

Causal, Casual and Curious

Judea Pearl*

Linear Models: A Useful “Microscope” for Causal Analysis

Abstract: This note reviews basic techniques of linear path analysis and demonstrates, using simple examples, how causal phenomena of non-trivial character can be understood, exemplified and analyzed using diagrams and a few algebraic steps. The techniques allow for swift assessment of how various features of the model impact the phenomenon under investigation. This includes: Simpson’s paradox, case–control bias, selection bias, missing data, collider bias, reverse regression, bias amplification, near instruments, and measurement errors.

Keywords: structural equation model, linear model, path analysis

*Corresponding author: **Judea Pearl**, Department of Computer Science, University of California Los Angeles, Los Angeles, CA, USA, E-mail: judea@cs.ucla.edu

1 Introduction

Many concepts and phenomena in causal analysis were first detected, quantified, and exemplified in linear structural equation models (SEM) before they were understood in full generality and applied to nonparametric problems. Linear SEM’s can serve as a “microscope” for causal analysis; they provide simple and visual representation of the causal assumptions in the model and often enable us to derive close-form expressions for quantities of interest which, in turns, can be used to assess how various aspects of the model affect the phenomenon under investigation. Likewise, linear models can be used to test general hypotheses and to generate counter-examples to over-ambitious conjectures.

Despite their ubiquity, however, techniques for using linear models in that capacity have all but disappeared from the main SEM literature, where they have been replaced by matrix algebra on the one hand, and software packages on the other. Very few analysts today are familiar with traditional methods of path tracing [1–4] which, for small problems, can provide both intuitive insight and easy derivations using elementary algebra.

This note attempts to fill this void by introducing the basic techniques of path analysis to modern researchers, and demonstrating, using simple examples, how concepts and issues in modern causal analysis can be understood and analyzed in SEM. These include: Simpson’s paradox, case–control bias, selection bias, collider bias, reverse regression, bias amplification, near instruments, measurement errors, and more.

2 Preliminaries

2.1 Covariance, regression, and correlation

We start with the standard definition of variance and covariance on a pair of variables X and Y . The *variance* of X is defined as

$$\sigma_x^2 = E[X - E(x)]^2$$

and measures the degree to which X deviates from its mean $E(X)$.

The *covariance* of X and Y is defined as

$$\sigma_{xy} = E[X - E(x)][Y - E(Y)]$$

and measures the degree to which X and Y covary.

Associated with the covariance, we define two other measures of association: (1) the regression coefficient β_{yx} and (2) the correlation coefficient ρ_{yx} . The relationships between the three is given by the following equations:

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad [1]$$

$$\beta_{yx} = \frac{\sigma_{xy}}{\sigma_x^2} = \frac{\sigma_y}{\sigma_x} \rho_{xy} \quad [2]$$

We note that $\rho_{xy} = \rho_{yx}$ is dimensionless and confined to the unit interval; $0 \leq \rho_{xy} \leq 1$. The regression coefficient, β_{yx} , represents the slope of the least square error line in the prediction of Y given X

$$\beta_{yx} = \frac{\partial}{\partial X} E(Y | X = x)$$

2.2 Partial correlations and regressions

Many questions in causal analysis concern the change in a relationship between X and Y conditioned on a given set Z of variables. The easiest way to define this change is through the *partial regression coefficient* $\beta_{yx \cdot z}$ which is given by

$$\beta_{yx \cdot z} = \frac{\partial}{\partial X} E(Y | X = x, Z = z)$$

In words, $\beta_{yx \cdot z}$ is the slope of the regression line of Y on X when we consider only cases for which $Z = z$. The partial correlation coefficient $\rho_{xy \cdot z}$ can be defined by normalizing $\beta_{yx \cdot z}$:

$$\rho_{xy \cdot z} = \beta_{yx \cdot z} \sigma_{x \cdot z} / \sigma_{y \cdot z}.$$

