

3.3 | Reconstruction of the Covariance Matrix

After representing the regression coefficients and variances depicted in a path diagram using a series of matrices, we can develop a system of matrix equations with which to solve for those parameters that were not fixed in the diagram. This solution, much like that used in exploratory factor analysis, is meant to provide the best possible reconstruction of the covariability between the original dimensions. But whereas factor analysis focuses on the correlation structure of the original data matrix, structural equation modeling instead focuses on the unstandardized covariance structure of the data.

3.3.1 | Replacing NAs with Starting Values

In the Bentler-Weeks approach to structural equation modeling, every dimension, factor, error, and dimension serves as either an independent or a dependent variable. The paths connecting these variables constitute the parameters to be estimated, which can be represented using the two matrices of regression coefficients and the one matrix of variances.

We can now use these matrices to express the structural relationships between and among the dependent and independent variables using matrix algebra:

$$\boldsymbol{\eta} = \boldsymbol{\beta}\boldsymbol{\eta} + \boldsymbol{\gamma}\boldsymbol{\zeta} \quad (3.5)$$

In this equation, $\boldsymbol{\eta}$ is a $q \times 1$ matrix of dependent variables and $\boldsymbol{\zeta}$ is an $r \times 1$ matrix of independent variables.

We can expand this equation into what is commonly referred to as the *Bentler-Weeks structural equation model*:

$$\begin{bmatrix} \eta_1 \\ \vdots \\ \eta_q \end{bmatrix} = \begin{bmatrix} 0 & \dots & * \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \times \begin{bmatrix} \eta_1 \\ \vdots \\ \eta_q \end{bmatrix} + \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ * & \dots & 0 \end{bmatrix} \times \begin{bmatrix} \zeta_1 \\ \vdots \\ \zeta_r \end{bmatrix} \quad (3.6)$$

Although it may seem strange that the matrix of dependent variables $\boldsymbol{\eta}$ appears on both sides of the equation, it is important to remember that in structural equations, dependent variables are themselves able to directly affect other variables, allowing them to appear on both sides.

This equation, which is very similar to the $\mathbf{Y} = \boldsymbol{\beta}\mathbf{X} + \boldsymbol{\varepsilon}$ equation used in linear regression, forms the basis of parameter estimation in structural equation modeling. However, because the goal of these parameter estimates is to capture as much of the covariance structure among the original measured dimensions as possible, this equation is first converted into a series of matrix equations that reconstruct the covariance matrix from these parameter estimates.

Before moving on to the equations themselves, we need to consider again the missing values in each of our matrices, which represent parameters that still need to be estimated. The goal is to find the best set of values to replace these **NA** values such that the resulting reconstructed covariance matrix is as close to the original as possible. As we will discuss in more detail below and again in Section A.3, this makes structural equation modeling an optimization problem and, as we've seen before, such problems require a starting point for each parameter.

In R, we can manually assign values to replace these missing values using matrix notation.

```
BT["MilksCows", "workswithCows"] <- 0.1
BT["BirthsCows", "workswithCows"] <- 0.2
BT["OpticalDensity", "workswithCows"] <- 0.3
BT["AgeBasedEmployment", "workswithCows"] <- 0.4
BT["worksGardens", "AgeBasedEmployment"] <- 0.5
BT["CleansBarns", "AgeBasedEmployment"] <- 0.6
BT["OpticalDensity", "AgeBasedEmployment"] <- 0.7
BT["workswithCows", "AgeBasedEmployment"] <- 0.4
BT
```

	ButchersCows	MilksCows	BirthsCows	Income
ButchersCows	0	0	0	0
MilksCows	0	0	0	0
BirthsCows	0	0	0	0
Income	0	0	0	0
DrainsLands	0	0	0	0
Age	0	0	0	0
WorksGarden	0	0	0	0
CleansBarns	0	0	0	0
OpticalDensity	0	0	0	0
WorkswithCows	0	0	0	0
AgeBasedEmployment	0	0	0	0

