

Using Path Diagrams as a Structural Equation Modelling Tool

by Peter Spirtes, Thomas Richardson, Chris Meek, Richard Scheines, and Clark Glymour¹

1. Introduction

Linear structural equation models (SEMs) are widely used in sociology, econometrics, biology, and other sciences. A SEM (without free parameters) has two parts: a probability distribution (in the Normal case specified by a set of linear structural equations and a covariance matrix among the “error” or “disturbance” terms), and an associated path diagram corresponding to the causal relations among variables specified by the structural equations and the correlations among the error terms. It is often thought that the path diagram is nothing more than a heuristic device for illustrating the assumptions of the model. However, in this paper, we will show how path diagrams can be used to solve a number of important problems in structural equation modelling.

There are a number of problems associated with structural equation modeling. These problems include:

- How much do sample data underdetermine the correct model specification? Of course, one must decide how much credence to give alternative explanations that afford different fits to any particular data set. There are a variety of techniques for that purpose, including Bayesian updating, and a variety of fit measures with well understood large sample properties. But what about two or more alternative models that fit a specific data set equally well, or, subject to certain restrictions, fit any data set meeting the restrictions equally well? The number of such equivalents for a given linear structural equation model may be very large. Even if there are sources of knowledge about structure from outside the data set, the number of equivalent models all meeting those knowledge constraints may be considerable, and the structures they postulate may have importantly different implications for policy. Unless we characterize such equivalencies, selection of a particular model can only involve an element of arbitrary choice.

¹ Spirtes, Glymour and Scheines are in the Department of Philosophy, Carnegie Mellon University. Richardson is in the Department of Statistics, University of Washington. Meek is at Microsoft Research. Thomas Richardson wishes to thank the Isaac Newton Institute, where he was a Rosenbaum fellow, while preparing this paper. The research was also supported under NSF Grants DMS-9704573, BES-940239, and IRI-9424378.

- Given that there are equivalent models, is it possible to extract the features common to those models? Under some circumstances, every member of a set of equivalent models may share some of the same linear coefficients or correlated errors. If that is the case, then it is possible that even though the data may not help us choose between the different models, the data may provide evidence for features common to all of the best models.

- When a modeler draws conclusions about coefficients in an unknown underlying structural equation model from a multivariate regression, precisely what assumptions are being made about the structural equation model? For example, when does a non-zero partial regression coefficient correspond to a non-zero coefficient in a structural equation?

These questions have been addressed many times, though usually only for models with special structures, and usually relying on linear algebra, the mathematics that seems most natural for a study of linear models. The aim of this paper is to explain how the path diagram provides much more than heuristics for special cases; the theory of path diagrams helps to clarify several of the issues just noted, issues that have been the focus of intelligent--if, in our judgment, ultimately too sweeping-- criticism of the use of structural equation models. What follows is a report that describes some of what has been learned about these issues by following a different set of mathematical ideas that exploit the graphical structure implicit in structural equation models.

In particular, we will present answers to these questions that depend upon an understanding of the relationship between the path diagram used to represent a structural equation model, and the zero partial correlations entailed by that path diagram (entailed in the sense that every structural equation model that shares the path diagram has a zero partial correlation). We will describe a graphical relation, the Pearl-Geiger-Verma d-separation criterion, among a pair of variables X and Y , and a set of variables Z , that is a necessary and sufficient condition for a structural equation model to entail a zero partial correlation. Such necessary and sufficient conditions have been known for path diagrams without correlated errors, but we will extend the conditions to path diagrams with correlated errors.

In section 2 we will motivate interest in the d-separation relation by describing the problems that it helps to solve in more detail. Then in section 3 we will show how the zero partial correlations entailed by a structural equation model can be read off from its path diagram, and in section 4 use the machinery developed in section 3 to provide some solutions to problems described in section 2. In section 5 we discuss the broader implications of this work for model selection, and illustrate this with two examples in section 6. In section 7 we prove the main theorem, hitherto unpublished, which justifies the

use of d-separation in path diagrams representing correlated errors (represented by edges of the form \rightleftarrows , which we call double-headed arrows).

2. Problems in SEM Modeling

In order to describe the problems listed in section 1 in more detail, we will first review how path diagrams are used to represent structural equation models without free parameters. The path diagram contains a directed edge from B to A if and only if there is a non-zero coefficient for B in the equation for A; and there is a double-headed arrow between A and B if and only if the error term for A and the error term for B have a non-zero correlation.² The path diagram associated with a SEM may contain directed cycles (representing feedback), and double-headed arrows (representing correlated errors.) We will call a path diagram which contains no double-headed arrows a **directed graph**. (We place sets of variables and defined terms in boldface.) In a SEM M , we will denote the correlation matrix among the non-error variables by $\Sigma(M)$, and the corresponding path diagram by $G(M)$. We will now review the problems mentioned in section 1 in more detail.

2.1. Covariance Equivalence

Consider the following example. The graph in Figure 1(a) is the path diagram of a SEM M proposed by Aberle (Blalock, 1961) as a model for evolutionary culture in American Indian tribes, where W is matridominant division of labor, X is matrilineal residence, Y is matricentered land tenure, and Z is matrilineal system of descent.

Suppose for the moment that there is a SEM with the path diagram in Figure 1(a) and the $p(\chi^2)$, the AIC (Aikake Information Criterion), and the BIC (Bayes Information Criterion) score for this SEM are all high³ (See Raftery (1995) for a discussion of the BIC score.) In order to evaluate how well the data supports this model, it is still necessary to know whether or not there are other models compatible with background knowledge that fit the data equally well (Lee and Hershberger, (1990), Stelzl (1986)). In this case, for each of the path diagrams in Figure 1, and for *any* data set D , there is a SEM with that path diagram that fits D as well as M does (in the sense that each SEM has the same $p(\chi^2)$ and the same

² This is slightly different than the usual convention in which if ϵ_A and ϵ_B are correlated, then they are explicitly included in the graph, there is a directed edge from ϵ_A to A , a directed edge from ϵ_B , and the double-headed arrow is placed between ϵ_A and ϵ_B . However, the convention adopted here will simplify later theorems and proofs.

³ In counting degrees of freedom, we will assume that a SEM with free parameters (and no latents) associates a linear coefficient parameter with each directed edge (i.e. β) in its path diagram, a correlation parameter with each double-headed arrow (i.e. ρ) in its path diagram, and a variance parameter with each vertex. We also assume that no extra constraints (such as equality constraints among parameters) are imposed.

BIC and AIC scores.) If \mathbf{O} represent the set of measured variables in path diagrams G_1 and G_2 , then G_1 and G_2 are **covariance equivalent over \mathbf{O}** if and only if for every SEM M such that $G(M) = G_1$, there is a SEM M' with path diagram $G(M') = G_2$, and the marginal of $\Sigma(M')$ over \mathbf{O} equals the marginal of $\Sigma(M)$ over \mathbf{O} , and vice-versa.⁴ (Informally, any covariance matrix over \mathbf{O} generated by a parameterization of path diagram G_1 can be generated by a parameterization of path diagram G_2 , and vice-versa.) If G_1 and G_2 have no latent variables, (i.e all of the variables in their path diagrams are in \mathbf{O}), then we will simply say that G_1 and G_2 are **covariance equivalent**. If two covariance equivalent models are equally compatible with background knowledge, and have the same degrees of freedom, the data does not help distinguish them, so it is important to be able to find the complete set of path diagrams that are covariance equivalent to a given path diagram. (Every SEM that contains a path diagram in Figure 1 has the same number of degrees of freedom.)

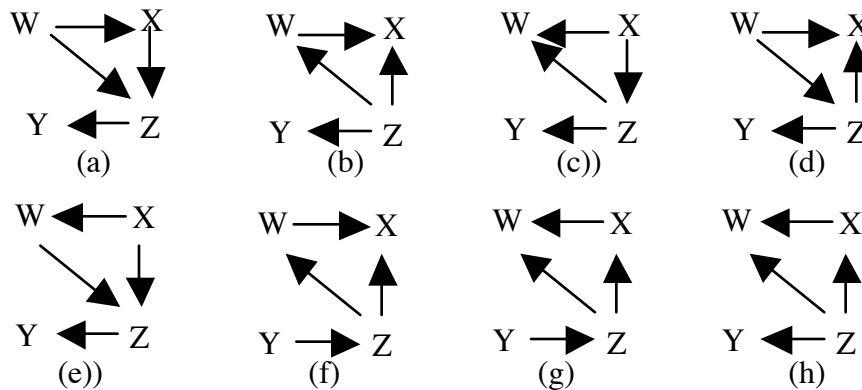


Figure 1

As we will illustrate below, it is often far from obvious what constitutes a complete set of path diagrams covariance equivalent to a given path diagram. We will call such a complete set a **covariance equivalence class over \mathbf{O}** . (Again, if we consider only SEMs without latent variables, we will call such a complete set a **covariance equivalence class**. If it is a complete set of path diagrams without correlated errors or directed cycles, i.e. directed acyclic graphs, that are covariance equivalent we will call it a **simple covariance equivalence class over \mathbf{O}** .) As shown in section 4, the path diagrams in Figure 1 are a simple covariance equivalence class.

⁴ For technical reasons, a more formal definition requires a slight complication. G is a **sub-path diagram** of G' when G and G' have the same vertices, and G has a subset of the edges in G' . G_1 and G_2 are **covariance equivalent over \mathbf{O}** if for every SEM M such that $G(M) = G_1$, there is a SEM M' with path diagram $G(M')$ that is a sub-path diagram of G_2 , and the marginal over \mathbf{O} of $\Sigma(M')$ equals the marginal over \mathbf{O} of $\Sigma(M)$, and for every SEM M' such that $G(M') = G_2$, there is a SEM M with path diagram $G(M)$ that is a sub-path diagram of G_1 , and the marginal over \mathbf{O} of $\Sigma(M)$ equals the marginal over \mathbf{O} of $\Sigma(M')$.

Another example of a case where it is not obvious whether or not two path diagrams are covariance equivalent over \mathbf{O} is shown below. It is often thought that the two path diagrams in Figure 2 (each of which is part of a just-identified SEM) are covariance equivalent over $\mathbf{O} = \{X, Y, Z\}$. However, as shown in Spirtes et al. (1996), there is a SEM with path diagram in Figure 2(b) with the covariance matrix Σ over $X, Y,$ and $Z,$ but there is no SEM that contains the path diagram in Figure 2 (a) with marginal covariance matrix Σ (where $T_1, T_2,$ and T_3 are latent variables).

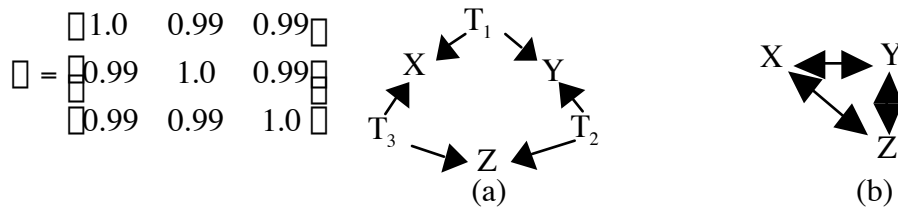


Figure 2

In section 4, we will describe how to efficiently test when two path diagrams without correlated errors or directed cycles are covariance equivalent. We will also give informative necessary conditions for two path diagrams with correlated errors, cycles, or latent variables to be covariance equivalent over \mathbf{O} . For related theorems see also Pearl (1997).

2.2. Features Common to a Covariance Equivalence Class

A second important question that arises with respect to covariance equivalence classes is whether it is possible to extract the features that the set of covariance equivalent path diagrams have in common. For example, every path diagram in Figure 1 has the same adjacencies, but the path diagrams do not have any edge with the same orientation in every member of the equivalence class (e.g. both $W \rightarrow X,$ and $W \leftarrow X$ occur in path diagrams in Figure 1).

However, there are other sets of covariance equivalent path diagrams in which a given edge always occurs with the same orientation in every member of the equivalence class. For example, Figure 3 shows another simple covariance equivalence class of graphs in which the orientation $X \rightarrow Z$ occurs in every member of the equivalence class.



Figure 3

This is informative because even though the data does not help choose between members of the equivalence class, insofar as the data is evidence for the disjunction of the members in the equivalence class, it is evidence for the orientation $X \rightarrow Z$.

In section 4 we will show how to extract all of the features common to a simple covariance equivalence class of path diagrams, and briefly indicate how it is possible to extract some features common to a covariance equivalence class of path diagrams with correlated errors, cycles, or latent variables.

2.3. Regression Coefficients and Structural Equation Coefficients

It is common knowledge among practising social scientists that for the coefficient of X in the regression of Y upon X to be interpretable as the effect of X on Y there should be no "confounding" variable Z which is a cause of both X and Y:

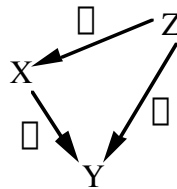


Figure 4

Simple calculations confirm this conclusion (using the notation in Figure 4):⁵

$$\text{Cov}(X, Y) = \alpha V(X) + \beta \gamma V(Z)$$

Hence

$$\frac{\text{Cov}(X, Y)}{V(X)} = \frac{\alpha V(X) + \beta \gamma V(Z)}{V(X)} \neq \alpha.$$

Thus the coefficient from the regression of Y on X alone will be a consistent estimator only if either β or γ is equal to zero. Further, observe that the bias term $\beta \gamma V(Z)/V(X)$ may be either positive or negative, and of arbitrary magnitude.