A well-known result in regression analysis [5] permits us to express $\rho_{xy \cdot z}$ recursively in terms of pair-wise regression coefficients. When Z is singleton, this reduction reads:

$$\rho_{yx \cdot z} = \frac{\rho_{yx} - \rho_{yz} \rho_{xz}}{[(1 - \rho_{xz}^2)(1 - \rho_{yz}^2)]^{\frac{1}{2}}} \quad [3]$$

Accordingly, we can also express $\beta_{yx \cdot z}$ and $\sigma_{yx \cdot z}$ in terms of pair-wise relationships, which gives:

$$\sigma_{yx \cdot z} = \sqrt{\sigma_{xx} - \sigma_{xz}^2 / \sigma_z^2} \sqrt{\sigma_{yy} - \sigma_{yz}^2 / \sigma_z^2} \rho_{yx \cdot z} \quad [4]$$

$$\sigma_{yx \cdot z} = \sigma_x^2 [\beta_{yx} - \beta_{yz} \beta_{zx}] = \sigma_{yx} - \frac{\sigma_{yz} \sigma_{zx}}{\sigma_z^2} \quad [5]$$

$$\beta_{yx \cdot z} = \frac{\beta_{yx} - \beta_{yz}\beta_{zx}}{1 - \beta_{zx}^2 \sigma_x^2 / \sigma_z^2} = \frac{\sigma_z^2 \sigma_{yx} - \sigma_{yz} \sigma_{zx}}{\sigma_x^2 \sigma_z^2 - \sigma_{zx}^2} = \frac{\sigma_y \rho_{yx} - \rho_{yz} \cdot \rho_{zx}}{\sigma_x \cdot 1 - \rho_{zx}^2} \tag{6}$$

Note that none of these conditional associations depends on the level z at which we condition variable Z ; this is one of the features that makes linear analysis easy to manage and, at the same time, limited in the spectrum of relationships it can capture.

2.3 Path diagrams and structural equation models

A linear structural equation model (SEM) is a system of linear equations among a set V of variables, such that each variable appears on the left hand side of at most one equation. For each equation, the variable on its left hand side is called the *dependent* variable, and those on the right hand side are called *independent* or *explanatory* variables. For example, the equation below

$$Y = \alpha X + \beta Z + U_Y \tag{7}$$

declares Y as the dependent variable, X and Z as explanatory variables, and U_Y as an “error” or “disturbance” term, representing all factors omitted from V that, together with X and Z determine the value of Y . A structural equation should be interpreted as an assignment process, that is, to determine the value of Y , nature consults the value of variables X, Z , and U_Y and, based on their linear combination in eq. [7], assigns a value to Y .

This interpretation renders the equality sign in eq. [7] non-symmetrical, since the values of X and Z are not determined by inverting eq. [7] but by other equations, for example,

$$X = \gamma Z + U_X \tag{8}$$

$$Z = U_Z \tag{9}$$

The directionality of this assignment process is captured by a *path-diagram*, in which the nodes represent variables, and the arrows represent the non-zero coefficients in the equations. The diagram in Figure 1(a) represents the SEM equations of eqs. [7]–[9] and the assumption of zero correlations between the U variables,

$$\sigma_{U_X, U_Y} = \sigma_{U_X, U_Z} = \sigma_{U_Z, U_Y} = 0$$

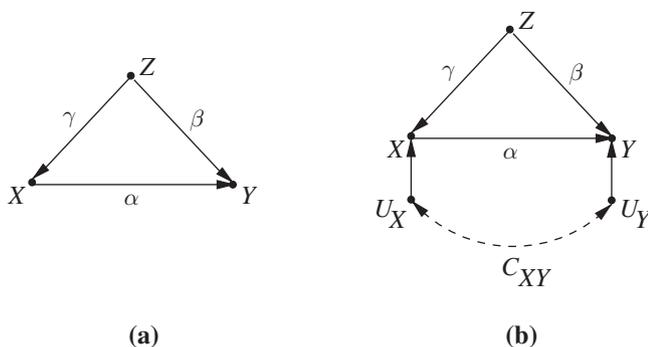


Figure 1 Path diagrams capturing the directionality of the assignment process of eqs. [7]–[9] as well as possible correlations among omitted factors.