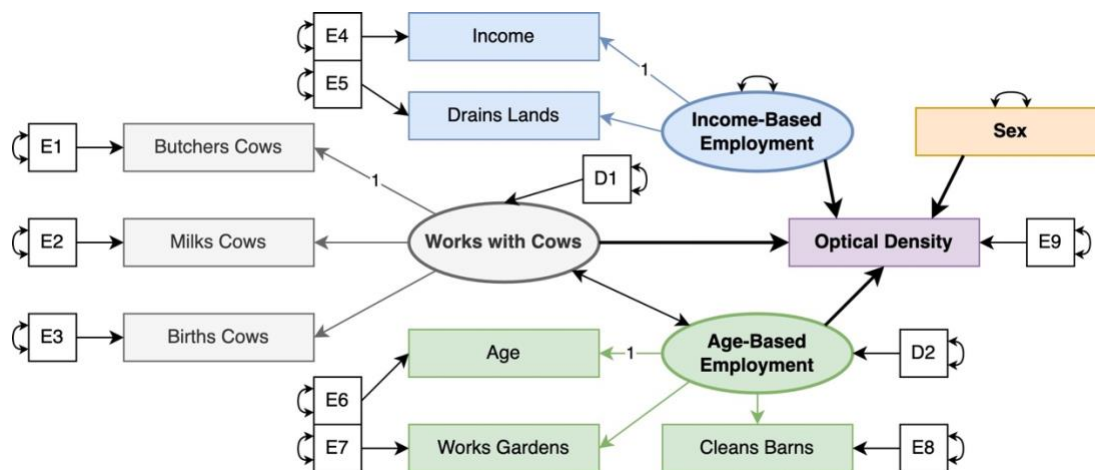
	DrainsLands	Age	WorksGardens	CleansBarns
ButchersCows	0	0	0	0
MilksCows	0	0	0	0
BirthsCows	0	0	0	0
Income	0	0	0	0
DrainsLands	0	0	0	0
Age	0	0	0	0
WorksGardens	0	0	0	0
CleansBarns	0	0	0	0
OpticalDensity	0	0	0	0
WorkswithCows	0	0	0	0
AgeBasedEmployment	0	0	0	0

	OpticalDensity	WorkswithCows	AgeBasedEmployment
ButchersCows	0	1.0	0.0
MilksCows	0	0.1	0.0
BirthsCows	0	0.2	0.0
Income	0	0.0	0.0
DrainsLands	0	0.0	0.0
Age	0	0.0	1.0
WorksGarden	0	0.0	0.5
CleansBarns	0	0.0	0.6
OpticalDensity	0	0.3	0.7
WorkswithCows	0	0.0	0.4
AgeBasedEmployment	0	0.4	0.0

Note that we have used row and column names in place of row and column numbers to simplify referencing specific elements. Because elements in a matrix are references as rows first and columns second, this naming notation can be read as “to variable one from variable two.”

Although any starting values (including identical ones) can be used to populate the NA values in this and other matrices, it is common to use slightly different values for each unique parameter. For elements that represent the covariability between complex dimensions or oblique factors and are denoted using paths with two arrows, however, the starting values should be identical.

For instance, in our example path diagram (recreated below), there is a single path between AgeBasedEmployment and WorksWithCows, indicating the presence of an oblique factor relationship. When this path is represented in a matrix, however, it can only be done so with two elements in the matrix – one from AgeBasedEmployment to WorksWithCows and another from WorksWithCows to AgeBasedEmployment. However, because these two elements represent the same path, the value should be the same between them, both in terms of the starting values and in terms of the eventually optimized parameter estimates.



We can use a similar approach to replace the NA values in the matrix of regression coefficients between the independent and dependent variables.

```
GM["OpticalDensity","Sex"] <- 2.1
GM["DrainsLands","IncomeBasedEmployment"] <- 2.2
GM["OpticalDensity","IncomeBasedEmployment"] <- 2.3
GM["ButchersCows","E1"] <- 2.4
GM["MilksCows","E2"] <- 2.5
GM["BirthsCows","E3"] <- 2.6
GM["Income","E4"] <- 2.7
GM["DrainsLands","E5"] <- 2.8
GM["Age","E6"] <- 2.9
GM["worksGardens","E7"] <- 3.0
GM["CleansBarns","E8"] <- 3.1
GM["OpticalDensity","E9"] <- 3.2
GM["workswithCows","D1"] <- 3.3
GM["AgeBasedEmployment","D2"] <- 3.4
```

Although we can take a similar approach to individually replace each of the NA values in the matrix of variances, the fact that all the missing values are, by definition, along the primary diagonal of the matrix, we can use the `diag()` function in R to accomplish this more efficiently.

```
diag(PH) <- c(4.1,4.2,4.3,4.4,4.5,4.6,4.7,4.8,4.9,5.0,5.1,5.2,5.2)
```

This is another reason to maintain the same order between errors and disturbances and their corresponding dimensions and factors – it allows us to use a single command to replace all the missing values in this matrix!