However, $\text{Cov}(X, Z) = \beta V(Z)$ and $\text{Cov}(Y, Z) = (\alpha \gamma + \gamma) V(Z)$, and hence

$$\begin{aligned} \text{Cov}(X, Y | Z) &\equiv \text{Cov}(X, Y) - \frac{\text{Cov}(X, Z)\text{Cov}(Y, Z)}{V(Z)} \\ &= \alpha V(X) + \beta \gamma V(Z) - \beta \gamma V(Z) (\alpha \gamma + \gamma) \\ &= \alpha (V(X) - \beta^2 V(Z)) \end{aligned}$$

⁵ Section 7 after Lemma 5 contains a simple rule for calculating covariances from a path diagram. This rule is related to Wright's use of path coefficients (Wright, 1934).

and

$$V(X|Z) \equiv V(X) - \frac{\text{Cov}(X, Z)^2}{V(Z)} = V(X) - \beta^2 V(Z),$$

so the coefficient of X in the regression of Y on X and Z is a consistent estimator of β since $\text{Cov}(X, Y|Z)/V(X|Z) = \beta$.

The danger presented by failing to include confounding variables is well understood by social scientists. Indeed, it is often used as the justification for considering a long “laundry list” of “potential confounders” for inclusion in a given regression equation.

What is perhaps less well understood is that including a variable which is not a confounder can also lead to biased estimates of the structural coefficient. We now consider a number of simple cases demonstrating this.



Figure 5

In the SEM with the path diagram depicted in Figure 5, $\text{Cov}(X, Y) = \beta V(X)$, hence the coefficient of X in the regression of Y upon X is a consistent estimator of β . However, $\text{Cov}(Y, Z) = \beta V(Y)$, and $\text{Cov}(X, Z) = \beta \beta V(X)$, so that

$$\frac{\text{Cov}(X, Y|Z)}{V(X|Z)} = \beta \frac{V(Z) - \beta^2 V(Y)}{V(Z) - \beta^2 \beta^2 V(X)} = \beta \frac{V(\epsilon_Z)}{\beta V(\epsilon_Z) + \beta^2 V(\epsilon_Y)}$$

Hence the coefficient of X in the regression of Y on X and Z is an inconsistent estimator of β . The estimate will have the same sign as β , but will have smaller absolute magnitude. Note that $\text{Cov}(X, Y|Z)/V(X|Z) = 0$ if and only if $\beta = 0$.

It might be objected that this type of error is unlikely to arise in practise since often information about time order would rule out Z as a potential unmeasured confounder. In the next example this response is not applicable since Z may temporally precede both X and Y. Let ϵ_X , ϵ_Y , and ϵ_Z be the error variables in Figure 6(a), and ϵ'_X , ϵ'_Y , and ϵ'_Z be the error variables in Figure 6 (b).

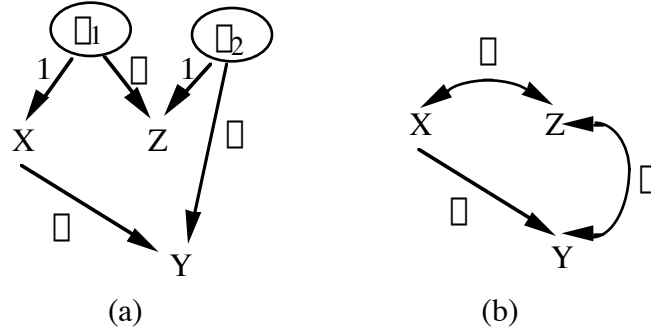


Figure 6

In the path diagram depicted in Figure 6(a) there are two unmeasured confounders T_1 and T_2 , which are uncorrelated with one another. Any SEM with this path diagram may be converted into a SEM with the path diagram depicted in Figure 6(b), letting $\alpha = \text{Cov}(X, Z) = \alpha V(T_1)$, $\beta = \alpha V(T_2)$, $V(\epsilon_X^*) = V(\epsilon_X) + V(T_1)$, $V(\epsilon_Y^*) = V(\epsilon_Y) + \alpha^2 V(T_1) + V(T_2)$, and $V(\epsilon_Z^*) = V(\epsilon_Z) + \alpha^2 V(T_2)$.

Note however, that the reverse is not in general true: not every model containing correlated errors ($X \sim Y$) can be converted into a SEM model with latent variables but without correlated errors by introducing a latent T that is a parent of X and Y ($X \leftarrow T \rightarrow Y$), as pointed out in section 2.1. (It is however always possible to convert a model with correlated errors into *some* latent variable model without correlated errors, but which may contain more than one latent common cause of each pair of variables. This is because every normal distribution is a linear transformation of a set of independent normal variables, which can play the role of the latent variables.)

Returning to the path diagram in Figure 6(b) note that the regression of Y on X yields a consistent estimate of α since $\text{Cov}(X, Y) = \alpha V(X)$. However,

$$\begin{aligned} \frac{\text{Cov}(X, Y|Z)}{V(X|Z)} &= \frac{\text{Cov}(X, Y)V(Z) - \text{Cov}(X, Z)\text{Cov}(Y, Z)}{V(X)V(Z) - \text{Cov}(X, Z)^2} \\ &= \frac{\alpha V(X)V(Z) - \alpha(\alpha + \beta)}{V(X)V(Z) - \alpha^2} \\ &= \alpha \frac{V(Z)}{V(X)V(Z) - \alpha^2} \end{aligned}$$

Hence the coefficient of X in the regression of Y on X and Z is not a consistent estimate of α , (unless $\beta = 0$ or $\alpha = 0$), and may even have a completely different sign. In the case where $\beta = 0$, the coefficient of X in the regression of Y on X will be zero in the population, but will become non-zero once Z is included.

SEM folklore often appears to suggest that it is better to include rather than exclude a variable from a regression. This notion is perhaps given support by reference to

“controlling for Z”, the implication being that controlling for Z eliminates a source of bias. The conclusion to be drawn from these examples is that there is no sense in which one is “playing safe” by including rather than excluding “potential confounders”; if they turn out not to be potential confounders then this could change a consistent estimate into an inconsistent estimate.

The situation is also made somewhat worse by the use of misleading definitions of 'confounder': sometimes a confounder is said to be a variable that is strongly correlated with both X and Y, or even a variable whose inclusion changes the coefficient of X in the regression. Since, for sufficiently large α and β , Z in Figure 6 would qualify as a confounder under either of these definitions, it follows that under either definition including confounding variables in a regression may make a hitherto consistent estimator inconsistent.

Finally, it is worth reiterating the well-known fact that in certain circumstances there may be no regression which will estimate the parameter of interest, (although some other consistent estimator may exist):

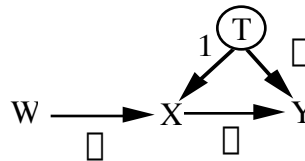


Figure 7

In the SEM shown in Figure 7, $Cov(X,Y) = \alpha V(X) + \beta V(T)$; hence the coefficient of X in the regression of Y on X is not a consistent estimator of α . Further

$$\frac{Cov(X, Y|W)}{V(X|W)} = \alpha + \frac{\beta V(T)}{V(X) - \alpha^2 V(W)} = \alpha + \frac{\beta V(T)}{V(T) + V(\alpha X)}$$

hence including W in the regression does not help matters. However, a consistent estimator exists, the so-called Instrumental Variable estimator:

$$\frac{Cov(Y, W)}{Cov(X, W)} = \frac{\beta V(W)}{\alpha V(W)} = \beta$$

In this discussion we have highlighted a number of problems that arise when estimating structural coefficients via regression. These examples raise the following general questions:

(a) If Y is regressed on a set of variables **W**, including X, in which SEMs will the partial regression coefficient of X be a consistent estimate of the structural coefficient α associated with the X \rightarrow Y edge?

(b) If Y is regressed on the set \mathbf{W} , which includes X , in which SEMs will the partial regression coefficient of X be zero if the structural coefficient associated with the $X \rightarrow Y$ edge is zero?

(c) Given a particular SEM in which there is an edge $X \rightarrow Y$ with coefficient β , is it possible to find a subset \mathbf{W} of observed variables (including X), such that when Y is regressed on the set \mathbf{W} , the coefficient of X in the regression is a consistent estimate of β ?

(d) Given a particular SEM and a structural coefficient β , is it possible to find a function $h(\mathbf{S})$ (where \mathbf{S} is the sample covariance matrix) that is a consistent estimator of β ?

We shall answer questions (a), (b) and (c), by applying the graphical criterion of d -separation. One advantage of a graphical criterion is that it can be applied simply by visual inspection of the path diagram, and does not require lengthy algebraic manipulations which become increasingly arduous when more variables are involved in the calculation. We do not know the answer to (d), which is one form of the well-known "identification problem"; it is possible that extensions of the graphical criteria we present may hold the key. For related theorems, see Pearl (1997).

2.4. *Other Applications*

In addition to the uses described above, there are a number of other applications that we do not have the space to describe here. The d -separation relation has proved useful in automated search for causal structure from data and background knowledge (Spirtes and Glymour, 1991, Spirtes, Glymour and Scheines, 1993, Pearl and Verma, 1991, Cooper, 1992), in calculating the effects of interventions on causal systems (Spirtes, Glymour and Scheines, 1993, and Pearl, 1995), and has shed light on a number of issues in statistics ranging from Simpson's Paradox to experimental design (Spirtes, Glymour and Scheines, 1993). See also the applications in Pearl (1997).

3. **Linear Structural Equation Models and d -separation**

In a linear SEM the random variables are divided into two disjoint sets, the substantive variables and the error variables. Corresponding to each substantive random variable V is a unique error term ϵ_V .⁶ A linear SEM contains a set of linear equations in which each substantive random variable V is written as a linear function of other substantive random variables together with ϵ_V , and a correlation matrix among the error terms. Initially, we will assume that the error variables are multi-variate Gaussian. However, many of the results that

⁶ There is an equivalent definition of a linear SEM in which parent-less or 'exogenous' substantive variables have no associated error variables.

we will prove are about partial correlations, which do not depend upon the distribution of the error terms, but depend only upon the linear equations and the correlations among the error terms.

Since we have no interest in first moments, without loss of generality each variable can be expressed as a deviation from its mean.

For example, the following is a linear SEM M , ϵ_A , ϵ_B , ϵ_C , ϵ_D , and ϵ_E are Gaussian "error terms", and A , B , C , D , and E are substantive random variables:

$$\begin{aligned} A &= \epsilon_A \\ B &= \epsilon_B \\ C &= .2B + .8D + \epsilon_C \\ D &= -.5C + .1E + \epsilon_D \\ E &= \epsilon_E \end{aligned}$$

Correlation Matrix Among Error Terms

	ϵ_A	ϵ_B	ϵ_C	ϵ_D	ϵ_E	
ϵ_A	1.0	0.5	0.0	0.0	0.0	
ϵ_B	0.5	1.0	0.0	0.0	0.0	
ϵ_C	0.0	0.0	1.0	0.0	0.0	
ϵ_D	0.0	0.0	0.0	1.0	0.0	
ϵ_E	0.0	0.0	0.0	0.0	1.0	

If the coefficients in the linear equations are such that the substantive variables are a unique linear function of the error variables alone, the set of equations is said to have a **reduced form**. A linear SEM with a reduced form also determines a joint distribution over the substantive variables. We will consider only linear SEMs which have coefficients for which there is a reduced form, all variances and partial variances among the substantive variables are finite and positive, and all partial correlations among the substantive variables are well defined (e.g. not infinite).

The path diagram of a linear SEM with uncorrelated errors is written with the conventions that it contains an edge $A \rightarrow B$ if and only if the coefficient for A in the structural equation for B is non-zero, and there is a double-headed arrow between two variables A and B if and only if the correlation between ϵ_A and ϵ_B is non-zero. Thus the path diagram for M is shown in Figure 8.

In order to define the d-separation relation, we need to introduce the following path diagram terminology. The concepts defined here are illustrated in Figure 8. A path diagram

consists of two parts, a set of vertices \mathbf{V} and a set of edges \mathbf{E} . Each edge in \mathbf{E} is between two distinct vertices in \mathbf{V} . There are two kinds of edges in \mathbf{E} , directed edges $A \rightarrow B$ or $A \leftarrow B$, and double-headed edges $A \leftrightarrow B$; in either case A and B are **endpoints** of the edge; further, A and B are said to be **adjacent**. There may be multiple edges between vertices. In Figure 8 the set of vertices is $\{A,B,C,D,E\}$ and the set of edges is $\{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow C, E \rightarrow D\}$. For a directed edge $A \rightarrow B$, A is the **tail** of the edge and B is the **head** of the edge, A is a **parent** of B , and B is a **child** of A .

An **undirected path** U between X_a and X_b is a sequence of edges $\langle E_1, \dots, E_m \rangle$ such that one endpoint of E_1 is X_a , one endpoint of E_m is X_b , and for each pair of consecutive edges E_i, E_{i+1} in the sequence, $E_i \neq E_{i+1}$, and one endpoint of E_i equals one endpoint of E_{i+1} . In Figure 8, $A \rightarrow B \rightarrow C \rightarrow D$ is an example of an undirected path between A and D . A **directed path** P between X_a and X_b is a sequence of directed edges $\langle E_1, \dots, E_m \rangle$ such that the tail of E_1 is X_a , the head of E_m is X_b , and for each pair of edges E_i, E_{i+1} adjacent in the sequence, $E_i \neq E_{i+1}$, and the head of E_i is the tail of E_{i+1} . For example, $B \rightarrow C \rightarrow D$ is a directed path. A **vertex occurs on a path** if it is an endpoint of one of the edges in the path. The set of vertices on $A \rightarrow B \rightarrow C \rightarrow D$ is $\{A, B, C, D\}$. A path is **acyclic** if no vertex occurs more than once on the path. $C \rightarrow D \rightarrow C$ is a cyclic directed path. The following is a list of all the acyclic directed paths in Figure 8: $B \rightarrow C, C \rightarrow D, E \rightarrow D, D \rightarrow C, B \rightarrow C \rightarrow D, E \rightarrow D \rightarrow C$.

