# Path Diagrams

James H. Steiger

*Path Diagrams* play a fundamental role in structural modeling. In this handout, we discuss aspects of path diagrams that will be useful to you in describing and reading about confirmatory factor analysis models and structural equation models.

## 1. An Introduction to Path Diagrams

Path diagrams are like flowcharts. They show variables interconnected with lines that are used to indicate causal flow. Each *path* involves two variables (in either boxes or ovals) connected by either *arrows* (lines, usually straight, with an arrowhead on one end) or *wires* (lines, usually curved, with no arrowhead), or "slings" (with two arrowheads).

Arrows are used to indicate "directed" relationships, or linear relationships between two variables. An arrow from $X$ to $Y$ indicates a linear relationship where $Y$ is the dependent variable and $X$ the independent variable.
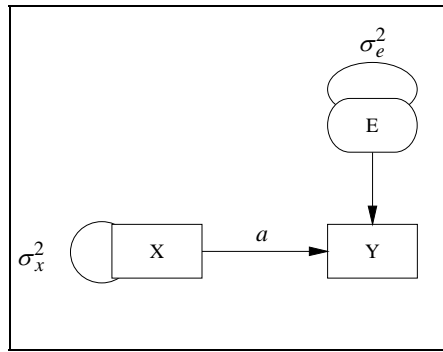
Wires or Slings are used to represent "undirected" relationships, which represent variances (if the line curves back from a variable to itself) or covariances (if the line curves from one variable to another).

One can think of a path diagram as a device for showing which variables cause changes in other variables. However, path diagrams need not be thought of strictly in this way. They may also be given a narrower, more specific interpretation.
Consider the classic linear regression equation

$$Y = aX + E$$

and its path representation shown below.

Such diagrams establish a simple isomorphism. All variables in the equation system are placed in the diagram, either in boxes or ovals. Each equation is represented on the diagram as follows: All independent variables (the variables on the right side of an equation) have *arrows* pointing to the dependent variable. The weighting coefficient is placed in clear proximity to the arrow.

Notice that, besides representing the linear equation relationships with arrows, the diagrams also contain some additional aspects. First, the variances of the independent variables, which must be specified in order to test the structural relations model, are shown on the diagrams using curved lines ("wires") without arrowheads attached, or two-headed arrows (sometimes called "slings"). Second, some variables are represented in ovals, others in rectangular boxes. *Manifest variables* (i.e., those that can be measured directly) are placed in boxes in the path diagram. *Latent variables* (i.e., those that cannot be measured directly, like factors in factor analysis, or residuals in regression) are placed in an oval or circle. For example, the variable $E$ in the above diagram can be thought of as a linear regression residual when $Y$ is predicted from $X$. Such a residual is not observed directly, but is in principle calculable from $Y$ and $X$ (if $a$ is known), so it is treated as a *latent* variable and placed in an oval.

The example discussed above is an extremely simple one. Generally, one is interested in testing much more complicated models. As the equation systems under examination become increasingly complicated, so do the *covariance structures* they imply. Ultimately, the complexity can become so bewildering that one loses sight of some very basic principles. For one thing, the train of reasoning which supports testing causal

models with linear structural equations testing has several weak links. The relationships between variables may be non-linear. They may be linearly related for reasons unrelated to what we commonly view as causality. The old statistical adage, "correlation is not causation" remains true, even if the correlation is complex and multivariate. What causal modeling *does* allow you to do is examine the extent to which data fail to agree with one consequence (viz., the implied covariance structure) of a model of causality. If the linear equations system isomorphic to the path diagram does fit the data well, it encourages continued belief in the model, but does not prove its correctness.

Although path diagrams can be used to represent causal flow in a system of variables, they need not imply such a causal flow. Path diagrams may be viewed as simply an isomorphic representation of a linear equations system. As such, they can convey linear relationships whether or not causal relations are assumed. Hence, although one *might* interpret the diagram in the above figure to mean that "$X$ causes $Y$," the diagram can also be interpreted as a visual representation of the linear regression relationship between $X$ and $Y$.