The diagram in Figure 1(b) on the other hand represents eqs. [7]–[9] together with the assumption

$$\sigma_{U_X, U_Z} = \sigma_{U_Z, U_Y} = 0$$

while $\sigma_{U_X, U_Y} = C_{XY}$ remains undetermined.

The coefficients α, β , and γ are called *path coefficients*, or *structural parameters* and they carry causal information. For example, α stands for the change in Y induced by raising X one unit, while keeping all other variables constant.¹

The assumption of linearity makes this change invariant to the levels at which we keep those other variables constant, including the error variables; a property called “effect homogeneity.” Since errors (e.g., U_X, U_Y, U_Z) capture variations among individual units (i.e., subjects, samples, or situations), effect homogeneity amounts to claiming that all units react equally to any treatment, which may exclude applications with profoundly heterogeneous subpopulations.

2.4 Wright’s path-tracing rules

In 1921, the geneticist Sewall Wright developed an ingenious method by which the covariance σ_{xy} of any two variables can be determined swiftly, by mere inspection of the diagram [1]. Wright’s method consists of equating the (standardized²) covariance $\sigma_{xy} = \rho_{xy}$ between any pair of variables to the sum of products of path coefficients and error covariances along all d -connected paths between X and Y . A path is d -connected if it does not traverse any collider (i.e., head-to-head arrows, as in $X \rightarrow Y \leftarrow Z$).

For example, in Figure 1(a), the standardized covariance σ_{xy} is obtained by summing α with the product $\beta\gamma$, thus, yielding $\sigma_{xy} = \alpha + \beta\gamma$, while in Figure 1(b) we get $\sigma_{xy} = \alpha + \beta\gamma + C_{XY}$. Note that for the pair (X, Z) , we get $\sigma_{xz} = \gamma$ since the path $X \rightarrow Y \leftarrow Z$ is not d -connected.

The method above is valid for standardized variables, namely variables normalized to have zero mean and unit variance. For non-standardized variables the method need to be modified slightly, multiplying the product associated with a path p by the variance of the variable that acts as the “root” for path p . For example, for Figure 1(a) we have $\sigma_{xy} = \sigma_x^2\alpha + \sigma_z^2\beta\gamma$, since X serves as the root for path $X \rightarrow Y$ and Z serves as the root for $X \leftarrow Z \rightarrow Y$. In Figure 1(b), however, we get $\sigma_{xy} = \sigma_x^2\alpha + \sigma_z^2\beta\gamma + C_{XY}$ where the double arrow $U_X \leftrightarrow U_Y$ serves as its own root.

2.5 Reading partial correlations from path diagrams

The reduction from partial to pair-wise correlations summarized in eqs. [4]–[6], when combined with Wright’s path-tracing rules permits us to extend the latter so as to read partial correlations directly from the diagram. For example, to read the partial regression coefficient $\beta_{yx:z}$, we start with a standardized model where all variances are unity (hence, $\sigma_{xy} = \rho_{xy} = \beta_{xy}$), and apply eq. [6] with $\sigma_x = \sigma_z = 1$ to get:

$$\beta_{yx:z} = \frac{(\sigma_{yx} - \sigma_{yz}\sigma_{zx})}{(1 - \sigma_{xz}^2)} \quad [10]$$

¹ Readers familiar with do -calculus [6] can interpret α as the experimental slope $\alpha = \frac{\partial}{\partial x} E[(Y|do(x), do(z))]$ while those familiar with counterfactual logic can write $\alpha = \frac{\partial}{\partial x} Y_{xz}(u)$. The latter implies the former, and the two coincide in linear models, where causal effects are homogeneous (i.e., unit-independent.)

² Standardized parameters refer to systems in which (without loss of generality) all variables are normalized to have zero mean and unit variance, which significantly simplifies the algebra.