3.3.2 | Reconstructing the Covariance Matrix

Now that we have replaced all the missing values with starting points, we can use the matrices of regression coefficients and variances to reconstruct the covariance matrix for the original measured dimensions. Although it is certainly possible to do so in a single equation, it is much simpler to use a series of four equations, each of which is meant to reconstruct one specific portion of the covariance matrix.

Consider again our example covariance matrix. If we arrange the corresponding dimensions such that the dependent variables precede the independent variables (referring to the former as Y and to the latter as X , as is commonly done in linear regression), we can separate this matrix into four quadrants representing the covariances among the dependent variables, between the independent and dependent variables, and among the independent variables:

$$\Sigma = \begin{matrix} & \begin{matrix} Y_1 & \cdots & Y_9 \end{matrix} & \begin{matrix} X_1 \end{matrix} \\ \begin{matrix} Y_1 \\ \vdots \\ Y_9 \end{matrix} & \left[\begin{array}{ccc|c} & & & \\ & YY & & YX \\ & & & \\ - & - & - & + \\ & XY & & \\ & & & XX \end{array} \right] & \end{matrix} \quad (3.7)$$

Following the construction of our path diagram, our example data had nine different dependent dimensions but only a single independent dimension that measured the participant's Sex.

Note that only measured dimensions were included in this covariance matrix, whereas our matrices of regression coefficients and variances included not only measured dimensions, but also factors, errors, and disturbances.

As such, our first step in using these parameter matrices to reconstruct the covariance matrix is to define two *selection matrices* that will isolate only the measured dimensions from the full lists of dependent and independent variables.

We construct the *selection matrix for the dependent variables* G_y by listing every dependent variable as a column but only using the measured dimensions among these dependent variables as the rows. We then populate the primary diagonal of this matrix with ones to indicate their selection from among the other variables.

In R, we can manually create a matrix of zeroes using the `as.matrix()` function and the `diag()` function to populate the primary diagonal of this matrix.

```
Gy <- matrix(0, nrow = 9, ncol = 11,
             dimnames = list(c("ButchersCows", "MilksCows",
                               "BirthsCows", "Income", "DrainsLands",
                               "Age", "WorksGardens", "CleansBarns",
                               "OpticalDensity"),
                             c("ButchersCows", "MilksCows",
                               "BirthsCows", "Income", "DrainsLands",
                               "Age", "WorksGardens", "CleansBarns",
                               "OpticalDensity", "WorkswithCows",
                               "AgeBasedEmployment"))))
diag(Gy) <- 1
```

We take a similar approach to create the *selection matrix for the independent variables* G_x in R.

```
Gx <- matrix(0, nrow = 1, ncol = 13,
             dimnames = list(c("Sex"),
                             c("Sex", "IncomeBasedEmployment",
                               "E1", "E2", "E3", "E4", "E5", "E6",
                               "E7", "E8", "E9", "D1", "D2"))))
diag(Gx) <- 1
```

Now that we have a way of isolating the measured dimensions from the full sets of dependent and independent variables, we can now set up the series of equations to reconstruct each of the four quadrants in the covariance matrix.

We begin with the upper left quadrant of the reconstructed covariance matrix $\hat{\Sigma}_{YY}$, which represents the covariance among the dimensions that serve as dependent variables in the structural equation model:

$$\hat{\Sigma}_{YY} = G_Y \times \left((I - \tilde{\beta})^{-1} \right) \times \tilde{\gamma} \times \tilde{\Phi} \times \tilde{\gamma}^T \times \left((I - \tilde{\beta})^{-1} \right)^T \times G_Y^T \quad (3.8)$$

The *tilde* above each of the parameter matrices is used to indicate that the missing values therein have been replaced with starting values for the optimization. The *identity matrix* I is a matrix of zeroes with ones along the primary diagonal that is the same size as the matrix of regression coefficients among the dependent variables β .