## 2. PATH1 Rules for Path Diagrams

In this section, rules for path diagrams are established that will guarantee that the diagram will represent accurately any model which fully accounts for all variances and covariances of all variables, both manifest and latent. These rules are based on the following considerations.

Path diagrams consist of variables connected by *wires* and *arrows*, representing, respectively, *undirected* and *directed* relationships between variables. These variables must be either *endogenous* or *exogenous*. (An *endogenous variable* is one that is a dependent variable in at least one linear equation in the equation system under consideration; an *exogenous variable* is one that is never a dependent variable. In a path diagram, endogenous variables have at least one arrow pointing to them, exogenous variables have no arrows pointing to them.) The variables must also be either *manifest* or *latent*. Hence any variable can be classified into 4 categories: (a)

manifest endogenous, (b) manifest exogenous, (c) latent endogenous, and (d) latent exogenous.

If random variables are related by linear equations, then variables which are endogenous have variances and covariances which are determinate functions of the variables on which they regress. For example, if $X$ and $Y$ are orthogonal and

$$W = aX + bY \, ,$$

then

$$\sigma^2_W = a^2 \sigma^2_X + b^2 \sigma^2_Y \, .$$

Hence, one way of guaranteeing that a diagram can account for variances and covariances among all its variables is to require:

(1)     representation of all variances and covariances among exogenous variables,

(2)     no variances or covariances to be directly represented in the diagram for endogenous variables, and

(3)     all variables in the diagram be involved in at least one relationship.

There is a significant practical problem with many path diagrams — lack of space. In many cases, there are so many exogenous variables that there is simply not enough room to represent, adequately, the variances and covariances among them. Diagrams which try often end up looking like piles of spaghetti.

One way of compensating for this problem is to include rules for default variances and covariances which allow a considerable number of them to be represented implicitly in the diagram.

These considerations lead to the following rules:

(1) *Manifest variables* are always represented in boxes (squares or rectangles) while *latent variables* are always in ovals or circles.

(2) Each directed relationship is represented explicitly by an *arrow* between two variables.

(3) Undirected relationships need not be represented explicitly. (See rule 9 below regarding implicit representation of undirected relationships.)

(4) Undirected relationships, when represented explicitly, are shown by a *wire* from a variable to itself, or from one variable to another.

(5) *Endogenous variables* may never have wires connected to them.

(6) *Free parameter* numbers for a *wire* or *arrow* are always represented with integers or labels placed on, near, or slightly above the middle of the wire or arrowline. A *free parameter* is a number whose value is estimated by the program. Two free parameters having the same parameter number or label are required to have the same value.

(7) *Fixed values* for a *wire* or *arrow* are always represented with a floating point number *containing a decimal point.* The number is generally placed on, near, or slightly above the middle of the *wire* or *arrow* line. A fixed value is assigned by the user. (There are default values that are applicable in various situations.)

(8) Different statistical populations are represented by a line of demarcation and the words *Group 1* (for the first population or group), *Group 2*, etc., in each diagram section.

(9) All *exogenous* variables must have their variances and covariances represented either explicitly or implicitly by either free parameters or fixed values. If variances and covariances are not represented explicitly, then the following rules hold:

(9a) Among latent exogenous variables, variances not explicitly represented in the diagram are assumed to be fixed values of *1.0*, and covariances not explicitly represented are assumed to be fixed values of *0*.