A vertex A is an **ancestor** of B (and B is a **descendant** of A) if and only if either there is a directed path from A to B or $A = B$. Thus the ancestor relation is the transitive, reflexive closure of the parent relation. The following table lists the child, parent, descendant and ancestor relations in Figure 8.

Vertex	Children	Parents	Descendants	Ancestors
A	\emptyset	\emptyset	$\{A\}$	$\{A\}$
B	$\{C\}$	\emptyset	$\{B,C,D\}$	$\{B\}$
C	$\{D\}$	$\{B,D\}$	$\{C,D\}$	$\{B,C,D,E\}$
D	$\{C\}$	$\{C,E\}$	$\{C,D\}$	$\{B,C,D,E\}$
E	$\{D\}$	\emptyset	$\{C,D,E\}$	$\{E\}$

A vertex X is a **collider** on undirected path U if and only if U contains a subpath $Y \rightarrow X \rightarrow Z$, or $Y \leftarrow X \rightarrow Z$, or $Y \rightarrow X \leftarrow Z$, or $Y \leftarrow X \leftarrow Z$; otherwise if X is on U it is a **non-collider** on U . For example, C is a collider on $B \rightarrow C \rightarrow D$ but a non-collider on $B \rightarrow C \leftarrow D$. X is an **ancestor of a set** of vertices \mathbf{Z} if X is an ancestor of some member of \mathbf{Z} .

For disjoint sets of vertices, \mathbf{X} , \mathbf{Y} , and \mathbf{Z} , \mathbf{X} is **d-connected** to \mathbf{Y} given \mathbf{Z} if and only if there is an acyclic undirected path U between some member X of \mathbf{X} , and some member Y of \mathbf{Y} , such that every collider on U is an ancestor of \mathbf{Z} , and every non-collider on U is not in \mathbf{Z} . For disjoint sets of vertices, \mathbf{X} , \mathbf{Y} , and \mathbf{Z} , \mathbf{X} is **d-separated** from \mathbf{Y} given \mathbf{Z} if and only if \mathbf{X} is not d-connected to \mathbf{Y} given \mathbf{Z} .

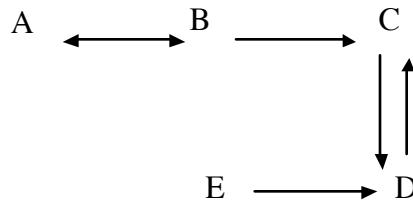


Figure 8

For example, the path $E \rightarrow D \rightarrow C$ d-connects E and C given \emptyset ; it also d-connects E and C given $\{A\}$, $\{B\}$, or $\{A,B\}$. $E \rightarrow D \rightarrow C$ d-connects E and C given $\{D\}$, given $\{D,B\}$, $\{D,A\}$, or $\{D,A,B\}$. The following is a list of all the pairwise d-separation relations in Figure 8 (where each pair is followed by a list of all of the sets that d-separate them):

- $\{A\}$ and $\{C\}$ are d-separated given: $\{B\}$, $\{B,D\}$, $\{B,E\}$, $\{B,D,E\}$
- $\{A\}$ and $\{D\}$ are d-separated given: $\{B\}$, $\{B,C\}$, $\{B,E\}$, $\{B,C,E\}$
- $\{A\}$ and $\{E\}$ are d-separated given: \emptyset , $\{B\}$, $\{B,C\}$, $\{B,D\}$, $\{B,C,D\}$, $\{C,D\}$
- $\{B\}$ and $\{E\}$ are d-separated given: \emptyset , $\{C,D\}$

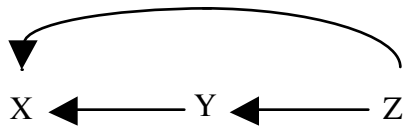
The first theorem states that d-separation in a path diagram G is a sufficient condition for G to entail that $\square(X,Y,Z) = 0$ (i.e. in every SEM with path diagram G , the partial correlation of X and Y given Z equals 0.)

Theorem 1: If M is a SEM, and $\{X\}$ and $\{Y\}$ are d-separated given Z in $G(M)$, then $\square(X,Y,Z) = 0$ in $\square(M)$.

The second theorem states that d-separation is a necessary condition for a path diagram to entail a zero partial correlation.

Theorem 2: If $\{X_i\}$ and $\{X_j\}$ are not d-separated given Z in path diagram G , then there is a SEM M such that $G(M) = G$, and $\square(X_i,X_j,Z) \neq 0$ in $\square(M)$.

Theorem 2 does *not* say that there might not be an individual SEM M with “extra” zero partial correlations among variables that are not d-separated in $G(M)$, as the following example shows.



$$X = .3 Y + .6 Z + \epsilon_x$$

$$Y = -2 Z + \epsilon_y$$

$$Z = \epsilon_z$$

Figure 9

(The errors are uncorrelated because there are no double-headed arrows in the path diagram.) In this case X and Y are independent, i.e. $\rho(X,Y) = 0$, even though $\{X\}$ and $\{Y\}$ are not d-separated given \emptyset . However, this zero correlation holds because of the particular linear coefficients. Thus, according to Theorem 2 there is some other SEM M such with the same path diagram in which $\rho(X,Y) \neq 0$. It has been shown (Spirtes et. al 1993) that the set of parameters which produce conditional independence relations among variables which are not d-separated in G has zero Lebesgue measure over the parameter space.

4. Applications

4.1. Covariance Equivalence for Path diagrams Without Correlated Errors or Directed Cycles

If for SEM M , there is another SEM M' with a different path diagram but the same number of degrees of freedom, and the same marginal distribution over the measured variables in M , then the $p(\chi^2)$ for M' equals $p(\chi^2)$ for M , and they have the same BIC scores and AIC scores. Such SEMs are guaranteed to exist if there are SEMs that have the same number of degrees of freedom and contain path diagrams which are covariance equivalent to each other. Stelzl (1986) and Lee and Hershberger (1990) discuss sufficient conditions for covariance equivalence (which they call simply equivalence). Theorem 3 states necessary, as well as sufficient conditions for covariance equivalence in path diagrams without correlated errors or directed cycles.

G_1 and G_2 are **d-separation equivalent** if for each disjoint sets \mathbf{X} , \mathbf{Y} , and \mathbf{Z} , \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} in G_1 if and only if \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} in G_2 .

Theorem 3: If G_1 and G_2 are directed acyclic graphs, G_1 and G_2 are covariance equivalent if and only if G_1 and G_2 are d-separation equivalent.

The test for covariance equivalence of two path diagrams described in Lee and Hershberger (1990) requires determining whether there is a series of edge replacements or reversals preserving equivalence that lead from one path diagram to the other. Because they

do not specify an ordering in which the tests are to be done, this could be a very slow process. The following theorem, due to Verma and Pearl (1990) shows how d-separation equivalence can be calculated in $O(E^2)$ time, where E is the number of edges in a path diagram. X is an **unshielded collider** in directed acyclic graph G if and only if G contains edges $A \rightarrow X \rightarrow B$, and A is not adjacent to B in G .

Theorem 4: Two directed acyclic graphs are d-separation equivalent if and only if they contain the same vertices, the same adjacencies, and the same unshielded colliders.

It is apparent from Theorem 4 that any two SEMs with covariance equivalent directed acyclic graphs have the same degrees of freedom.

4.2. *Covariance Equivalence for Path Diagrams with Correlated Errors or Directed Cycles*

Necessary conditions for covariance equivalence for path diagrams with correlated errors or cycles, and for path diagrams with latent variables follow from Theorem 1 and Theorem 2. If \mathbf{O} is a subset of the vertices in G_1 and a subset of the vertices in G_2 , then G_1 and G_2 are **d-separation equivalent over \mathbf{O}** if for each disjoint \mathbf{X} , \mathbf{Y} , and \mathbf{Z} included in \mathbf{O} , \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} in G_1 if and only if \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} in G_2 .

Theorem 5: If G_1 and G_2 are path diagrams that are covariance equivalent over \mathbf{O} , then G_1 and G_2 are d-separation equivalent over \mathbf{O} .

The converse is not generally true because while d-separation equivalence guarantees that the conditional independence constraints imposed by two path diagrams are the same, there are other, non-conditional independence constraints, that can be imposed by one path diagram but not the other. The path diagrams in Figure 2 are examples of path diagrams that are d-separation equivalent, but not covariance equivalent over $\mathbf{O} = \{X, Y, Z\}$.

If V is the maximum of the number of variables in G_1 or G_2 , and M is the number of variables in \mathbf{O} , Spirtes and Richardson (1996) presents an $O(M^3 \times V^2)$ algorithm for checking whether two acyclic path diagrams G_1 and G_2 (which may contain latent variables and correlated errors) are d-separation equivalent over \mathbf{O} . Richardson (1996) presents an $O(V^7)$ algorithm for determining when two cyclic path diagrams without latent variables are d-separation equivalent.

4.3. *Extracting Features Common to a Covariance Equivalence Class*

Theorem 4 is also the basis of a simple representation (called a pattern in Verma and Pearl 1990) of the entire set of path diagrams without correlated errors or cycles covariance equivalent to a given path diagram without correlated errors or cycles. The pattern for each

path diagram in Figure 1 is shown in Figure 10(a), and the pattern for each path diagram in Figure 3 is shown in Figure 10(b).

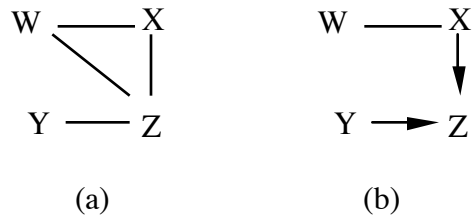


Figure 10

A pattern has the same adjacencies as the path diagrams in the covariance equivalence class that it represents. In addition, an edge is oriented as $X \rightarrow Z$ in the pattern if and only if it is oriented as $X \rightarrow Z$ in every path diagram in the simple covariance equivalence class. Meek (1995), Andersson et al. (1995), and Chickering (1995) show how to generate a pattern from an acyclic graph in $O(E)$ time (where E is the number of edges.)

In the case of acyclic path diagrams which may also contain latent variables, and the case of cyclic path diagrams which do not contain latent variables, there is an object analogous to a pattern called a Partial Ancestral Graph (PAG), which contains only measured variables but represents some of the features common to the members of a covariance equivalence class over \mathbf{O} . Spirtes and Verma (1992) shows how to create a PAG⁷ from an acyclic path diagram in $O(V^5)$ time (where V is the number of vertices in the path diagram). Richardson (1996c) presents an $O(V^7)$ algorithm for constructing a PAG from a (possibly cyclic) graph.

4.4. Solutions to the questions on regression

In this section we apply d-separation in order to answer three questions about the use of regression to estimate structural coefficients that we raised earlier. We introduce the following notation first: Given a SEM with path diagram G , we define $G \setminus \{X \rightarrow Y\}$ as the path diagram in which the $X \rightarrow Y$ edge is removed.

(a) *If Y is regressed on a set of variables W , including X , in which SEMs will the partial regression coefficient of X be a consistent estimate of the structural coefficient β associated with the $X \rightarrow Y$ edge?*

⁷ The algorithm given by Spirtes and Verma was designed to output an object called a partially oriented inducing path graph (POIPG); however, it has subsequently been shown that the output can be re-interpreted as a PAG.

The coefficient of X is a consistent estimator of β if \mathbf{W} does not contain any descendant of Y in G , and X is d-separated from Y given \mathbf{W} in $G \setminus \{X \sqsubseteq Y\}$.⁸ If this condition does not hold, then for almost all instantiations of the parameters in the SEM, the coefficient of X will fail to be a consistent estimator of β .

It follows directly from this that (almost surely) β cannot be estimated consistently via any regression equation if either there is an edge $X \sqsubseteq Y$ (i.e. β_X and β_Y are correlated) or if X is a descendant of Y (so that the path diagram is cyclic). The result itself follows from the fact that under the conditions stated,

$$\text{Cov}(X, \beta_Y \mid \mathbf{W} \setminus \{X\}) = \text{Cov}(X, Q \mid \mathbf{W} \setminus \{X\}) = 0$$

for each $Q \sqsubseteq \mathbf{Parents}(Y, G) \setminus \{X\}$. ($\mathbf{Parents}(Y, G)$ is the set of parents of Y in G .) It follows that

$$\text{Cov}(X, Y \mid \mathbf{W} \setminus \{X\}) = \text{Cov}(X, \beta_X + \sum_{Q_i \sqsubseteq \mathbf{Parents}(Y)} a_i Q_i + \beta_Y \mid \mathbf{W} \setminus \{X\}) = \beta_Y \text{V}(X \mid \mathbf{W} \setminus \{X\})$$

and hence $\frac{\text{Cov}(X, Y \mid \mathbf{W} \setminus \{X\})}{\text{V}(X \mid \mathbf{W} \setminus \{X\})} = \beta_Y$.

(b) *If Y is regressed on the set \mathbf{W} , including X , in which SEMs will the partial regression coefficient of X be zero if there is no edge between X and Y ?*

The coefficient of X will be zero if X and Y are d-separated given $\mathbf{W} \setminus \{X\}$. (See Scheines (1994) and Glymour (1994)). This follows directly from the fact that the coefficient of X in the regression equation is proportional to $\beta(X, Y, \mathbf{W} \setminus \{X\})$, which in turn will be zero if $\{X\}$ is d-separated from $\{Y\}$ given $\mathbf{W} \setminus \{X\}$. As before, if $\{X\}$ and $\{Y\}$ are not d-separated given $\mathbf{W} \setminus \{X\}$, then, even if there is no edge between X and Y , for almost all assignments of values to the model parameters the coefficient of X will be non-zero.