We can use matrix algebra in R to implement this equation, first using the `diag()` function to quickly create the needed identity matrix.

```

I <- diag(nrow(BT))
SIG_yy <- Gy %%% solve(I - BT) %%% GM %%% PH %%%
          t(GM) %%% t(solve(I - BT)) %%% t(Gy)
SIG_yy

```

	ButchersCows	MilksCows	BirthsCows	Income	DrainsLands
ButchersCows	118.653941	9.388594	18.777188	0.000	0.000
MilksCows	9.388594	28.438859	1.877719	0.000	0.000
BirthsCows	18.777188	1.877719	34.175438	0.000	0.000
Income	0.000000	0.000000	0.000000	37.734	9.240
DrainsLands	0.000000	0.000000	0.000000	9.240	57.176
Age	66.179138	6.617914	13.235828	0.000	0.000
WorksGarden	33.089569	3.308957	6.617914	0.000	0.000
CleansBarns	39.707483	3.970748	7.941497	0.000	0.000
OpticalDensity	74.491179	7.449118	14.898236	9.660	21.252

	Age	worksGardens	CleansBarns	OpticalDensity
ButchersCows	66.179138	33.089569	39.707483	74.491179
MilksCows	6.617914	3.308957	3.970748	7.449118
BirthsCows	13.235828	6.617914	7.941497	14.898236
Income	0.000000	0.000000	0.000000	9.660000
DrainsLands	0.000000	0.000000	0.000000	21.252000
Age	138.401560	49.016780	58.820136	88.477234
WorksGardens	49.016780	68.608390	29.410068	44.238617
CleansBarns	58.820136	29.410068	83.342082	53.086340
OpticalDensity	88.477234	44.238617	53.086340	176.804417

Next is the upper right quadrant of the reconstructed covariance matrix $\hat{\Sigma}_{YX}$, which represents the covariance between the dimensions that serve as dependent variables and those that serve as independent variables in the structural equation model:

$$\hat{\Sigma}_{YX} = G_Y \times \left((I - \tilde{\beta})^{-1} \right) \times \tilde{\gamma} \times \tilde{\Phi} \times G_X^T \quad (3.9)$$

Although there is a similar equation for finding the lower left quadrant of the reconstructed covariance matrix $\hat{\Sigma}_{XY}$, it is important to remember that covariance matrices are, by definition, symmetric. As such, we can simply take the transpose of the upper right quadrant to find the lower left quadrant:

$$\hat{\Sigma}_{XY} = (\hat{\Sigma}_{YX})^T \quad (3.10)$$

The last piece is the lower right quadrant of the reconstructed covariance matrix $\hat{\Sigma}_{XX}$, which represents the covariance among the dimensions that serve as independent variables:

$$\hat{\Sigma}_{XX} = G_X \times \tilde{\Phi} \times G_X^T \quad (3.11)$$

We can again use matrix algebra to implement each of these three equations in R.

```

SIG_yx <- Gy %%% solve(I - BT) %%% GM %%% PH %%% t(Gx)
SIG_xy <- t(SIG_yx)
SIG_xx <- Gx %%% PH %%% t(Gx)

```

Now that we have reconstructed each of these quadrants, we can use the `cbind()` and `rbind()` functions in R to bind these quadrants together into the reconstructed covariance matrix $\hat{\Sigma}$.

```
SIG_hat <- rbind(cbind(SIG_yy, SIG_yx),
                 cbind(SIG_xy, SIG_xx))
SIG_hat
```

	ButchersCows	MilksCows	BirthsCows	Income	DrainsLands
ButchersCows	118.653941	9.388594	18.777188	0.000	0.000
MilksCows	9.388594	28.438859	1.877719	0.000	0.000
BirthsCows	18.777188	1.877719	34.175438	0.000	0.000
Income	0.000000	0.000000	0.000000	37.734	9.240
DrainsLands	0.000000	0.000000	0.000000	9.240	57.176
Age	66.179138	6.617914	13.235828	0.000	0.000
WorksGardens	33.089569	3.308957	6.617914	0.000	0.000
CleansBarns	39.707483	3.970748	7.941497	0.000	0.000
OpticalDensity	74.491179	7.449118	14.898236	9.660	21.252
Sex	0.000000	0.000000	0.000000	0.000	0.000