At this point, each pair-wise covariance can be computed from the diagram through path-tracing and, substituted in eq. [10], yields an expression for the partial regression coefficient $\beta_{yx \cdot z}$.

To witness, the pair-wise covariances for Figure 1(a) are:

$$\sigma_{yx} = \alpha + \beta\gamma \quad [11]$$

$$\sigma_{xz} = \gamma \quad [12]$$

$$\sigma_{yz} = \beta + \alpha\gamma \quad [13]$$

Substituting in (10), we get

$$\begin{aligned} \beta_{yx \cdot z} &= [(\alpha + \beta\gamma) - (\beta + \gamma\alpha)\gamma]/(1 - \gamma^2) \\ &= \alpha(1 - \gamma^2)/(1 - \gamma^2) \\ &= \alpha \end{aligned} \quad [14]$$

Indeed, we know that, for a confounding-free model like Figure 1(a) the direct effect α is identifiable and given by the partial regression coefficient $\beta_{yx \cdot z}$. Repeating the same calculation on the model of Figure 1(b) yields:

$$\beta_{yx \cdot z} = \alpha + C_{XY}$$

leaving α non-identifiable.

Armed with the ability to read partial regressions, we are now prepared to demonstrate some peculiarities of causal analysis.

3 The microscope at work: examples and their implications

3.1 Simpson's paradox

Simpson's paradox describes a phenomenon whereby an association between two variables reverses sign upon conditioning on a third variable, regardless of the value taken by the latter. The history of this paradox and the reasons it evokes surprise and disbelief are described in Chapter 6 of [7]. The conditions under which association reversal appears in linear models can be seen directly in Figure 1(a). Comparing eqs. [12] and [14] we obtain

$$\beta_{yx} = \alpha + \beta\gamma \quad \beta_{yx \cdot z} = \alpha$$

Thus, if α has a different sign from $\beta\gamma$, it is quite possible to have the regression of Y on X , β_{yx} , change sign upon conditioning on $Z = z$, for every z . The magnitude of the change depends on the product $\beta\gamma$ which measures the extent to which X and Y are confounded in the model.

3.2 Conditioning on intermediaries and their proxies

Conventional wisdom informs us that, in estimating the effect of one variable on another, one should not adjust for any covariate that lies on the pathway between the two [8]. It took decades for epidemiologists to discover that similar prohibition applies to proxies of intermediaries [9]. The amount of bias introduced by such adjustment can be assessed from Figure 2.

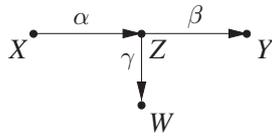


Figure 2 Path diagram depicting an intermediate variable (Z) and its proxy (W). Conditioning on W would distort the regression of Y on X .

Here, the effect of X on Y is simply $\alpha\beta$ as is reflected by the regression slope $\beta_{yx} = \alpha\beta$. If we condition on the intermediary Z , the regression slope vanishes, since the equality $\sigma_{yx} = \alpha\beta = \sigma_{yz}\sigma_{zx}$ renders β_{xy-z} zero in eq. [10]. If we condition on a proxy W of Z , eq. [10] yields

$$\beta_{yx-w} = \frac{\beta_{yx} - \beta_{yw}\beta_{wx}}{1 - \beta_{wx}^2} = \frac{\alpha\beta - \beta\gamma\alpha}{1 - \alpha^2\gamma^2} = \frac{\alpha\beta(1 - \gamma^2)}{1 - \alpha^2\gamma^2} \quad [15]$$

which unveils a bias of size

$$\beta_{yx-z} - \alpha\beta = \alpha\beta\gamma^2(1 - \alpha^2)/(1 - \alpha^2\gamma^2)$$

As expected, the bias disappears for $\gamma = 0$ and intensifies for $\gamma = 1$, where conditioning on W amounts to suppressing all variations in Z .

