 = 1, 152 conditions. Each condition
was replicated 9 times, resulting in a total of 10,368 data sets, each with N = 1,000
simulated participants. Prior to simulating the data, a normally distributed random
number (μ = 0, σ = 0.1) was added to each parameter to better cover the parameter space.



Figure 6. Path diagram and simulation outcomes for a bivariate latent interaction model. Again we see some sign reversals in  $b_1$  due to means of the latent variables being near zero.

<sup>6</sup> The simulation resulted in 98.2% convergence. Mean relative bias for  $b_1$  was 0.083, <sup>7</sup> for  $b_6$  was -0.001 and for  $b_7$  was < -0.013. Assuming an alpha level of 0.95, coverage as <sup>8</sup> calculated by 1.96 time the parameter standard error was 0.992 for  $b_1$ , 0.983 for  $b_6$  and <sup>9</sup> 0.984 for  $b_7$ . Thus, standard errors were conservative for the product effect,  $b_1$ , and the two <sup>10</sup> direct effects,  $b_6$  and  $b_7$ .

### <sup>11</sup> Simulation Five: Squared Predictor Variable

A model with a predictor variable, x and its square,  $x^2$  on a latent outcome as shown in Figure 7 was simulated where the parameters took on one of the values  $b_1 = \{-0.5, 0.5\}$ ,  $b_2 = \{-0.5, 0.5\}$ ,  $b_3 = \{-0.5, 0.5\}$ ,  $b_4 = \{-0.5, 0.5\}$ ,  $\mu_x = \{-1.0, -0.5, 0, 0.5, 1.0\}$ , and  $V_e = \{0.1, 0.35, 0.6, 0.95, 1.1\}$ . The variance of the predictor variable was set to  $V_x = 3.0$ and the residual variance of the outcome latent variable was set to  $V_z = V_e$ . The unique

<sup>1</sup> variances were all set to the current value of  $V_e$ . This resulted in

 $_2$  2 × 2 × 2 × 5 × 5 = 400 conditions. Each condition was replicated 25 times, resulting in

a total of 10,000 data sets, each with N = 1,000 simulated participants. Prior to

simulating the data, a normally distributed random number ( $\mu = 0, \sigma = 0.1$  was added to

<sup>5</sup> each parameter to better cover the parameter space.



Figure 7. Path diagram and simulation outcomes for a squared predictor.

<sup>6</sup> The simulation resulted in 99.9% convergence. Mean relative bias for  $b_1$  was 0.008 <sup>7</sup> and for  $b_4$  was 0.001. Assuming an alpha level of 0.95, coverage as calculated by 1.96 time <sup>8</sup> the parameter standard error was 0.999 for  $b_1$  and 0.972 for  $b_4$ . Thus, standard errors were <sup>9</sup> conservative for the product effect,  $b_1$ , and the direct effect,  $b_4$ .

10

## Discussion

Estimation of SEM models with Products of Variables is a general method that can provide unbiased estimates of parameters when predictor variables are normally distributed. Coverage of parameters of products is higher than expected when calculated using standard error estimates, but the likelihood ratio test appears to perform as expected. The primary advantage of the PoV method is in its generality. Most alternative methods for estimating parameters of SEM models containing products focus on special

cases and involve complicated constraints or estimation of parameters with non-normally
distributed variability. The PoV method is easy to use as it has been implemented in
OpenMx and Ωnyx. The method can handle interaction and moderator variables that are
either latent or have missingness. We expect that there are many more models other than
the obvious moderation and interaction models for which PoV estimation can be used.

## 6 Interaction versus Moderation

The use of the product calculation node in path diagrams has clarified for us the 7 difference between interaction between variables and moderation of one variable by 8 another. Examine the difference between Figure 5 and Figure 6. These two diagrams are 9 exactly the same other than the direct effect between y and z that appears in the 10 interaction diagram, Figure 6, but does not appear in the moderation diagram, Figure 5. 11 Thus, the only thing that distinguishes a moderating variable from the other multiplicand 12 is its lack of a direct effect on the outcome. It becomes clear that when one has a product 13 of two variables, say  $z = x \cdot y$  in a model, either x or y or both may be considered to be a 14 moderator depending on presence or absence of a direct effect of x or y on z. 15

The reader may wonder why we did not use a classic interaction model with direct 16 effects as an example model to illustrate the procedure. Figure 8 illustrates why this model 17 is unidentified using the current method. The reason why it is unidentified as a structural 18 model predicting covariances is apparent by counting seven free parameters while there are 19 only six degrees of freedom in a  $3 \times 3$  covariance matrix. One must watch for local 20 underidentification in models that include products of variables. The overall model may be 21 globally identified as determined by counting parameters and degrees of freedom in the 22 covariance matrix, but there still be areas of local underidentification resulting in sign 23 reversals such as those seen in simulations 1 and 4. 24