(c) *Given a particular SEM, with path diagram G , in which there is an edge $X \sqsubseteq Y$, with coefficient β , is it possible to find a subset \mathbf{W} of observed variables, (including X), such that when Y is regressed on the set \mathbf{W} , the coefficient of X in the regression is a consistent estimate of β ?*

From (a), we know that if there is a subset \mathbf{W} of the observed variables which contains no descendant of Y , but which d-separates X from Y in $G \setminus \{X \sqsubseteq Y\}$, then the regression coefficient of X in the regression of Y on \mathbf{W} will be a consistent estimate of β .

⁸Note this criterion is similar to Pearl's back door criterion (Pearl, 1993), except that the back-door criterion was proposed as a means of estimating the *total* effect of X on Y .

5. Implications for Model Selection

In this section, we discuss some of the methodological implications of the results presented in the previous sections.

Social scientists who construct SEMs face many problems - among others, what variables to measure, how to construct measurement scales, how to remove outliers, and how to transform the variables. Often the ultimate goal of SEM construction is to achieve understanding of the causal relations among the variables, or to estimate the coefficients in an underlying structural equation model. As we have demonstrated in 4.4, in the absence of very strong background knowledge, regression is not a reliable technique for either of these purposes. This leaves the social scientist with the task of selecting SEMs from among a vast array of possibilities.

SEM selection can be thought of as a search problem - it is a search among a space of SEMs for the simplest SEMs that are compatible with background knowledge and fit the data. The search problem is very difficult because of measurement error, non-random samples, missing values, etc. However, here we wish to concentrate on the problems of the sheer size of the search space, and the existence of many plausible alternatives.

Even if latent variables are excluded, the search space is enormous (the number of different SEMs grows super-exponentially with the number of variables.) If latent variables are allowed, the number of possible models becomes infinite. Of course background knowledge, such as time order, can vastly reduce the search space. Nevertheless, even given background knowledge, the number of a priori plausible alternatives is often orders of magnitude too large to search by hand.

The problem is made even more difficult by the need to find not just one good SEM in the search space, but all of the good SEMs. As we showed in 4.1, there are often many SEMs that have the same p-value (as well as the same BIC and AIC scores.) Although all of these models receive the same scores, they can produce very different estimates of underlying parameters, and represent very different theories of the causal relations among the variables. In the absence of background knowledge to distinguish among these alternatives, it is important to present *all* of the simplest alternatives compatible with the background knowledge and data, rather to simply arbitrarily choose one. This suggests that the proper output of a search procedure should at least include a set of covariance equivalent SEMs compatible with background knowledge, rather than a single SEM.

Our approach to solving the problems of the large search space and the existence of many SEMs that may receive a high BIC score has been to search the space of PAGs, rather

than searching the space of SEMs with latent variables. One advantage of searching the space of PAGs is that for a fixed number of observed variables, there are a finite number of PAGs, but an infinite number of latent variable SEMs (since, in theory, one could add an arbitrary number of latent variables). In addition, if a PAG represents a SEM that receives a high BIC score, it also represents all of the SEMs that are covariance equivalent to that SEM; hence a search algorithm that outputs a PAG is not making an arbitrary choice among a set of covariance equivalent SEMs⁹. Further, while it is known that the BIC score is an $O(1)$ approximation of the posterior for a PAG, it is not known whether this is the case for latent variable SEMs. Finally, it is much easier to calculate a BIC score for a PAG than it is for many latent variable SEMs.

The FCI algorithm takes as input a covariance matrix, distributional assumptions, and background knowledge (e.g. time order), and outputs a PAG. The search proceeds by performing a sequence of conditional independence tests. In the large sample limit, the search is guaranteed to be correct under assumptions described in Spirtes et al., 1993. In the worst case (many direct connections among the observed variables, or many latent common causes of pairs of observed variables) the time it takes to perform the search grows exponentially as the number of variables. However, in some cases, it can perform searches on up to 100 measured variables. How large a set of SEMs is represented by the output depends upon what the true SEM is (if such exists), and how many latent variables it contains. (When the FCI algorithm is restricted to the case where there are no latent variables, the algorithm may be simplified, in which case it is called the PC algorithm.) See Spirtes, et al. (1993) for details.

We have also devised a greedy BIC score algorithm, that at each stage makes the one change to the PAG that most improves the score of the PAG. The greedy BIC score algorithm takes as input a covariance matrix, distributional assumptions, and background knowledge (e.g. time order), and outputs a set of PAGs, along with their BIC scores. See Spirtes, Richardson and Meek (1996) for details. (Instructions for downloading and using a program, TETRAD II, that contains both the FCI algorithm and the greedy BIC score algorithm can be found on the world wide web at <http://hss.cmu.edu/philosophy/TETRAD/tetrad.html>.)

There are a number of uses of the PAGs or set of PAGs output by these search procedures. They can be used to answer some, but not all, questions about the effects of

⁹ While the set of SEMs represented by a PAG is not too small in the sense that if it represents a SEM it also represents all of the SEMs covariance equivalent to it, it is larger than strictly necessary in that it generally does not contain a single covariance equivalence class. While this does not affect the correctness of the output, it does mean that the output is less informative than is theoretically possible.

interventions upon causal systems. In addition, as we will illustrate in the next section, they can be used as a starting point for selecting a particular latent variable SEM. See Spirtes et al. 1993 for details.

6. Applications of PAG Searches

6.1. Foreign Investment

The first example illustrates how the PC and FCI algorithms can be used to generate alternative models which cast doubt upon conclusions drawn from a regression. Timberlake and Williams (1984) used regression to claim foreign investment in third-world countries promotes dictatorship. They measured political exclusion (PO) (i.e., dictatorship), foreign investment penetration in 1973 (FI), energy development in 1975 (EN), and civil liberties (CV) for 72 countries. Civil liberties was measured on an ordered scale from 1 to 7, with lower values indicating greater civil liberties.

Their inference is unwarranted. Their model (with the relations between the regressors omitted) and the model obtained from the PC algorithm using a .12 significance level to test for vanishing partial correlations) are shown in Figure 1.¹⁰ (Because the algorithm performs a sequence of tests, one cannot read the reliability of the algorithm off of the significance level. We typically run the algorithms at a variety of different significance levels, and compare the results to see if any of the features of the output are constant.)

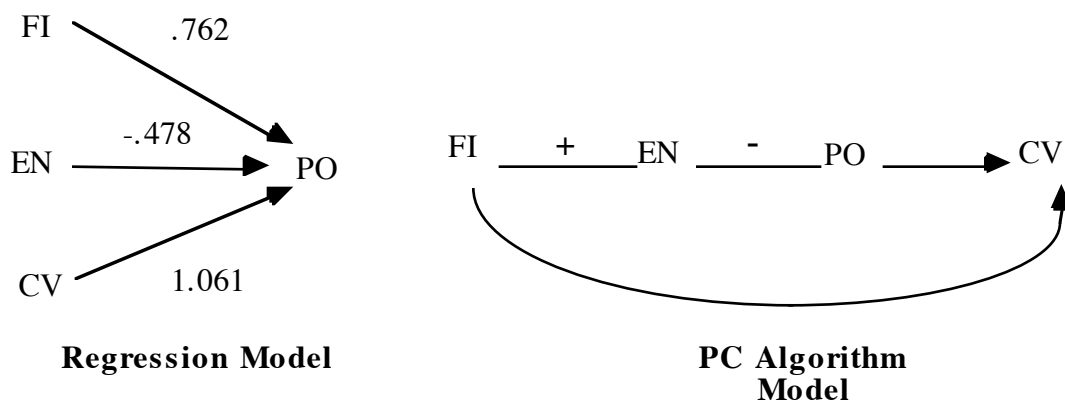


Figure 11

The PC Algorithm will not orient the FI-EN and EN-PO, edges, or determine whether they are due to at least one unmeasured common cause. Maximum likelihood estimates of any of the SEMs represented by the pattern output by the PC algorithm require that the

¹⁰Searches at lower significance levels remove the adjacency between FI and EN.

influence of FI on PO (if any) be negative, and the models easily pass a likelihood ratio test with the EQS program. If any of these SEMs is correct, Timberlake and William's regression model appears to be a case in which an effect of the outcome variable is taken as a regressor.

This analysis of the data assumes there are no unmeasured common causes. If we run the correlations through the FCI algorithm using the same significance level, we obtain a PAG that, together with the required signs of the dependencies, says that foreign investment and energy consumption have a common cause, as do foreign investment and civil liberties, that energy development has no influence on political exclusion, but political exclusion may have a negative effect on energy development, and that foreign investment has no influence, direct or indirect, on political exclusion.

Given the small sample size, and the uncertainty about the distributional assumptions, we do not present the alternative models suggested by the PC and FCI algorithms as particularly well-supported by the evidence. However, we do think that they are at least as well-supported as the regression model, and hence serve to cast doubt upon conclusions drawn from that model.

6.2. *Lead and IQ*

The next example shows how the FCI algorithm can be used to find a PAG, which can then be used as a starting point for a search for a latent variable DAG model. It also illustrates how such a procedure produces different results than simply applying regression or using regression to generate more sophisticated models, such as errors-in-variables models.

By measuring the concentration of lead in a child's baby teeth, Herbert Needleman was the first epidemiologist to even approximate a reliable measure of cumulative lead exposure. His work helped convince the United States to eliminate lead from gasoline and most paint (Needleman, et. al., 1979). In their 1985 article in *Science*, Needleman, Geiger and Frank gave results for a multivariate linear regression of children's IQ on lead exposure. Having started their analysis with almost 40 covariates, they were faced with a variable selection problem to which they applied backwards elimination regression, arriving at a final regression equation involving lead and five covariates. The covariates were measures of genetic contributions to the child's IQ (the parent's IQ), the amount of environmental stimulation in the child's early environment (the mother's education), physical factors that might compromise the child's cognitive endowment (the number of previous live births), and the parent's age at the birth of the child, which might be a proxy for many factors. The measured variables they used are as follows:

ciq	- child's verbal IQ score	piq	- parent's IQ scores
lead	- measured concentration in baby teeth	mab	- mother's age at child's birth
med	- mother's level of education in years	fab	- father's age at child's birth
nlb	- number of live births previous to the sampled child		

The standardized regression solution¹¹ is as follows, with t-ratios in parentheses. Except for fab, which is significant at 0.1, all coefficients are significant at 0.05, and R² = .271.

$$\hat{ciq} = \begin{matrix} \square & .143 & \square & .219 & \square & .247 & \square & .237 & \square & .204 & \square & .159 \\ & (2.32) & & (3.08) & & (3.87) & & (1.97) & & (1.79) & & (2.30) \end{matrix} \quad [1]$$

This analysis prompted criticism from Steve Klepper, an economist at Carnegie Mellon (see Klepper, 1988; Klepper, Kamlet, & Frank, 1993). Klepper correctly argued that Needleman's statistical model (a linear regression) neglected to account for measurement error in the regressors. That is, Needleman's measured regressors were in fact imperfect proxies for the actual but latent causes of variations in IQ, and in these circumstances a regression analysis gives a biased estimate of the desired causal coefficients and their standard errors. He constructed an errors-in-variables model to take into account the measurement error. See Figure 12, where the latent variables are in boxes, and the relations between the regressors are unconstrained.

Unfortunately, an errors-in-variables model that explicitly accounts for Needleman's measurement error is "underidentified," and thus cannot be estimated by classical techniques without making additional assumptions. Klepper, however, worked out an ingenious technique to bound the estimates, provided one could reasonably bound the amount of measurement error contaminating certain measured regressors (Klepper, 1988, 1993). The required measurement error bounds vary with each problem, however, and those required in order to bound the effect of actual lead exposure below 0 in Needleman's model seemed wholly unreasonable. Klepper concluded that the statistical evidence for Needleman's hypothesis was indeed weak. A Bayesian analysis, based on Gibbs sampling techniques, found that several posteriors corresponding to different priors lead to similar results. Although the size of the Bayesian point estimate for lead's influence on IQ moved

¹¹ The covariance data for this reanalysis was originally obtained from Needleman by Steve Klepper, who generously forwarded it. In this, and all subsequent analyses, the correlation matrix is used.

up and down slightly, its sign and significance (the 95% central region in the posterior over the lead-iq connection always included zero) were robust.

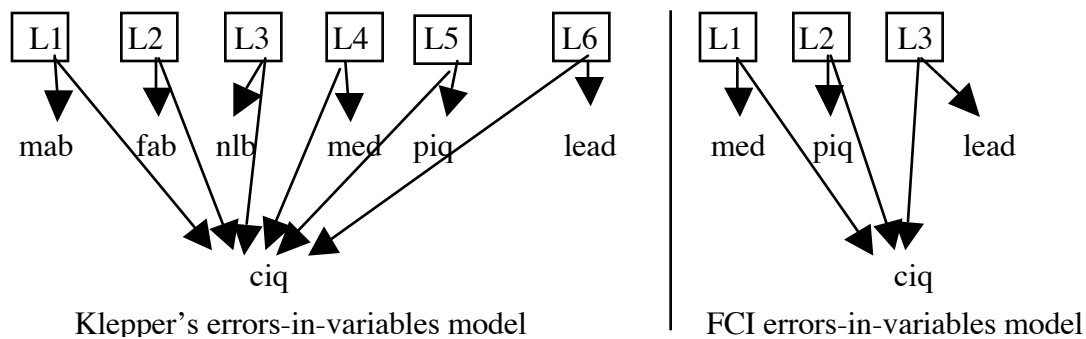


Figure 12

A reanalysis using the FCI algorithm produced different results. Scheines first used TETRAD II to first generate a PAG which was subsequently used as the basis for constructing an errors-in-variables model. The FCI algorithm produced a PAG that indicated that mab, fab, and nlb are *not* causes of ciq, contrary to Needleman’s regression.¹² If we construct an errors-in-variables model compatible with the PAG produced by the FCI algorithm, the model does not contain mab, fab, or nlb. See Figure 12. (We emphasize that there are other models compatible with the PAG which are not errors-in-variables models; the selection of an error-in-variables model from the set of models represented by the PAG is an assumption.) In fact the variables that the FCI algorithm eliminated were precisely those which required unreasonable measurement error assumptions in Klepper's analysis. With the remaining regressors, Scheines specified an errors-in-variables model to parameterize the effect of actual lead exposure on childrens’ IQ. This model is still underidentified but under several priors, nearly all the mass in the posterior was over negative values for the effect of actual lead exposure--now a latent variable--on measured IQ. In addition, applying Klepper’s bounds analysis to this model indicated that the effect of actual lead exposure on iq was bounded above zero given reasonable assumptions about the degree of measurement error.