	Age	WorksGardens	CleansBarns	OpticalDensity
ButchersCows	66.179138	33.089569	39.707483	74.491179
MilksCows	6.617914	3.308957	3.970748	7.449118
BirthsCows	13.235828	6.617914	7.941497	14.898236
Income	0.000000	0.000000	0.000000	9.660000
DrainsLands	0.000000	0.000000	0.000000	21.252000
Age	138.401560	49.016780	58.820136	88.477234
WorksGardens	49.016780	68.608390	29.410068	44.238617
CleansBarns	58.820136	29.410068	83.342082	53.086340
OpticalDensity	88.477234	44.238617	53.086340	176.804417
Sex	0.000000	0.000000	0.000000	8.610000

	Sex
ButchersCows	0.00
MilksCows	0.00
BirthsCows	0.00
Income	0.00
DrainsLands	0.00
Age	0.00
WorksGardens	0.00
CleansBarns	0.00
OpticalDensity	8.61
Sex	4.10

Given that the goal of structural equation models is to find the set of parameter estimates that best recreates the original covariance structure, the next step is to compare this reconstructed covariance matrix $\hat{\Sigma}$ to the original observed covariance matrix Σ . To that end, we require a means of quantifying the differences between these two matrices.

Several different methods have been proposed to quantify these differences, including simply calculating a matrix of residuals in much the same way as we did with the reconstructed correlation matrix in exploratory factor analysis. The most common approach, however, is based on the *maximum likelihood* principle.

Under this principle, we define an objective function Q that quantifies the differences between the reconstructed and observed covariance matrices, with the eventual goal (as we will see in

Sections 3.4 and A.3) of finding the set of parameter estimates that best minimizes this difference:

$$Q = \log|\hat{\Sigma}| - \log|\Sigma| + \text{tr}(\Sigma \times \hat{\Sigma}^{-1}) - d \quad (3.12)$$

Recall that $|\Sigma|$ represent the determinant of the observed covariance matrix, the $\text{tr}()$ operator refers to the trace of a matrix, and d refers to the dimensionality of the original data matrix itself.

The last thing we need to do before implementing this comparison in R is to ensure that the order in which the dimensions appear in the observed covariance matrix is the same as the order in which they appear in the reconstructed one. Because the order in the reconstructed covariance matrix is based primarily on whether a dimension serves as a dependent or an independent variable in the structural equation and path diagram, this order is very often not the same as that of the original data matrix.

Thankfully, we can use simple matrix notation in R to rearrange both the rows and columns in the original covariance matrix such that they align with the order in the reconstructed one, using either the name or number of each column.

```
SIG <- SIG[rownames(SIG_hat), colnames(SIG_hat)]
```

Here we note an additional benefit of taking the time and effort to properly name the rows and columns in each of the matrices of regression coefficients and variances, as well as the selection matrixes – it makes rearranging the order of the original covariance matrix much faster!

Now we are ready to calculate our starting value of the objective function Q .

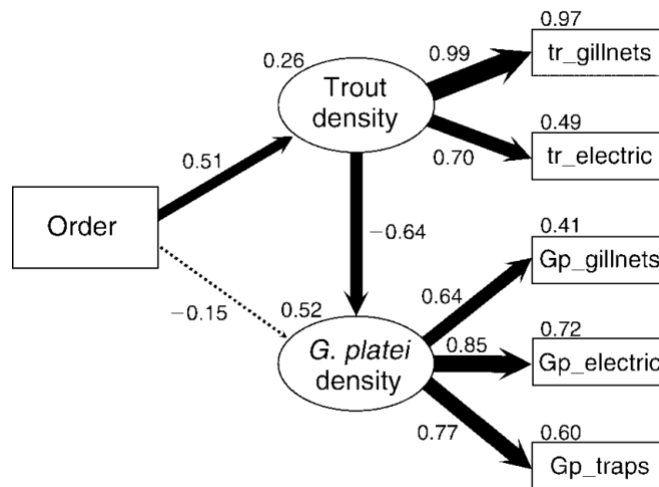
```
Q <- log(det(SIG_hat)) - log(det(SIG)) +  
      sum(diag(SIG %*% solve(SIG_hat))) - nrow(SIG)  
Q  
[1] 27.31325
```

The best set of estimates for each of the unfixed parameters in the matrices of regression coefficients and variances is the one that is able to produce the absolute smallest value of this objective function Q . How to find and then interpret the values in this best set of estimates is the main topic of Sections 3.4 and A.3.