(9b) Among manifest exogenous variables, variances and covariances not explicitly represented are assumed to be *free parameters* each having a different parameter number. These parameter numbers are not equal to any number appearing explicitly in the diagram.
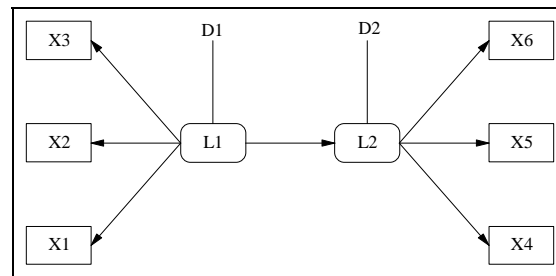
By adopting a consistent standard for path diagrams, we can facilitate clear communication of path models, regardless of what system is used to analyze them.

The typical beginning student of SEM will attempt to reproduce results from published papers employing a wide variety of standards for their path diagrams. In some cases

this approach will create no problems. However, experience indicates that it is often useful to translate published diagrams into diagram that obeys rules 1-9 above, before specifying the model for estimation. Frequently the translation process will draw attention to errors or ambiguities in the published diagram. This issue will be discussed in the following section.

## 3. Resolving Ambiguities in Path Diagrams

The figure below shows a portion of a path diagram which is quite typical of what is found in the literature. This is not a complete diagram and it does not conform to diagramming rules in the preceding section.
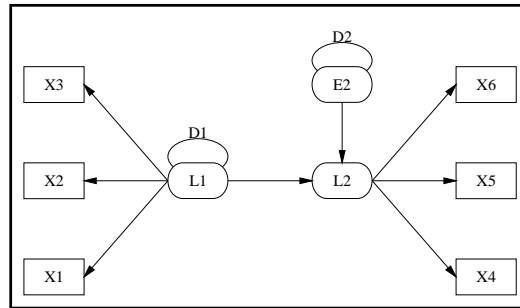


Some of the diagram is clear and routine, but what do we make of the symbols *D1* and *D2*? Variable *L1* is a latent exogenous variable. It has arrows pointing away from it and no arrows pointing to it. Since, by rule 9 for diagrams (see above), all exogenous variables must have their variances and covariances explained, the most reasonable assumption is that *D1* stands for the variance of latent variable *L1*. Hence, the diagram is modified to make *D1* a parameter attached to a wire from *L1* to itself.

But what is the status of *D2*? In the diagram it looks just like *D1*, but closer inspection reveals it must mean something different. *D2* is connected to *L2*, and *L2* is an endogenous latent variable. Consequently, the most reasonable interpretation is that *D2* represents an error variance for latent variable *L2*. It is represented with an error latent variable *E2* with variance *D2*.

The revised path diagram, more accurately reflecting the author's model, is shown in the figure below. Notice, however, that the diagram is still not fully explicit. Each of the

manifest variables is endogenous, and, as such, needs an error (or residual) variance. Many path diagrams, for the sake of compactness, will not include these paths.



In some cases you will have to be creative, tenacious, and lucky to figure out what the author of a path diagram intended. Even the most accomplished and generally careful authors will leave out paths, forget to mention that some values were fixed rather than free parameters, or simply misrepresent the model actually tested. Sometimes the only way to figure out what the author actually did is to try several models, until you find coefficients which agree with the published values. These difficulties are compounded by the occasional typographical errors that appear in published covariance and correlation matrices.

It seems reasonable to conclude that if authors were to adopt diagramming rules and/or report their models in the *PATH1* language, these problems would be reduced.

Some path diagrams do not represent the error variance attached to endogenous latent variables at all — they leave this to the reader to figure out for him/herself. Whenever an endogenous latent variable has no error term, you should suspect that an error latent variable has been left out, especially if your degrees of freedom don't agree with those of the published paper.

## 4. The RAM Diagramming System

The later version of the RAM system developed by Jack McArdle adds an additional twist to this. A *residual variable* is an exogenous variable that has a directed path to

one (and only one) endogenous variable. In the RAM system, *residual variables* are not represented explicitly in the diagram. Rather, their variances are shown as two-headed arrows (or "slings") attached to the variable they point to. For an example of how this works, see the handout on *Confirmatory Factor Analysis with R*.