¹² The fact that mab had a significant regression coefficient indicates that mab and ciq are correlated conditional on the other variables; the FCI algorithm concluded that mab is not a cause of ciq because mab and ciq are unconditionally uncorrelated. See Spirtes et al. 1993 for an explanation of the FCI algorithm in more detail.

7. Appendix

We will prove Theorem 1 and Theorem 2 in two steps. First we will prove them for the case where $G(M)$ contains no double-headed arrows; then for the case where $G(M)$ does contain double-headed arrows.

A probability measure P over \mathbf{V} satisfies the **global directed Markov property** for path diagram G if and only if for any three disjoint sets of variables \mathbf{X} , \mathbf{Y} , and \mathbf{Z} included in \mathbf{V} , if \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} , then \mathbf{X} is independent of \mathbf{Y} given \mathbf{Z} in P .

The following lemma relates the global directed Markov property to factorizations of a density function. Denote a density function over \mathbf{V} by $f(\mathbf{V})$, where for any subset \mathbf{X} of \mathbf{V} , $f(\mathbf{X})$ denotes the marginal of $f(\mathbf{V})$. Let $\mathbf{An}(\mathbf{X})$ be the set of ancestors of members of \mathbf{X} . If $f(\mathbf{V})$ is the density function for a probability measure over a set of variables \mathbf{V} , say that $f(\mathbf{V})$ **factors according to directed graph** G with vertices \mathbf{V} if and only if for every subset \mathbf{X} of \mathbf{V} ,

$$f(\mathbf{An}(\mathbf{X})) = \prod_{V \in \mathbf{An}(\mathbf{X})} g_V(V, \mathbf{Parents}(V))$$

where g_V is a non-negative function.

Lemma 1 was proved in Lauritzen et al. (1990) for the acyclic case, and the proof carries over essentially unchanged for the cyclic case.

Lemma 1: If \mathbf{V} is a set of random variables with a probability measure P that has a density function $f(\mathbf{V})$ and $f(\mathbf{V})$ factors according to directed graph G , then P satisfies the global directed Markov property for G .

Lemma 2 was proved in Spirtes (1995) and Koster (1995).

Lemma 2: If M is a SEM, and $\{X\}$ and $\{Y\}$ are d-separated given \mathbf{Z} in directed graph $G(M)$, then $\Pr(X, Y, \mathbf{Z}) = 0$ in $\Pr(M)$.

Lemma 3 was proved in Spirtes (1995).

Lemma 3: For any directed graph G , if $\{X\}$ and $\{Y\}$ are not d-separated given \mathbf{Z} in $G(M)$, there is a SEM M , $G = G(M)$ and $\Pr(X, Y, \mathbf{Z}) \neq 0$ in $\Pr(M)$.

We will now show that Theorem 1 and Theorem 2 hold even when G contains double-headed arrows. Let the set of vertices in G be \mathbf{V} . For a given triple X , Y , and \mathbf{Z} , if $\{X\}$ is d-separated from $\{Y\}$ given \mathbf{Z} in $G(M)$ and $G(M)$ contains double-headed arrows, the

strategy is to convert M into another SEM $M'(M, X, Y, Z)$ such that $G(M'(M, X, Y, Z))$ has additional latent variables, but no double-headed arrows, the marginal over V of $\square(M'(M, X, Y, Z))$ is equal to $\square(M)$, and $\{X\}$ and $\{Y\}$ are d-separated given Z in $G(M'(M, X, Y, Z))$. (We write $M'(M, X, Y, Z)$ in order to emphasize that the SEM M' constructed from M is a function of the path diagram of M , and the vertices X , Y , and Z in the d-separation relation being considered.) It will then follow from Lemma 2 that $\square(X, Y, Z) = 0$ in $\square(M)$.

If $\{X\}$ is d-separated from $\{Y\}$ given Z in $G(M)$, the graph $G(M'(M, X, Y, Z))$ is constructed by the following algorithm, where a **trek** between X_i and X_j is an undirected path between X_i and X_j that contains no colliders. (We will illustrate the application of the algorithm to the path diagram in Figure 13.)

Algorithm: Construct Latent Directed Graph

Inputs - Path Diagram G with vertex set V , Vertices X, Y, Z ;

Output - Directed Graph $G_{\text{Construct}}(G, X, Y, Z)$, with vertex set $V \sqcup T$;

1. Order the variables so that X is first, Y is second, followed by each variable with a descendant in Z , followed by any remaining variables that have X or Y as descendants in $G(M)$, followed by the rest of the variables. Given this ordering, we will now refer to the variables as X_1, \dots, X_n , where for all i , X_i is the i^{th} variable in the ordering.

2. For each variable X_i , add to the existing graph G , a variable T_i , and edges from T_i to X_j , for each $j \geq i$. Call the resulting graph, which has vertex set $(X_1, \dots, X_n, T_1, \dots, T_n)$ $G_{\text{Construct}(0)}$.

3. Let $G_{\text{Construct}(i)}$ be the the graph constructed after the i^{th} iteration of the following step, starting with $i = 1$: If $r > i$, and there is no trek between X_r and X_i in $G_{\text{Construct}(i-1)}$ containing a variable T_j , where $j < i$, and \square and \square are uncorrelated in \square , then remove the $T_i \rightarrow X_r$ edge.

For inputs G , X , Y , and Z , we will refer to the output of this algorithm as $G_{\text{Construct}}(G, X, Y, Z)$. Note that it follows from step 2 of the construction algorithm that if there is a trek $X_i \rightarrow T_j \rightarrow X_k$, then $j \leq \min(i, k)$.

Suppose for the graph in Figure 13 we are interested in whether $\square(X, Y) = 0$ (i.e. $Z = \emptyset$).

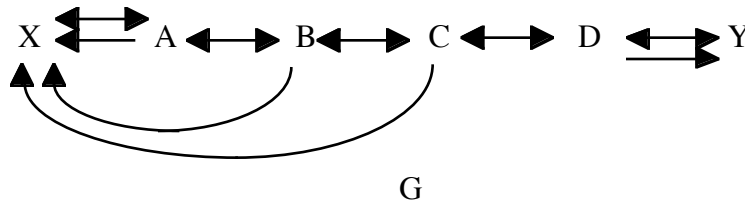


Figure 13

Applying the first step of Algorithm Construct Latent Directed Graph to G in Figure 13 with vertex inputs X, Y, \emptyset , results in the naming of the vertices shown in Figure 14.

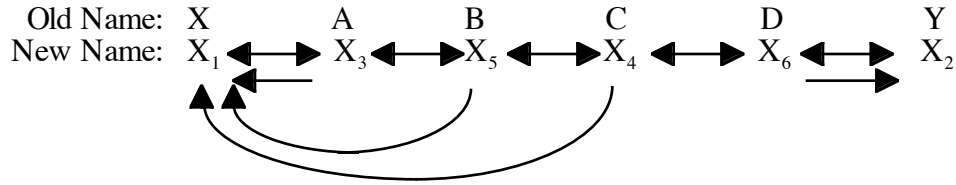


Figure 14: G with vertices renamed

Applying steps 2 and 3 of Algorithm Construct Latent Directed Graph results in the directed graph shown in Figure 15.

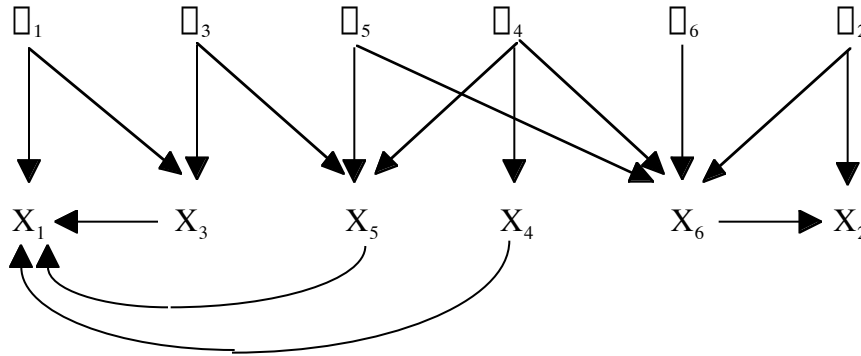


Figure 15: $G_{\text{Construct}}(X, Y, \emptyset)$

As an example of an application of step 3, the edge from T_3 to X_4 is removed because in $G_{\text{Construct}(2)}$ there is no trek between X_3 and X_4 that contains T_1 or T_2 , and there is no double-headed arrow between X_3 and X_4 in G .

The next series of lemmas shows how to construct a SEM $M'(M, X, Y, Z)$ with measured variables \mathbf{V} and latent variables \mathbf{T} , so that the marginal over \mathbf{V} of $\square(M'(M, X, Y, Z)) = \square(M)$, and $G(M'(M, X, Y, Z)) = G_{\text{Construct}}(G(M), X, Y, Z)$.

Lemma 4: If \square is a positive definite matrix, then there exists a positive definite matrix $\square' = \square - \square\mathbf{I}$, where \square is a real positive number.

Proof. Suppose that \square is a positive definite matrix. It follows then that for all solutions of $\det(\square - \square\mathbf{I}) = 0$, \square is positive. Let the smallest solution of $\det(\square - \square\mathbf{I}) = 0$ be \square_1 . Let \square be less than \square_1 and greater than 0. Let $\square' = \square - \square\mathbf{I}$. We will now show that all of the solutions of $\det(\square' - \square'\mathbf{I}) = 0$ are positive. $\square' - \square'\mathbf{I} = \square - \square\mathbf{I} - \square'\mathbf{I} = \square - (\square' + \square)\mathbf{I}$. If we set $\square'' = \square - \square$, then for each solution of $\det(\square - \square\mathbf{I}) = 0$, there is a solution of $\det(\square - (\square' + \square)\mathbf{I}) = 0$.

Since $\Sigma' = \Sigma - \Sigma_1$, and Σ_1 is less than Σ , the smallest solution of $\det(\Sigma' - \Sigma_1'I) = 0$ is greater than 0. \square

A linear transformation of a set of random variables is **lower triangular** if and only if there is an ordering of the variables such that the matrix representing the transformation is zero for all entries a_{ij} , when $j > i$.

Lemma 5: If X_1, \dots, X_n have a joint normal distribution $N(0, \Sigma)$, where Σ is positive definite, then there is a set of n mutually independent standard normal variables T_1, \dots, T_n , such that X_1, \dots, X_n are a lower triangular linear transformation of T_1, \dots, T_n and for each i , the coefficient of T_i in the equation for X_i is not equal to zero.

Proof. For every positive definite correlation matrix Σ , there is a SEM M with correlation matrix $\Sigma(M) = \Sigma$, and directed acyclic graph $G(M)$ that has each pair of vertices in $G(M)$ adjacent (Spirtes et al. 1993). The reduced form of a complete directed acyclic graph is a lower triangular transformation of independent error variables (in this case the T variables) that is non-zero on the diagonal, because Σ is positive definite. \square

There is a simple rule for calculating $\text{Cov}(X, Y)$ from a path diagram with no directed cycles that is used in the following lemmas. There is an edge directed **into** a vertex A on a path P if and only if P contains an edge $A \rightarrow B$ or $A \rightarrow B$. A vertex on a trek with no edges of the trek directed into it is called the **source** of a trek. Each trek has at most one source. (For example B is the source of the trek $A \rightarrow B \rightarrow C$, A is the source of $A \rightarrow B \rightarrow C$, and the trek $A \rightarrow B \rightarrow C \rightarrow D$ has no source.) Associated with each edge $A \rightarrow B$ in a graph is a label that corresponds to the coefficient of A in the equation for B , and associated with each edge $A \rightarrow B$ is a label that corresponds to the correlation of the error terms for A and B . $\text{Cov}(X, Y)$ is equal to the sum over all of the treks, of the product of the edge labels on the trek, times the variance of the source of the trek (if there is one). For example, in Figure 4, $\text{Cov}(Y, Z) = (\alpha\beta + \gamma)V(Z)$. For a proof of the case without correlated errors, see Glymour et al. (1987); the case with correlated errors is a simple modification of the latter proof.

Lemma 6: There is a SEM $M'(M, X, Y, Z)$ with measured variables \mathbf{V} and latent variables \mathbf{T} , such that $G(M'(M, X, Y, Z)) = G_{\text{Construct}}(G(M), X, Y, Z)$, and the marginal over \mathbf{V} of $\Sigma(M'(M, X, Y, Z))$ is equal to $\Sigma(M)$.

Proof. Let the correlation matrix among the error terms of M be Σ . If the equations in M are:

$$X_i = \sum_{j \neq i} b_{ij} X_j + \epsilon_i \quad (1)$$

(where some of the b_{ij} may equal zero, and some of the ϵ may be correlated) we will construct equations in $M'(M, X, Y, Z)$ that are:

$$X_i = \sum_{j \neq i} b_{ij} X_j + \sum_{j \in I} a_{ij} T_j + \epsilon'_i \quad (2)$$

by showing that there is a latent variable model of ϵ of the form

$$\epsilon_i = \sum_{j \in I} a_{ij} T_j + \epsilon''_i \quad (3)$$

where each of the T_i and ϵ''_i are uncorrelated.