It may be more difficult to see why identification works when one premultiplies two manifest variables and creates a third variable for an interaction model. In this case there



*Figure 8*. A manifest variable interaction model with both direct effects has more parameters than statistics and thus is unidentified.

is a 4 × 4 covariance matrix and so there are not negative degrees of freedom. But why is
the 4 × 4 covariance matrix of full rank? This is due to the contributions of higher order
moments introduced by the multiplication formed at the individual data rows when the
premultiplication is performed. This model is identified because a product of two normal
distributions is not itself normal.

# <sup>6</sup> Product Calculation Nodes and the Identity Function

One feature of PoV is its reliance on a product node extension of RAM model 7 matrices. Although OpenMx has chosen to call this a ProductVar, it is not really a 8 variable. As such it does not have a mean or a variance, it is simply a node that indicates 9 the multiplication operation. Thus although a product node occupies a row and column in 10 the RAM matrices, OpenMx prevents the user from connecting a variance, covariance, or 11 mean path to this calculation node. However, if one examines the RAM matrices, it 12 becomes apparent that there is an automatically specified mean path to every product 13 node with a fixed value of 1.0. Why is that? 14

In RAM, removing a path and setting a path to zero are identical operations. 1 However, this is not the case in PoV. Note that the identity function for addition is to add 2 zero whereas the identity function for multiplication is multiply by one. If there is no 3 desired effect between predictor and outcome in a normal SEM, a regression path from the 4 predictor to the outcome is set to zero. However, if a regression path pointing from one 5 multiplicand to a multiplication operator symbol is set to zero in an PoV diagram, then the 6 other multiplicand is multiplied by zero. This may not be the intention of the modeler. 7 When one wishes to use a likelihood ratio to test the difference between a model with and 8 without a product, we recommend to use the outgoing path from the product node as the 9 path to set to zero for the one degree of freedom minus two log likelihood comparison. 10

### **11** Assumptions and Limitations

The algorithm in this article makes the strong assumption that the variables in an 12 PoV allow the transformation of the model into an equivalent model in which all variance 13 sources are represented as linear combinations of independent normally distributed 14 standardized variables. This is a relaxation of the usual SEM assumption of multivariate 15 normality. The outcome of a product of variables will not be normally distributed. 16 However, the PoV assumption is that the non-normality of a product outcome variable can 17 be completely accounted by the non-normality generated by multiplying multivariate 18 normal variance sources. Thus, all variables that are not outcome variables must be 19 multivariate normal. To be explicit, this means that all residuals must be normally 20 distributed, including the residuals of variables that are outcomes of products of variables. 21

This article presents a novel addition to the already complex infrastructure of structural equation modeling. We acknowledge that there may be special cases that need to be tested before products of variables models can be recommended in those cases. We recommend using simulation and either OpenMx or Ωnyx to verify models that include PoV structures that go beyond the common use cases presented here. For instance, since the PoV relies on normality of predictor variables, ordinal predictors might not be able to be used as a multiplicand. If one has ordinal manifest variables that one wishes to use in a product, we recommend a method such as OpenMx definition variables (Neale, 1998; Neale et al., 2016). In principal, the method for estimating a normally distributed latent variable from ordinal indicators as implemented in OpenMx should be able to be used to create a latent multiplicand, but we have not yet fully tested this case.

Conclusions

Products of variables can be included in structural equation models and model expectations can be calculated using the methods introduced here. This means that products of combinations of manifest and latent variables are now possible to implement in standard SEM software in an automated and easy-to-specify way. Three important and commonly used models that will immediately benefit from PoV estimation are: i) models with interactions between latent variables, ii) latent moderator models, and iii) moderator models with missingness in the moderator variable.

Future work will certainly find novel and interesting uses for the PoV method. One 16 case that we have begun to explore is the use of PoV to extract a "factor of paths", 17 creating a latent variable that accounts for the commonality between coefficients that 18 include person-specific variance. This could be construed to be a special case of multilevel 19 models with random coefficients and thus we may find that PoV estimation may prove 20 useful for multilevel models. A second case that we have begun to explore is nonlinear 21 dynamical systems models. Differential equations with squared and cubic terms are at the 22 heart of dynamical systems that exhibit bifurcation and/or chaotic dynamics. We expect 23 that PoV estimation will prove useful in fitting these sorts of models to complex 24 physiological and behavioral timeseries. 25

26

8

We believe that the algorithm presented here represents a paradigm shift for

representing theories within the context of structural equation models. We look forward to
seeing PoV estimation appearing as a feature in software other than OpenMx and Ωnyx
and with that in mind, we refer SEM software authors to our open source code available on
GitHub.