Speaking of suppressing variations, the model in Figure 3

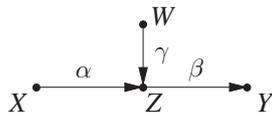


Figure 3 Conditioning on W does not distort the regression of Y on X .

may carry some surprise. Conditioning on W in this model also suppresses variations in Z , especially for high γ and, yet, it introduces no bias whatsoever; the partial regression slope is [eq. 10]:

$$\beta_{yx-w} = \frac{\sigma_{yx} - \sigma_{yw}\sigma_{xw}}{1 - \sigma_{xw}^2} = \frac{\alpha\beta - 0}{1 - 0} = \alpha\beta \quad [16]$$

which is precisely the causal effect of X on Y . It seems as though no matter how tightly we “clamp” Z by controlling W , the causal effect of X on Y remains unaltered. Appendix I explains this counter-intuitive result.

3.3 Case-control bias

In the last section, we explained the bias introduced by conditioning on an intermediate variable (or its proxy) as a restriction on the flow of information between X and Y . This explanation is not entirely satisfactory, as can be seen from the model of Figure 4.

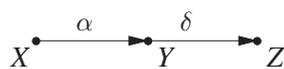


Figure 4 Conditioning on Z , a descendant of Y , biases the regression of Y on X .

Here, Z is not on the pathway between X and Y , and one might surmise that no bias would be introduced by conditioning on Z , but analysis dictates otherwise. Path tracing combined with eq. [10] gives:

$$\begin{aligned} \beta_{yx \cdot z} &= (\sigma_{yx} - \sigma_{yz}\sigma_{zx}) / (1 - \sigma_{xz}^2) \\ &= (\alpha - \delta^2\alpha) / (1 - \alpha^2\delta^2) \\ &= \alpha(1 - \delta^2) / (1 - \alpha^2\delta^2) \end{aligned} \tag{17}$$

and yields the bias

$$\beta_{yx \cdot z} - \alpha = \alpha\delta^2(\alpha^2 - 1) / (1 - \alpha^2\delta^2) \tag{18}$$

This bias reflects what economists called “selection bias” [10] and epidemiologists “case-control bias” [11], which occurs when only patients for whom the outcome Y is evidenced (e.g., a complication of a disease) are counted in the database. An intuitive explanation of this bias (invoking virtual colliders) is given in [7, p. 339]. In contrast, conditioning on a proxy of the explanatory variable X , as in Figure 5, introduces no bias, since

$$\beta_{yx \cdot z} = \frac{(\sigma_{yx} - \sigma_{yz}\sigma_{zx})}{(1 - \sigma_{xz}^2)} = \frac{a - (ab)b}{1 - b^2} = a \tag{19}$$

This can also be deduced from the conditional independence $Z \perp\!\!\!\perp Y|X$ which is implied by the diagram in Figure 5, but not in Figure 4. However, to assess the size of the induced bias, as we did in eq. [18], requires an algebraic analysis of path tracing.



Figure 5 Conditioning on Z , a descendant of X , does not bias the regression of Y on X .

3.4 Sample selection bias

The two examples above are special cases of a more general phenomenon called “selection bias” which occurs when samples are preferentially selected to the data set, depending on the values of some variables in the model [12–15]. In Figure 6, for example, if

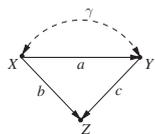


Figure 6 Conditioning on $Z = 1$ represents inclusion in the dataset and biases the regression of Y on X , unless $c = 0$.

$Z = 1$ represents the inclusion in the data set, and $Z = 0$ exclusion, the selection decision is shown to be a function of both X and Y . Since inclusion ($Z = 1$) amounts to conditioning on Z , we may ask what the regression of Y on X is in the observed data, $\beta_{yx,z}$, compared with the regression in the entire population, $\beta_{yx} = a + \gamma$.