By hypothesis, Σ is a positive definite matrix. By Lemma 4 there is a set of variables $\epsilon'_1, \dots, \epsilon'_n$ with positive definite matrix $\Sigma' = \Sigma - \Sigma I$, where $\Sigma > 0$. So we can write

$$\epsilon = \epsilon' + \epsilon'' \quad (4)$$

where the ϵ'' are uncorrelated with each other and the ϵ' variables, each ϵ'' is normally distributed with mean zero and variance Σ . The ϵ' variables will serve as the uncorrelated error terms in the new model that we construct; the ϵ variables are used only in intermediate stages of construction, and have the same covariance matrix as the ϵ variables, except that the variances of the variables have been decreased by a small amount Σ , i.e. $\Sigma' = \Sigma - \Sigma I$. As a first step to constructing a latent variable model of V , we will construct a latent variable model of ϵ .

By Lemma 5, there is a set of variables $T = \{T_1, \dots, T_n\}$ such that $\epsilon'_1, \dots, \epsilon'_n$ with correlation matrix Σ' are a lower triangular linear transformation of T_1, \dots, T_n and for each i , the coefficient of T_i in the equation for ϵ'_i is not equal to zero. That is

$$\epsilon'_i = \sum_{j \in I} a_{ij} T_j \quad (5)$$

where $a_{ii} \neq 0$.

There is a directed graph H that represents this latent variable model of the ϵ'_i variables, in which there is an edge from T_j to ϵ'_i only if $j \leq i$. From the construction of H , there are no edges from T_j to ϵ'_1 unless $j = 1$. Hence, for every $j \neq 1$, in H every trek between ϵ'_1 and ϵ'_j contains T_1 . It follows that there is at most one trek between ϵ'_1 and ϵ'_j . The edge from T_1 to ϵ'_1 is not zero. Hence it follows from the trek rule for calculating covariances from a path diagram, that if ϵ'_1 and ϵ'_j are not correlated in Σ' (i.e. there is no double-headed arrow between X_1 and X_j in $G(M)$) then the edge from T_1 to ϵ'_j is zero. (In the example from Figure 14, $a_{12} = a_{14} = a_{15} = a_{16} = 0$.)

Applying this strategy to each of the T_i variables in turn, we can now show that for each i and $r > i$, if there is no trek between \square_r and \square_i containing a variable T_j , where $j < i$, and \square_i and \square_j are uncorrelated in \square , then there is no $\square_i \square \square_r$ edge in H . Suppose on the contrary that in H there is no trek between \square_r and \square_i containing a variable T_j , where $j < i$, and \square and \square are uncorrelated in M , but the $T_i \square \square_r$ edge is in H . By the construction of H , if $k > i$, then there is no edge from T_k to \square_i . It follows that if in H there is no trek between \square_r and \square_i containing a variable T_j , where $j < i$, then every trek between \square_i and any other variable contains the edge from T_i to \square_i , which is in H since $a_{ii} \neq 0$. The $T_i \square \square_r$ edge exists by hypothesis, so there is exactly one trek between \square_i and \square_r in H . Hence, in every SEM L with vertices $\{\square_1, \dots, \square_n\}$ and directed graph $G(L) = H$, \square_i and \square_r are correlated in $\square(L)$. (Note that this could not be claimed if there were more than one trek between \square_i and \square_r since in that case the treks might cancel each other.) Since the covariances between distinct \square variables are equal to the correlations between the corresponding \square variables, it follows that \square and \square are correlated in \square , and hence there is a double-headed arrow between \square_i and \square_j in $G(M)$. This is a contradiction.

The graph H for the path diagram in Figure 14 is shown in Figure 16.

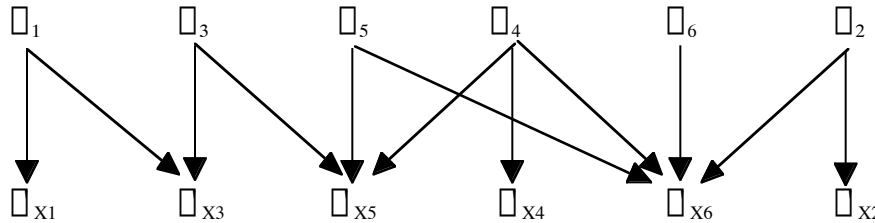


Figure 16: H

From the latent variable model of the \square variables, we can now form a model $M'(M, X, Y, Z)$ with measured variables \mathbf{V} and latent variables T_1, \dots, T_n , but without correlated errors.

$$X_i = \sum_{j \neq i} b_{ij} X_j + \sum_{j \in \mathbf{T}} a_{ij} T_j + \square'_i$$

It follows from equations (1), (4), and (5) that the marginal distribution of $\mathbf{V} = \{X_1, \dots, X_n\}$ in $M'(M, X, Y, Z)$ is the same as the distribution of \mathbf{V} in M .

We will now show that $G(M'(M, X, Y, Z)) = G_{\text{Construct}}(G(M), X, Y, Z)$. For variables A and B in \mathbf{V} , by the construction of M' , there is an edge between A and B in $G(M'(M, X, Y, Z))$ if and only if there is an edge between A and B in G , and hence an edge between A and B in $G_{\text{Construct}}(G(M), X, Y, Z)$. (Hence the ancestor relations among the variables in $G(M)$ are the same as the ancestor relations among the corresponding variables in $G(M'(M, X, Y, Z))$.) There is an edge between a variable T in \mathbf{T} and a variable A in \mathbf{V} in $G(M'(M, X, Y, Z))$ if and

only if there is an edge between T and \square_A in H . We have already shown that for each i and $r > i$, if there is no trek between \square_r and \square_i containing a variable T_j , where $j < i$, and \square_i and \square_j are uncorrelated in \square , then there is no $\square_i \square \square_r$ edge in H . It follows that for each i and $r > i$, if there is no trek between X_r and X_i containing a variable T_j , where $j < i$, and \square and \square are uncorrelated in \square , then there is no $T_i \square X_r$ edge in $G(M'(M,X,Y,Z))$. (This latter property is the property obtaining in $G_{\text{Construct}}(G(M),X,Y,Z)$ by application of steps 2 and 3.) Hence $G_{\text{Construct}}(G(M),X,Y,Z) = G(M'(M,X,Y,Z)) \square$

The next series of lemmas show that if X_1 and X_2 are d-separated given Z in $G(M)$, then X_1 and X_2 are d-separated given Z in $G(M'(M,X,Y,Z))$.

We will call a trek $X_j \square T_m \square X_i$ that contains a T variable a **latent trek** in $G_{\text{Construct}}(G(M),X,Y,Z)$. In $G(M)$, a **correlated error trek sequence** is a sequence of vertices $\langle X_1, \dots, X_k \rangle$ such that no pair of vertices adjacent in the sequence are identical, and for each consecutive pair of vertices X_r and X_s in the sequence, there is an edge $X_r \square X_s$. For example in Figure 13, the sequence of vertices $\langle X, A, B, C, D, Y \rangle$ is a correlated error trek sequence between X and Y .

Lemma 7: If there is a latent trek between X_i and X_j in $G_{\text{Construct}}(G(M),X,Y,Z)$ that contains a variable T_r , i.e. $X_i \square T_r \square X_j$, then in $G(M)$ there is a correlated error trek sequence between X_i and X_j , such that every variable in the correlated error trek sequence, with the possible exception of the endpoints, X_i and X_j , has index (i.e. subscript) less than or equal to r (henceforth referred to as the correlated error trek sequence in $G(M)$ corresponding to the latent trek between X_i and X_j in $G_{\text{Construct}}(G(M),X,Y,Z)$.)

Proof. The proof is by induction on r . Suppose first that $r = 1$. From the construction algorithm for $G_{\text{Construct}}(G(M),X,Y,Z)$, if there is a latent trek between X_i and X_j in $G_{\text{Construct}}(G(M),X,Y,Z)$ that contains T_1 then there are edges $X_i \square X_1$ and $X_j \square X_1$ in $G(M)$. The concatenation of these two edges forms a correlated error trek sequence in which (trivially) every variable in the sequence except for the endpoints has an index less than or equal to 1. The induction hypothesis is that for all $r \leq k$, if there is a latent trek between X_i and X_j in $G_{\text{Construct}}(G(M),X,Y,Z)$ that contains T_r , then in $G(M)$ there is a correlated error trek sequence between X_i and X_j , such that every variable in the sequence, with the possible exception of the endpoints has an index less than r . Suppose now that in $G_{\text{Construct}}(G(M),X,Y,Z)$ there is a latent trek between X_i and X_j such that the trek contains T_{k+1} , where $i, j \geq k+1$.

Suppose first that both $i, j > k+1$. Since the edge between T_{k+1} and X_i exists in $G_{\text{Construct}}(G(M),X,Y,Z)$, it follows from the construction algorithm for $G_{\text{Construct}}(G(M),X,Y,Z)$ that either there is a latent trek between X_i and X_{k+1} in $G_{\text{Construct}}(G(M),X,Y,Z)$ that contains some T_r , $r < k+1$, or there is a double-headed arrow between X_{k+1} and X_i in $G(M)$. In the

former case, by the induction hypothesis there is a correlated error trek sequence between X_i and X_{k+1} that, except for the endpoints, contains only vertices whose indices are less than or equal to r , and hence less than or equal to $k+1$. In the latter case, $\langle X_i, X_{k+1} \rangle$ is a correlated error trek sequence between X_i and X_{k+1} . Similarly, there is a correlated error trek sequence between X_{k+1} and X_j that, except for the endpoints, contains only vertices whose indices are less than or equal to $k+1$. These two correlated error trek sequences can be concatenated to form a correlated error trek sequence between X_i and X_j that, except for the endpoints, contains only vertices whose indices are less than or equal to $k+1$.

Suppose now that one either $i = k+1$ or $j = k+1$. Suppose without loss of generality that $j = k+1$. Since the edge between T_j and X_i exists in $G_{\text{Construct}}(G(M), X, Y, Z)$, it follows from the construction algorithm for $G_{\text{Construct}}(G(M), X, Y, Z)$ that either there is a latent trek between X_i and X_j in $G_{\text{Construct}}(G(M), X, Y, Z)$ that contains some T_r , $r < j$, or there is a double-headed arrow between X_j and X_i in $G(M)$. In either case there is a correlated error trek sequence between X_i and X_j that, except for the endpoints, contains only vertices whose indices are less than or equal to $k+1$. \square

For $G_{\text{Construct}}(G(M), X, Y, Z)$ shown in Figure 15, there is a latent trek between $X_5 \square T_5 \square X_6$, and a corresponding correlated error trek sequence $\langle X_5, X_4, X_6 \rangle$ in the graph G in Figure 14.

We will make use of the following Lemma which is a simple extension to path diagrams with directed cycles of Lemma 3.3.1 in Spirtes *et al.* (1993). This Lemma allows us to concatenate ‘small’ d-connecting paths to form a larger d-connecting path. We say a path is **into** endpoint X if the path contains some edge $X \square Y$ or $X \square Y$.

Lemma 8: In a path diagram G over a set of vertices V , if:

- (a) Q is a sequence of vertices in V from A to B , $Q \equiv \langle A \equiv X_0, \dots, X_{n+1} \equiv B \rangle$, such that $\square i, \square i \leq n, X_i \neq X_{i+1}$ (the X_i are only *pairwise distinct*, i.e. not necessarily distinct),
 - (b) $Z \square V \setminus \{A, B\}$,
 - (c) P is a set of undirected paths such that
 - (i) \square for each pair of consecutive vertices in Q , X_i and X_{i+1} , there is a unique undirected path in P that d-connects X_i and X_{i+1} given $Z \setminus \{X_i, X_{i+1}\}$,
 - (ii) \square if some vertex X_k in Q , is in Z , then the paths in P that contain X_k as an endpoint collide at X_k , (i.e. all such paths are directed into X_k)
 - (iii) \square if for three vertices X_{k-1}, X_k, X_{k+1} occurring in Q the d-connecting paths in P between X_{k-1} and X_k , and X_k and X_{k+1} , collide at X_k then X_k has a descendant in Z ,
- then there is a path U in G that d-connects $A \equiv X_0$ and $B \equiv X_{n+1}$ given Z .

Note that we do not require that a vertex occur only once in Q . Hence one occurrence of a vertex in Q may be a collider, and another occurrence of the same vertex in Q may be a

non-collider. (We say that Y_k is a collider (non-collider) in \mathbf{Q} if the pair of consecutive paths in \mathbf{P} that contain Y_k as an endpoint collide (do not collide) at Y_k .)

Lemma 9: If $X_1 \equiv X$ and $X_2 \equiv Y$ are d-connected given \mathbf{Z} in the directed graph $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$, then X and Y are d-connected given \mathbf{Z} in the path diagram $G(M)$. ($G(M)$ has vertex set \mathbf{V} , $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$ has vertex set $\mathbf{V} \sqcup \mathbf{T}$, and $\{X, Y\} \sqsubseteq \mathbf{V}$.)

Proof. Suppose that there is an undirected path U that d-connects X_1 and X_2 given \mathbf{Z} in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$. We will prove that X and Y are d-connected given \mathbf{Z} in $G(M)$ by constructing a sequence of vertices \mathbf{Q} and a set \mathbf{P} of paths in G between pairs of consecutive vertices in \mathbf{Q} satisfying the conditions of Lemma 8.