The last thing to note is that much like the sum of squared residuals in linear regression or the objective function in orthogonal rotations, the actual value of this objective function is almost never interpreted or otherwise used. Instead, only the difference between values derived using different sets of parameter estimates is of interest, wherein a decrease indicates that the new set of estimates provides a more accurate reconstruction of the covariance structure.

Exercises for Section 3.3

3.09 Invasive Salmonids, Lake Order, and Galaxiid Decline. Salmonid fishes such as trout, which are native to the northern hemisphere, have become naturalized in many southern countries and appear linked to the decline of native fishes such as galaxiids. An article published in 2012 in *Ecological Applications* investigated the effects of hydrological position (i.e., lake order) and trout density on the density of *Galaxias platei* using a structural equation model.



The path diagram above represents an optimized set of path parameters based on a sample of 25 lakes in the Aysén region of Chile that were surveyed in 2007 and 2009. Note that the number next to each construct represents the path of the corresponding error or disturbance.

- Use the parameter estimates in the path diagram to construct the matrix of regression coefficients among the dependent variables in R with the proper row and column names
- Construct the matrix of regression coefficients between the independent and the dependent variables in R with the proper row and column names
- Construct the matrix of variances among the independent variables in R, assuming that all variances have been fixed to a value of one, with the proper row and column names
- Use matrix algebra to reconstruct the upper left quadrant of the covariance matrix (you will first need to construct the selection matrix for the dependent variables)
- Use matrix algebra to reconstruct the upper right, lower left, and lower right quadrants of the covariance matrix (you will first need to construct the selection matrix for the independent variables)
- Use the `cbind()` and `rbind()` functions to find the full reconstructed covariance matrix

3.10 Cavity Trees, Den Order, and Fishers (continued). Exercise 3.02 introduced additional data on the characteristics of eight cavity trees used as dens by fishers in Minnesota that were published in 2020 in the *Canadian Journal of Forest Research*, which are recreated below:

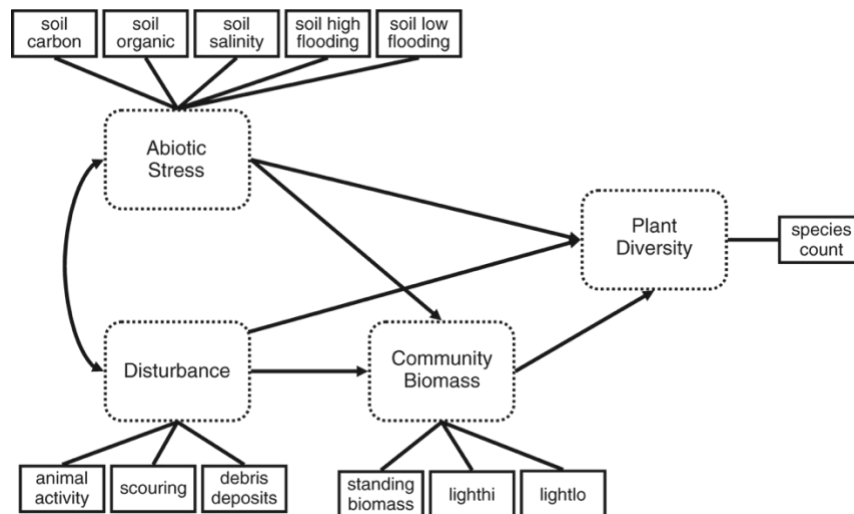
	DBH	Type	Slope	Aspect	Order
DBH	78.45	-2.96	-4.73	-183.75	-1.96

Type	-2.96	0.21	0.62	12.30	0.07
Slope	-4.73	0.62	34.23	-147.01	2.07
Aspect	-183.75	12.30	-147.01	2942.23	-12.38
Order	-1.96	0.07	2.07	-12.38	0.21

The covariance matrix above was calculated using the **DBH** (diameter at breast height in cm) and **Type** (0 for coniferous and 1 for deciduous) of the tree itself, as well as the **Slope** (in degrees from horizontal) and **Aspect** (in degrees from due South) of the ground the tree is on and the **Order** in which it was used by fishers at the Camp Ripley military facility in Minnesota.