Applying our path-tracing analysis in eq. [10] we get:

$$\beta_{yx,z} = \frac{\sigma_{yx} - \sigma_{yz}\sigma_{zx}}{1 - \sigma_{zx}^2} = \frac{(a + \gamma) - [(a + \gamma)b + c][b + (a + \gamma)c]}{1 - [b + (a + \gamma)c]^2} = \frac{(a + \gamma)[1 - b^2 - c^2] - bc}{1 - [b + (a + \gamma)c]^2}. \tag{20}$$

We see that a substantial bias may result from conditioning on Z , persisting even when X and Y are not correlated, namely when $\sigma_{xy} = a + \gamma = 0$. Note also that the bias disappears for $c = 0$, as in Figure 5, but not for $b = 0$, which returns us to the case-controlled model of Figure 4.

Selection bias is symptomatic of a general phenomenon associated with conditioning on collider nodes (Z in our example). The phenomenon involves spurious associations induced between two causes upon observing their common effect, since any information refuting one cause should make the other more probable. It has been known as Berkson Paradox [16], “explaining away” [17] or simply “collider bias.”³

3.5 Missing data

In contrast to selection bias, where exclusion ($S = 0$) removes an entire unit from the dataset, in missing data problems a unit may have each of its variables masked independently of the others [20, p. 89]. Therefore, the diagram representing the missingness process should assign each variable V_i a “switch” R_i , called “missingness mechanism” which determines whether V_i is observed ($R_i = 0$) or masked ($R_i = 1$). The arrows pointing to R_i tells us which variables determine whether R_i fires ($R_i = 1$) or not ($R_i = 0$). In Figure 7(a), for example, the missingness of X , denoted R_x , depends only on the latent variable L , while the missingness of Y is shown to depend on both L and X .

Assume we wish to estimate the covariance σ_{xy} from partially observed data generated by the model of Figure 7(a); can we obtain an unbiased estimate of σ_{xy} ? The question boils down to expressing σ_{xy} in

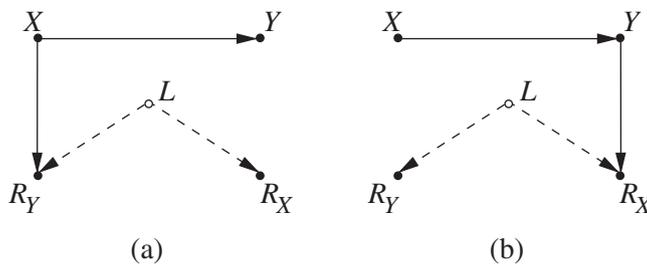


Figure 7 Missingness diagrams in which conditioning on $R_x = 0$ or $R_y = 0$ represents unmasking the values of X and Y , respectively. The parameter σ_{xy} , σ_x^2 and σ_y^2 can all be estimated bias-free from data generated by either model, through each model requires a different estimation procedure.

³ It has come to my attention recently, and I feel responsibility to make it public, that seasoned reviewers for highly reputable journals reject papers because they are not convinced that such bias can be created; it defies, so they claim, everything they have learned from statistics and economics. A typical resistance to accepting Berkson’s Paradox is articulated in [18, 19].

terms of the information available to us, namely the values of X and Y that are revealed to us whenever $R_x = 0$ or $R_y = 0$ (or both). If we simply estimate σ_{xy} from samples in which both X and Y are observed, that would amount to conditioning on both $R_x = 0$ and $R_y = 0$ which would introduce a bias since the pair $\{X, Y\}$ is not independent of the pair $\{R_x, R_y\}$ (owed to the unblocked path from Y to R_y).

The graph reveals, however, that σ_{xy} can nevertheless be estimated bias-free from the information available, using two steps. First, we note that X is independent of its missingness mechanism R_x , since the path from X to R_x is blocked (by the collider at R_y). Therefore, $\sigma_x^2 = (\sigma_x^2 | R_x = 0)$.⁴ This means that we can estimate σ_x from the samples in which X is observed, regardless of whether Y is missing. Next, we note that the regression slope β_{yX} can be estimated (e.g., using OLS) from samples in which both X and Y are observed. This is because conditioning on $R_x = 0$ and $R_y = 0$ is similar to conditioning on Z in Figure 5, where Z is a proxy of the explanatory variable X .