Our first step will be to use U to construct a sequence \mathbf{Q}' and a set of paths \mathbf{P}' in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$ from which we will then construct \mathbf{P} and \mathbf{Q} . Intuitively, we form \mathbf{Q}' and \mathbf{P}' by breaking U into pieces, such that each latent trek occurs as a separate piece. More formally, form a sequence \mathbf{Q}' of vertices and an associated sequence \mathbf{P}' of paths in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$ with the following properties: (i) every vertex in \mathbf{Q}' is in \mathbf{V} and occurs on U ; (ii) no vertex occurs in \mathbf{Q}' more than once; (iii) if X_i occurs before X_j in \mathbf{Q}' , then X_i occurs before X_j on U ; (iv) if the subpath of U between X_i and X_j is a latent trek, $X_i \sqsubseteq \mathbf{T}_r \sqsubseteq X_j$, then X_i and X_j both occur in that order in \mathbf{Q}' . The path in \mathbf{P}' associated with a pair X_i and X_j of consecutive vertices in \mathbf{Q}' is the subpath of U between X_i and X_j . In the example in Figure 15, in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$ the d-connecting path between X_1 and X_2 given $\mathbf{Z} = \emptyset$ is $X_1 \sqsubseteq X_5 \sqsubseteq \mathbf{T}_4 \sqsubseteq X_6 \sqsubseteq X_2$, $\mathbf{Q}' = \langle X_1, X_5, X_6, X_2 \rangle$, and $\mathbf{P}' = \langle X_1 \sqsubseteq X_5, X_5 \sqsubseteq \mathbf{T}_4 \sqsubseteq X_6, X_6 \sqsubseteq X_2 \rangle$. In this example, there are no colliders in \mathbf{Q}' .

Because U is a path that d-connects X_1 and X_2 given \mathbf{Z} in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$, it is clear that the paths in \mathbf{P}' have the following properties in $G_{\text{Construct}}(G(M), X, Y, \mathbf{Z})$: (i) Each path in \mathbf{P}' d-connects its endpoints X_i and X_j given $\mathbf{Z} \setminus \{X_i, X_j\}$; (ii) if paths in \mathbf{P}' collide at X_i then X_i has a descendant in \mathbf{Z} ; and (iii) if X_i is in \mathbf{Z} then the paths in \mathbf{P}' collide at X_i .

We will now show how to construct a sequence of vertices \mathbf{Q} and a set \mathbf{P} of paths in $G(M)$ between pairs of consecutive vertices in \mathbf{Q} satisfying the conditions of Lemma 8; it follows then that X and Y are d-connected given \mathbf{Z} in G .

We will create \mathbf{Q} by several modifications of \mathbf{Q}' . Step (1) in creating \mathbf{Q} is to replace each subsequence $\langle X_r, X_s \rangle$ of \mathbf{Q}' such that X_r and X_s are the endpoints of a latent trek in \mathbf{P}' , with the corresponding correlated error trek sequence $\langle X_r, \square, X_s \rangle$ in $G(M)$. Then replace the latent trek in \mathbf{P}' with the corresponding correlated error trek sequence in \mathbf{P}' . Note that each occurrence of X_k between $\langle X_r, \square, X_s \rangle$ is a collider in \mathbf{Q} . In the example, after the first step $\mathbf{Q} = \langle X_1, X_5, X_4, X_6, X_2 \rangle$ and $\mathbf{P} = \langle X_1 \sqsubseteq X_5, X_5 \sqsubseteq X_4, X_4 \sqsubseteq X_6, X_6 \sqsubseteq X_2 \rangle$, i.e. we replaced the subsequence $\langle X_5, X_6 \rangle$ in \mathbf{Q}' by $\langle X_5, X_4, X_6 \rangle$, and the latent trek $X_5 \sqsubseteq \mathbf{T}_4 \sqsubseteq X_6$ by $X_5 \sqsubseteq X_4$ and $X_4 \sqsubseteq X_6$ in \mathbf{Q}' .

Recall that the ancestor relations among the variables in \mathbf{V} (which includes the variables in \mathbf{Z}) in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$ are the same as the ancestor relations among the variables in $G(\mathbf{M})$. After stage (1) in creating \mathbf{Q} , if X_k is not an ancestor of \mathbf{Z} in $G(\mathbf{M})$ (or in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$), but has an occurrence in \mathbf{Q} that is a collider, it follows that X_k was added to \mathbf{Q} by replacing a subsequence $\langle X_r, X_s \rangle$ of \mathbf{Q}' by a corresponding correlated error trek sequence $\langle X_r, \square, X_s \rangle$ in $G(\mathbf{M})$. Hence any such X_k lies between some pair of vertices X_r and X_s that are adjacent in \mathbf{Q}' . Because every vertex in $\langle X_r, \square, X_s \rangle$ in \mathbf{Q} (except for X_r and X_s) has an index less than r and s , and X_k is not an ancestor of \mathbf{Z} in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$, it follows from the ordering of the variables that we chose, that X_r and X_s are not ancestors of \mathbf{Z} in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$. If a path U d -connects X_1 and X_2 given \mathbf{Z} , then every vertex on U is an ancestor of X_1 or X_2 or \mathbf{Z} . Because X_r and X_s are on U , but not ancestors of \mathbf{Z} in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$, and U d -connects X_1 and X_2 given \mathbf{Z} , X_r and X_s are both ancestors of $\{X_1, X_2\}$. Because in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$, X_r and X_s are both ancestors of $\{X_1, X_2\}$, and $k < r$ and s , it follows from the ordering of the variables that X_k is also an ancestor of $\{X_1, X_2\}$ in $G_{\text{Construct}}(G(\mathbf{M}), \mathbf{X}, \mathbf{Y}, \mathbf{Z})$. Hence X_k is an ancestor of $\{X_1, X_2\}$ in $G(\mathbf{M})$. In the example, in $G(\mathbf{M})$, X_4 is not an ancestor of the empty set but is an ancestor of X_1 , and it is between two vertices X_5 and X_6 which also are not ancestors of the empty set but are ancestors of X_1 or X_2 .

Thus, if there is some vertex X_k in \mathbf{Q} that is not an ancestor of \mathbf{Z} , but occurs in \mathbf{Q} as a collider then X_k is an ancestor of X_1 or X_2 . Let X_a be the last occurrence of a collider in \mathbf{Q} that is an ancestor of X_1 but not of \mathbf{Z} , if there is one, otherwise let $X_a = X_1$. Step (2) in forming \mathbf{Q} and \mathbf{P} is to replace the subsequence $\langle X_1, \square, X_a \rangle$ by $\langle X_1, X_a \rangle$ if $X_a \neq X_1$, and replacing the corresponding paths in \mathbf{P} by a directed path from X_a to X_1 if $X_a \neq X_1$. (Such a directed path exists if $X_a \neq X_1$ because X_a is an ancestor of X_1 .) This removes all occurrences of vertices between X_1 and X_a that are not ancestors of \mathbf{Z} , but are colliders in \mathbf{Q} . In the example, $X_a = X_4$, and after step 2, $\mathbf{Q} = \langle X_1, X_4, X_6, X_2 \rangle$ and $\mathbf{P} = \langle X_1 \square X_4, X_4 \square X_6, X_6 \square X_2 \rangle$.

By definition, every vertex that occurs as a collider between X_a and X_2 in \mathbf{Q} is an ancestor of \mathbf{Z} or of X_2 . Let X_b be the first vertex after X_a in \mathbf{Q} that is an ancestor of X_2 but not of \mathbf{Z} , if there is one, otherwise let $X_b = X_2$. Step (3) in forming \mathbf{Q} and \mathbf{P} is to replace the subsequence $\langle X_b, \square, X_2 \rangle$ by $\langle X_b, X_2 \rangle$ if $X_b \neq X_2$, and replacing the corresponding paths in \mathbf{P} by a directed path from X_b to X_2 if $X_b \neq X_2$. This removes all occurrences of colliders between X_b and X_2 that are not ancestors of \mathbf{Z} . Note that all occurrences of colliders that are left are between X_a and X_b , and every occurrence of a collider between X_a and X_b is an ancestor of \mathbf{Z} by construction. In the example, $X_b = X_2$, and after step (3), \mathbf{Q} and \mathbf{P} are unchanged.

We will now show that every path between a pair of variables X_u and X_v in \mathbf{P} d-connects X_u and X_v given $\mathbf{Z} \setminus \{X_u, X_v\}$. If the path between X_u and X_v is also in \mathbf{P}' , then it d-connects X_u and X_v given $\mathbf{Z} \setminus \{X_u, X_v\}$ because every path in \mathbf{P}' has this property. If the path between X_u and X_v is not in \mathbf{P}' , but was added in step (1) of the formation of \mathbf{P} , then the path between X_u and X_v is a correlated error trek $X_u \square X_v$, which clearly d-connects X_u and X_v given $\mathbf{Z} \setminus \{X_u, X_v\}$. If the path between X_u and X_v is not in \mathbf{P}' , but was added in step (2) of the formation of \mathbf{P} , then $X_u = X_1$, $X_v = X_a$, and the path between X_u and X_v is a directed path from X_a to X_1 that does not contain any member of \mathbf{Z} since either $X_a = X_1$ or X_a is not an ancestor of \mathbf{Z} . Hence the path d-connects X_u and X_v given \mathbf{Z} . Similarly, if the path between path between X_u and X_v is not in \mathbf{P}' , but was added in step (3) of the formation of \mathbf{P} , then $X_u = X_b$, $X_v = X_2$, and the path between X_u and X_v is a directed path from X_b to X_2 that does not contain any member of \mathbf{Z} . Hence the path d-connects X_u and X_v given \mathbf{Z} .

We will now show that every vertex that occurs as a collider in \mathbf{Q} has a descendant in \mathbf{Z} , and every vertex that occurs as a non-collider in \mathbf{Q} is not in \mathbf{Z} . Every vertex that occurs as a collider in \mathbf{Q} is an ancestor of \mathbf{Z} , because steps (2) and (3) in the formation of \mathbf{Q} removed all occurrences of colliders that were not ancestors of \mathbf{Z} . Every vertex that occurs as a non-collider in \mathbf{Q}' and as a non-collider in \mathbf{Q} is not in \mathbf{Z} , because every vertex that occurs as a non-collider in \mathbf{Q} is not in \mathbf{Z} . The only vertices that may occur as non-colliders in \mathbf{Q} but not in \mathbf{Q}' are X_a and X_b . X_a is not in \mathbf{Z} , because either it is equal to X_1 or X_2 , neither of which is in \mathbf{Z} , or it is not an ancestor of \mathbf{Z} by construction. Similarly, X_b is not in \mathbf{Z} .

Hence \mathbf{Q} is a sequence of paths that satisfy properties (i), (ii), and (iii) of Lemma 8. It follows from Lemma 8 that $X_1 \equiv X$ and $X_2 \equiv Y$ are d-connected given \mathbf{Z} in $G(\mathbf{M})$. \square

Theorem 1: If \mathbf{M} is a SEM, and $\{X\}$ and $\{Y\}$ are d-separated given \mathbf{Z} in $G(\mathbf{M})$, then $\square(X, Y, \mathbf{Z}) = 0$ in $\square(\mathbf{M})$.

Proof. By Lemma 6 and Lemma 9 there is a SEM $\mathbf{M}'(\mathbf{M}, X, Y, \mathbf{Z})$ with the marginal of $\square(\mathbf{M}'(\mathbf{M}, X, Y, \mathbf{Z})) = \square(\mathbf{M})$, and $\{X\}$ and $\{Y\}$ d-separated given \mathbf{Z} in $G(\mathbf{M}'(\mathbf{M}, X, Y, \mathbf{Z})) = G_{\text{Construct}}(\mathbf{M}, X, Y, \mathbf{Z})$. Because $G_{\text{Construct}}(\mathbf{M}, X, Y, \mathbf{Z})$ is the directed graph of a latent variable model $\mathbf{M}'(\mathbf{M}, X, Y, \mathbf{Z})$ with correlation matrix that has marginal $\square(\mathbf{M})$, no correlated errors, and X and Y are d-separated given \mathbf{Z} in $G_{\text{Construct}}(\mathbf{M}, X, Y, \mathbf{Z})$, it follows from Lemma 2 that $\square(X, Y, \mathbf{Z}) = 0$ in \square . \square

Theorem 2: If $\{X_i\}$ and $\{X_j\}$ are not d-separated given \mathbf{Z} in path diagram G , then there is a SEM \mathbf{M} such that $G(\mathbf{M}) = G$, and $\square(X_i, X_j, \mathbf{Z}) \neq 0$ in $\square(\mathbf{M})$.

Proof. Suppose that $\{X_i\}$ and $\{X_j\}$ are d-connected given \mathbf{Z} in G , and the set of vertices in G is \mathbf{V} . Form a graph Transform(G) with vertices $\mathbf{T} \subseteq \mathbf{V}$ in the following way. For a pair of vertices X_k and X_m in \mathbf{V} , there is a directed edge $X_k \square X_m$ in Transform(G) if and only if

there is a directed edge $X_k \rightarrow X_m$ in G . For vertices X_k and X_m in \mathbf{V} , there is a vertex $T(X_k, X_m)$ in \mathbf{T} , and edges $X_m \rightarrow T(X_k, X_m) \rightarrow X_k$ if and only if there is a double-headed arrow $X_k \leftrightarrow X_m$ in G . (For convenience in writing equations, for each latent variable $T(X_k, X_m)$ in $\text{Transform}(G)$, we will also refer to it as $T(X_m, X_k)$.)

For $\{X_i, X_j\} \perp\!\!\!\perp \mathbf{Z} \perp\!\!\!\perp \mathbf{V}$, if $\{X_i\}$ and $\{X_j\}$ are d-connected given \mathbf{Z} in G , then they are d-connected given \mathbf{Z} in $\text{Transform}(G)$. By Lemma 3 there is a SEM M' , with $G(M') = \text{Transform}(G(M))$, and $\perp\!\!\!\perp(X_i, X_j, \mathbf{Z}) \neq \perp\!\!\!\perp$.