### References

- Bauer, D. J. (2005, oct). A semiparametric approach to modeling nonlinear relations among latent variables. Structural Equation Modeling, 12(4), 513–535.
- Bauer, D. J., & Hussong, A. M. (2009). Psychometric approaches for developing commensurate measures across independent studies: traditional and new models. *Psychological methods*, 14(2), 101.
- Bohrnstedt, G. W., & Goldberger, A. S. (1969). On the exact covariance of products of random variables. Journal of the American Statistical Association, 64 (328), 1439–1442.
- Bohrnstedt, G. W., & Marwell, G. (1978). The reliability of products of two random variables. Sociological Methodology, 9, 254–273.
- Boker, S. M., & McArdle, J. J. (2005). Path analysis and path diagrams. In B. Everitt &
  D. Howell (Eds.), *Encyclopedia of statistics in behavioral science (vol. 3)* (pp. 1529–1531). New York: John Wiley & Sons.
- Boker, S. M., McArdle, J. J., & Neale, M. C. (2002). An algorithm for the hierarchical organization of path diagrams and calculation of components of covariance between variables. *Structural Equation Modeling*, 9(2), 174–194.
- Goodman, L. A. (1960). On the exact variance of products. Journal of the American Statistical Association, 55(292), 708–713.
- Kenny, D. A., & Judd, C. M. (1984). Estimating the nonlinear and interactive effects of latent variables. *Psychological Bulletin*, 96(1), 201–210.
- Klein, A., & Moosbrugger, H. (2000). Maximum likelihood estimation of latent interaction effects with the lms method. *Psychometrika*, 65(4), 457–474.
- Marsh, H. W., Wen, Z., & Hau, K.-T. (2004). Structural equation models of latent interactions: Evaluation of alternative estimation strategies and indicator construction. *Psychological Methods*, 9(3), 275.
- Marsh, H. W., Wen, Z., Hau, K.-T., Little, T. D., Bovaird, J. A., & Widaman, K. F.

(2007). Unconstrained structural equation models of latent interactions: Contrasting residual-and mean-centered approaches. Structural Equation Modeling: A Multidisciplinary Journal, 14 (4), 570–580.

- McArdle, J. J., & McDonald, R. P. (1984). Some algebraic properties of the Reticular Action Model for moment structures. British Journal of Mathematical and Statistical Psychology, 37, 234–251.
- Mehta, P. D., & Neale, M. C. (2005). People are variables too: Multilevel structural equations modeling. *Psychological Methods*, 10(3), 259–284.
- Moosbrugger, H., Schermelleh-Engel, K., & Klein, A. (1997). Methodological problems of estimating latent interaction effects. Methods of Psychological Research Online, 2(2), 95–111.
- Neale, M. C. (1998). Modeling interaction and nonlinear effects with mx: A general approach. In G. Marcoulides & R. Schumacker (Eds.), *Interaction and non-linear effects in structural equation modeling* (pp. 43–61). Mahwah, NJ: Lawrence Erlbaum Associates.
- Neale, M. C., Hunter, M. D., Pritikin, J. N., Zahery, M., Brick, T. R., Kirkpatrick, R., et al. (2016). Openmx 2.0: Extended structural equation and statistical modeling. *Psychometrika*, 81(2), 535–549. (PMCID: 25622929)
- Oertzen, T. von, Brandmaier, A., & Tsang, S. (2015). Structural equation modeling with Ωnyx. Structural Equation Modeling: A Multidisciplinary Journal, 22(1), 148–161.
- Papoulis, A., & Pillai, S. U. (2002). Probability, random variables, and stochastic processes. New York: Tata McGraw-Hill Education.
- Schumacker, R. E. (2002). Latent variable interaction modeling. Structural Equation Modeling, 9(1), 40–54.
- Wall, M. M., & Amemiya, Y. (2001). Generalized appended product indicator procedure for nonlinear structural equation analysis. *Journal of Educational and Behavioral Statistics*, 26(1), 1–29.

- Wall, M. M., & Amemiya, Y. (2003). A method of moments technique for fitting interaction effects in structural equation models. British Journal of Mathematical and Statistical Psychology, 56, 47–63.
- Wright, S. (1921). Correlation and causation. Journal of agricultural research, 20(7), 557–585.
- Wright, S. (1934). The method of path coefficients. The Annals of Mathematical Statistics, 5, 161–215.