Putting the two together (using eq. [2]) we can write:

$$\sigma_{xy} = \sigma_x^2 \beta_{yX} = (\sigma_x^2 | R_x = 0) \cdot (\beta_{yX} | R_x = 0, R_y = 0)$$

which guarantees that the product of the two estimates on the right hand side would result in an unbiased estimate of σ_{xy} . Note that a similar analysis of Figure 7(b) would yield

$$\sigma_{xy} = (\sigma_y^2 | R_y = 0) \cdot (\beta_{yX} | R_x = 0, R_y = 0)$$

which instructs us to estimate σ_y^2 using samples in which Y is observed and estimate the regression of X on Y from samples in which both X and Y are observed. Remarkably, the two models are statistically indistinguishable and yet each dictates a different estimation procedure, thus demonstrating that no model-blind estimator can guarantee to deliver an unbiased estimate, even when such exists. If the path diagram permits no decomposition of σ_{xy} into terms conditioned on $R_x = 0$ and $R_y = 0$ (as would be the case, for example, if an arrow existed from X to R_x in Figure 7(a)) we would conclude then that σ_{xy} is not estimable by any method whatsoever. A general analysis of missing data problems using causal graphs is given in Mohan et al. [21].

3.6 The M -bias

The M -bias is another instant of Berkson’s paradox where the conditioning variable, Z , is a pre-treatment covariate, as depicted in Figure 8.

The parameters γ_1 and γ_2 represent error covariances C_{XZ} and C_{ZY} , respectively, which can be generated, for example, by latent variables effecting each of these pairs.

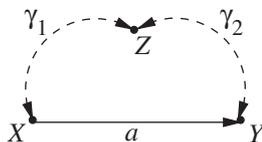


Figure 8 Adjusting for Z , which may be either pre-treatment or post-treatment covariate, introduces bias where none exists. The better the predictor the higher the bias.

⁴ $(\sigma_x^2 | R_x = 0)$ stands for the conditional variance of X given $R_x = 0$. We take the liberty of treating R_y as any other variable in the linear system, even though it is binary, hence the relationship $X \rightarrow R_y$ must be nonlinear. The linear context simplifies the intuition and the results hold in nonparametric systems as well.

To analyze the size of this bias, we apply eq. [10] and get:

$$\beta_{yx \cdot z} = \frac{a - (\gamma_2 + a\gamma_1)\gamma_1}{1 - \gamma_1^2} = a - \frac{\gamma_1\gamma_2}{1 - \gamma_1^2} \quad [21]$$

Thus, the bias induced increases substantially when γ_1 approaches one, that is, when Z becomes a good predictor of X . Ironically, this is precisely when investigators have all the textbook reasons to adjust for Z . Being pre-treatment, the collider Z cannot be distinguished from a confounder (as in Figure 1(a)) by any statistical means, and has alluded some statisticians to conclude that “there is no reason to avoid adjustment for a variable describing subjects before treatment” [22, p. 76].

3.7 Reverse regression

Is it possible that men would earn a higher salary than equally qualified women, and *simultaneously*, men are more qualified than women doing equally paying job? This counter-intuitive condition can indeed exist, and has given rise to a controversy called “Reverse Regression;” some sociologists argued that, in salary discrimination cases, we should not compare salaries of equally qualified men and women, but, rather, compare qualifications of equally paid men and women [23]. The phenomenon can be demonstrated in Figure 9.

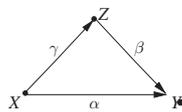


Figure 9 Path diagram in which Z acts as a mediator between X and Y , demonstrating negative reverse regression $\beta_{zx \cdot y}$ for positive α, β and γ .

Let X stand for gender (or age, or socioeconomic background), Y for job earnings and Z for qualification. The partial regression $\beta_{yx \cdot z}$ encodes the differential earning of males ($X = 1$) over females ($X = 0$) having the same qualifications ($Z = z$), while $\beta_{zx \cdot y}$ encodes the differential qualification of males ($X = 1$) over females ($X = 0$) earning the same salary (y).