Let $\text{Double}(X_i)$ be the set of vertices X_m in G such that there is an edge $X_m \leftrightarrow X_i$ in G . In M' ,

$$X_i = \prod_{X_m \in \text{Parents}(X_i)} a_{im} X_m + \prod_{X_m \in \text{Double}(X_i)} b_{ij} T(X_i, X_m) + \varGamma_i$$

Now define

$$\varGamma_i = \prod_{X_m \in \text{Double}(X_i)} b_{ij} T(X_i, X_m) + \varGamma_i$$

It follows then that

$$X_i = \prod_{X_m \in \text{Parents}(X_i)} a_{ij} X_m + \varGamma_i$$

is a SEM M , with $G(M) = G$, and $\perp\!\!\!\perp(X_i, X_j, \mathbf{Z}) \neq \perp\!\!\!\perp$ in $\perp\!\!\!\perp(M)$. \square

We will prove Theorem 5 before Theorem 3 because we will use Theorem 5 in the proof of Theorem 3.

Theorem 5: If G_1 and G_2 are path diagrams that are covariance equivalent over \mathbf{O} , then G_1 and G_2 are d-separation equivalent over \mathbf{O} .

Proof. Suppose that G_1 and G_2 are not d-separation equivalent over \mathbf{O} . Suppose without loss of generality that there is some $\{X\}$, $\{Y\}$ and \mathbf{Z} included in \mathbf{O} , such that $\{X\}$ and $\{Y\}$ are d-connected given \mathbf{Z} in G_1 , but not in G_2 . By Theorem 2, there is some SEM M with $G(M) = G_1$ such that $\perp\!\!\!\perp(X, Y, \mathbf{Z}) \neq \perp\!\!\!\perp$. By Theorem 1, there is no SEM M' with $G(M') = G_2$, in which $\perp\!\!\!\perp(X, Y, \mathbf{Z}) \neq \perp\!\!\!\perp$. Hence G_1 and G_2 are not covariance equivalent over \mathbf{O} . \square

Let $\text{Ancestors}^*(X, G)$ be the set of ancestors of X , excluding X , in directed graph G , and $\text{Descendants}^*(X, G)$ be the set of descendants of X excluding X in G .

Lemma 10: In a directed acyclic graph G , if X and Y are not adjacent, Y is not an ancestor of X , $\text{Ancestors}^*(Y, G) \setminus \{X\} \perp\!\!\!\perp \mathbf{Z}$, and $\mathbf{Z} \perp\!\!\!\perp \text{Descendants}^*(Y, G) = \emptyset$, then $\{X\}$ and $\{Y\}$ are d-separated given \mathbf{Z} .

Proof. Suppose that X and Y are not adjacent, but there is a path U that d-connects $\{X\}$ and $\{Y\}$ given \mathbf{Z} . Suppose that U contains an edge $A \rightarrow Y$. Then $\text{Ancestors}^*(Y, G) \setminus \{X\} \perp\!\!\!\perp \mathbf{Z}$, so $A \perp\!\!\!\perp \mathbf{Z}$. Since A is a non-collider on U ($A \neq X$), it follows that U does not d-connect

$\{X\}$ and $\{Y\}$ given Z , contrary to hypothesis. Suppose then that U contains an edge $A \rightarrow Y$. It follows that U contains a collider, because by hypothesis, Y is not an ancestor of X . Let C be the collider on U closest to Y . C is a descendant of Y , so $\text{Descendants}(C,G) \subseteq \text{Descendants}^*(Y,G)$. Hence $\text{Descendants}(C,G) \cap Z = \emptyset$, so again in this case U does not d-connect $\{X\}$ and $\{Y\}$ given Z . \square

Theorem 3: If G_1 and G_2 are directed acyclic graphs, G_1 and G_2 are covariance equivalent if and only if G_1 and G_2 are d-separation equivalent

Proof. By Theorem 5, if G_1 and G_2 are covariance equivalent G_1 and G_2 are d-separation equivalent.

Suppose that G_1 and G_2 are d-separation equivalent, and M is a SEM with directed acyclic graph $G(M) = G_1$. We can form a SEM M'' where $G(M'')$ is a subgraph of G_2 ¹³ and $\square(M'') = \square(M)$ in the following way.

Order the variables in G_2 so that X comes before Y in the ordering only if X is not a descendant of Y . Form a directed acyclic graph G_2' that has G_2 as a subgraph by putting an edge between X and Y if and only if X precedes Y in the ordering. This is proportional to cov of Y conditional on parents. In G_2 , Y independent of X conditional on parents. So coefficient is zero.

Because G_2' is a complete graph there is a SEM M' such that $G(M')$ is a subgraph of G_2' , and $\square(M') = \square(M)$. In any SEM M' with graph $G(M')$ that is a subgraph of G_2' , the error term for a variable Y is independent of the parents of Y . Hence if X is a parent of Y in $G(M')$, the regression coefficient of X when Y is regressed on its parents in $G(M')$ using $\square(M)$ is equal to the linear coefficient of X in the equation for Y in M' . If X is an ancestor of Y in G_2' , there is no edge from X to Y in $G(M')$ if and only if the linear coefficient of X in the equation for Y is 0 in M' .

Suppose X and Y are not adjacent in G_2 . In G_2 , either X is not an ancestor of Y or Y is not an ancestor of X ; suppose without loss of generality that Y is not an ancestor of X . Then X and Y are d-separated given $\text{Parents}(Y,G_2)$ in G_2 . By Lemma 10, in G_2 , $\{X\}$ and $\{Y\}$ are d-separated given $\text{Parents}(Y,G_2')$, because $\text{Parents}(Y,G_2')$ contains all of the ancestors of Y in G_2 , and no descendants of Y in G_2 . Because G_1 and G_2 are d-separation equivalent, $\{X\}$ and $\{Y\}$ are d-separated given $\text{Parents}(Y,G_2')$ in G_1 . By Theorem 1, $\square(X,Y,\text{Parents}(Y,G_2')) = 0$ in $\square(M)$. Hence the regression coefficient of X when Y is regressed on $\text{Parents}(Y,G_2')$ using $\square(M)$, is equal to 0. It follows that there is no edge between X and Y in $G(M')$. Hence $\square(M') = \square(M)$, and $G(M')$ is a subgraph of G_2 . \square

¹³ Note this includes the possibility that $G(M'') = G_2$.

Bibliography

Andersson, S., Madigan, D., and Perlman, M. (1995) A Characterization of Markov Equivalence Classes for Acyclic Digraphs, Technical Report 287, Department of Statistics, University of Washington.

Blalock, H., 1961, Causal Inferences in Nonexperimental Research, (W. W. Norton and Co., New York).

Bollen, K., 1989, Structural Equations with Latent Variables. (Wiley, New York).

Chickering, D. (1995) A Transformational Characterization of Equivalent Bayesian Network Structures, Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, Philippe Besnard and Steve Hanks (Eds.), Morgan Kaufmann Publishers, Inc., San Mateo, CA.

Frydenberg, M., 1990, The chain graph Markov property, Scandinavian Journal of Statistics, **17**, 333-353.

Geiger, D., and Pearl, J., 1988, Logical and Algorithmic properties of Conditional Independence. Technical Report R-97, Cognitive Systems Laboratory, University of California, Los Angeles.

C. Glymour, R. Scheines, P. Spirtes, and K. Kelly (1987) Discovering Causal Structure: Artificial Intelligence, Philosophy and Statistical Modeling, Academic Press, San Diego, CA.

C. Glymour, P. Spirtes, and R. Scheines (1994), In Place of Regression (in Patrick Suppes: Scientific Philosopher, Paul Humphreys (editor), Vol. 1, Kluwer Academic Publishers, Dordrecht, Holland.

Goldberger, A., Duncan, O. (eds.), 1973, Structural Equation Models in the Social Sciences (Seminar Press, New York).

Haavelmo, T., 1943, The statistical implications of a system of simultaneous equations, *Econometrica*, **11**, 1-12.

Kiiveri, H. and Speed, T., 1982, Structural analysis of multivariate data: A review, *Sociological Methodology*, Leinhardt, S. (ed.). Jossey-Bass, San Francisco.

Kiiveri, H., Speed, T., and Carlin, J., 1984, Recursive causal models, *Journal of the Australian Mathematical Society*, **36**, 30-52.

Klepper, S. (1988). Regressor diagnostics for the classical errors-in-variables model. *Journal of Econometrics*, **37**, 225-250.

Klepper, S., Kamlet, M., and Frank, R. (1993) Regressor Diagnostics for the Errors-in-Variables Model - An Application to the Health Effects of Pollution, *Journal of Environmental Economics and Management*. **24**, 190-211.

Koster, J., (1995) Markov Properties of Non-Recursive Causal Models, *Annals of Statistics*, November 1995.

Lee, S., and Hershberger, S. (1990). A simple rule for generating equivalent models in covariance structure modeling. *Multivariate Behavioral Research*, 25, 313-334.

Lauritzen, S., Dawid, A., Larsen, B., Leimer, H., 1990, Independence properties of directed Markov fields, *Networks*, **20**, 491-505.

Meek, C. (1995) Causal inference and causal explanation with background knowledge, *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, Philippe Besnard and Steve Hanks (Eds.), Morgan Kaufmann Publishers, Inc., San Mateo, CA, pp. 403-410.

Needleman, H., Geiger, S., and Frank, R. (1985). "Lead and IQ Scores: A Reanalysis," *Science*, 227, pp. 701-704.

Pearl, J., (1986) Fusion, propagation, and structuring in belief networks, *Artificial Intelligence* **29**, 241-88.

Pearl, J., (1988). *Probabilistic Reasoning in Intelligent Systems*, (Morgan Kaufman: San Mateo, CA).

Pearl, J. and Verma, T. (1991). A theory of inferred causation, in *Principles of Knowledge Representation and Reasoning: Proceedings of the Second International Conference* (Morgan Kaufmann, San Mateo, CA).

Pearl, J. (1997). *Graphs, Causality and Structural Equation Models*, Technical Report R-253, Cognitive Science Laboratory, UCLA.

Raftery, A. (1995), *Baysian Model Selection in social research*, in Marsden, ed. *Sociological Methodology*, Blackwells, Cambridge MA.

Richardson, T. (1994) *Properties of Cyclic Graphical Models*, Master's Thesis, Carnegie Mellon University.

Richardson, T. (1996a). A Polynomial Algorithm for Deciding Equivalence in Directed Cyclic Graphical Models. Technical Report PHIL-63, Department of Philosophy, Carnegie Mellon University.

Richardson, T. (1996b) A Discovery Algorithm for Directed Cyclic Graphs. In *Uncertainty in Artificial Intelligence: Proceedings of the Twelfth Conference* (F.Jensen and E.Horvitz, eds.), 462-469, Morgan Kaufmann, San Francisco.

Richardson, T. (1996c) *Feedback Models: Interpretation and Discovery*. Ph.D. Thesis, Dept. of Philosophy, Carnegie-Mellon University.

Scheines, R. (1994) Causation, Indistinguishability, and Regression, in *Sofstat '93: Advances in Statistical Software 4*, Frank Faulbaum (editor), Gustav Fischer Verlag, pp. 89-98.

Scheines, R. (1997) Estimating Latent Causal Influences: TETRAD II Model Selection and Bayesian Parameter Estimation, in Proceedings of the 6th International Workshop on Artificial Intelligence and Statistics.

Spirtes, P. and Glymour, C., 1990, Causal Structure Among Measured Variables Preserved with Unmeasured Variables. Technical Report CMU-LCL-90-5, Laboratory for Computational Linguistics, Carnegie Mellon University.

Spirtes, P., and Glymour, C.(1991) An algorithm for fast recovery of sparse causal graphs, *Social Science Computer Review*, **9**, 62-72.

Spirtes, P., Verma, T. (1992) Equivalence of Causal Models with Latent Variables. Technical Report CMU-PHIL-33, Department of Philosophy, Carnegie Mellon University, October, 1992.

Spirtes, P., Glymour, C., and Scheines, R.(1993) Causation, Prediction, and Search, (Springer-Verlag Lecture Notes in Statistics 81, New York).

Spirtes, P., (1995) Directed Cyclic Graphical Representation of Feedback Models, in Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, ed. by Philippe Besnard and Steve Hanks, Morgan Kaufmann Publishers, Inc., San Mateo.

Spirtes, P., Richardson, T., Meek, C., Scheines, R., and Glymour, C., (1996). Using D-separation to Calculate Zero Partial Correlations in Linear Models with Correlated Errors, Technical Report CMU-72-Phil.

Spirtes, P., and Richardson, T. (1996), A Polynomial Time Algorithm For Determining DAG Equivalence in the Presence of Latent Variables and Selection Bias, Proceedings of the 6th International Workshop on Artificial Intelligence and Statistics.

Spirtes, P., Richardson, T., and Meek, C. (1996). Heuristic Greedy Search Algorithms for Latent Variable Models, Proceedings of the 6th International Workshop on Artificial Intelligence and Statistics.

Spirtes, P., Richardson, T., Meek, C. (1997). The Dimensionality of Mixed Ancestral Graphs, Technical Report CMU-83-Phil.

Stelzl, I. (1986) Changing causal relations without changing the fit: Some rules for generating equivalent LISREL-models. *Multivariate Behavioral Research*, **21**, 309-331.

Timberlake, M. and Williams, K. (1984). Dependence, political exclusion, and government repression: Some cross-national evidence. *American Sociological Review* **49**, 141-146.

Verma, T. and Pearl, J. (1990b). Equivalence and synthesis of causal models in Proc. Sixth Conference on Uncertainty in AI. Association for Uncertainty in AI, Inc., Mountain View, CA.

Whittaker, J.,1990, Graphical Models in Applied Multivariate Statistics (Wiley, New York).

Wright, S. (1934). The method of path coefficients, *Annals of Mathematical Statistics* **5**, 161-215.