- Recreate the matrix of regression coefficients among the dependent variables from Exercise 3.06, using starting values of 0.5 for every unfixed parameter
- Recreate the matrix of regression coefficients between the independent and dependent variables from Exercise 3.06, using starting values of 0.7 for every unfixed parameter
- Recreate the matrix of variances among the independent variables from Exercise 3.06, using starting values of 1.1 for every unfixed parameter
- Construct the selection matrix for the dependent variables
- Use matrix algebra to reconstruct the upper left quadrant of the covariance matrix
- Explain why there is no need to reconstruct any other quadrant of the covariance matrix

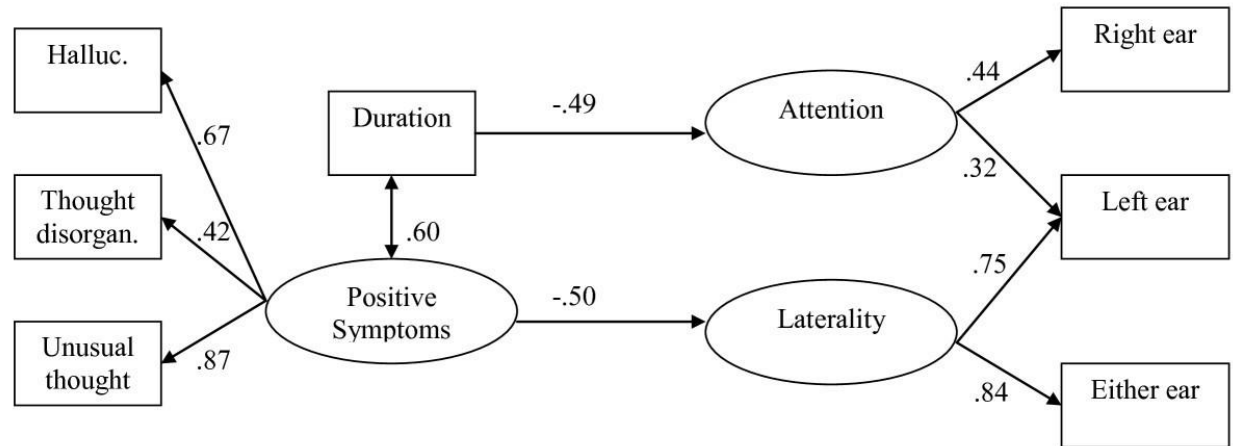
3.11 Plant Diversity and Environmental Factors (continued). Exercise 3.03 introduced data on how structural equation modeling can be used to support robust generalizations that were published in 2010 in *Ecological Monographs*, which are recreated below:



The path diagram above represents major categories of influences on spatial variations in plant diversity based on a theoretical framework developed in 1997.

- Use this path diagram to explain which dimensions (if any) would appear in the upper left quadrant of the reconstructed covariance matrix
- Explain which dimensions would appear in the upper right quadrant of this matrix
- Explain which dimensions would appear in the lower left quadrant of this matrix
- Explain which dimensions would appear in the lower right quadrant of this matrix

3.12 Schizophrenia and Language Processing (continued). Exercise 3.04 introduced data on impairment in left temporal lobe language processing and positive symptoms of schizophrenia that were first published in 2010 in the *BMC Research Notes*, which are recreated below:



The path diagram above includes the hypothesized relationships among positive symptoms, duration of schizophrenia, and dichotic listening, using a sample of 129 patients from clinics in Norway and California who were diagnosed with schizophrenia.

	1	2	3	4
1. Hallucinations	4.49	0.92	2.66	7.24
2. Thought disorganization	0.92	2.22	2.61	4.55
3. Unusual thoughts	2.66	2.61	4.37	9.58
4. Duration	7.24	4.55	9.58	80.82

	1	2	3	4
1. Hallucinations	2.49	0.93	1.94	1.96
2. Thought disorganization	0.93	1.59	1.21	1.23
3. Unusual thoughts	1.94	1.21	3.51	2.55
4. Duration	1.96	1.23	2.55	3.32

The observed (top) and reconstructed (bottom) covariance matrices for the predictor variables included in the path diagram above were calculated using the corresponding structural equation model and a set of starting values for each unfixed parameter in the model.

- Construct this observed covariance matrix in R with the proper row and column names
- Construct this reconstructed covariance matrix in R with the proper row and column names
- Describe if it appears that this structural equation has sufficiently reconstructed the covariance structure among the original predictor dimensions
- Find the difference between the observed and reconstructed covariance matrices using the maximum likelihood principle