For the model in Figure 9, we have

$$\beta_{yx \cdot z} = \alpha$$

$$\beta_{zx \cdot y} = (\sigma_{zx} - \sigma_{zy}\sigma_{yx}) / (1 - \sigma_{zy}^2) = [(\gamma - (\beta + \gamma\alpha)(\alpha + \beta\gamma))] / [1 - (\beta + \gamma\alpha)^2]$$

Surely, for any $\alpha > 0$ and $\beta > 0$ we can choose γ so as to make $\beta_{zx \cdot y}$ negative. For example, the combination $\alpha = \beta = 0.8$ and $\gamma = 0.1$ yields

$$\beta_{zx \cdot y} = [(0.1 - (0.8 + 0.1 \times 0.8)(0.8 + 0.8 \times 0.1)] / [1 - (0.8 + 0.1 \times 0.9)^2] = -5.8545$$

Thus, there is no contradiction in finding men earning a higher salary than equally qualified women, and simultaneously, men being more qualified than women doing equally paying job. A negative $\beta_{zx \cdot y}$ may be a natural consequence of male-favoring hiring policy ($\alpha > 0$), male-favoring training policy ($\gamma > 0$) and qualification-dependent earnings ($\beta > 0$).

The question of whether standard or reverse regression is more appropriate for proving discrimination is also clear. The equality $\beta_{yx \cdot z} = \alpha$ leaves no room for hesitation, because α coincides with the counterfactual

definition of “direct effect of gender on hiring had qualification been the same,” which is the court’s definition of discrimination.

The reason the reverse regression appeals to intuition is because it reflects a model in which the employer decides on the qualification needed for a job on the basis of both its salary level and the applicant sex. If this were a plausible model, it would indeed be appropriate to persecute an employer who demands higher qualifications from men as opposed to women. But such a model should place Z as a post-salary variable, for example, $X \rightarrow Z \leftarrow Y$.

3.8 Bias amplification

In the model of Figure 10, Z acts as an instrumental variable, since $\sigma_{zu} = 0$. If U is unobserved, however, Z cannot be distinguished from a confounder, as in Figure 1(a), in the sense that for every set of parameters (α, β, γ) in Figure 1(a) one can find a set (a, b, c, d) for the model in Figure 10 such that the observed covariance matrices of the two models are the same. This indistinguishability, together with the fact that Z may be a strong predictor of X may lure investigators to condition on Z to obtain an unbiased estimate of d [24]. Recent work has shown, however, that such adjustment would amplify the bias created by U [25–27]. The magnitude of this bias and its relation to the pre-conditioning bias, ab , can be computed from the diagram of Figure 10, as follows:

$$\beta_{yx \cdot z} = \frac{\sigma_{xy} - \sigma_{yz} \cdot \sigma_{xz}}{1 - \sigma_{xz}^2} = \frac{(ab + \gamma_0) - c\gamma_0 c}{1 - c^2} = \gamma_0 + \frac{ab}{1 - c^2} \tag{22}$$

We see the bias created, $\frac{ab}{(1-c^2)}$, is proportional to the pre-existing bias ab and increases with c ; the better Z predicts X , the higher the bias. An intuitive explanation of this phenomenon is given in Pearl [26]

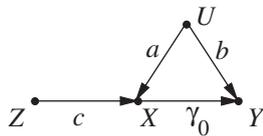


Figure 10 Bias amplification, $\beta_{yx \cdot z} - \gamma_0 > ab$, produced by conditioning on an instrumental variable (Z).

3.9 Near instruments – amplifiers or attenuators?

The model in Figure 11 is indistinguishable from that of Figure 10 when U is unobserved. However, here Z acts both as an instrument and as a confounder. Conditioning on Z is beneficial in blocking the confounding path $X \leftarrow Z \rightarrow Y$ and harmful in amplifying the baseline bias $cd + ab$. The trade-off between these two tendencies can be quantified by computing $\beta_{yx \cdot z}$, yielding

$$\begin{aligned} \beta_{yx \cdot z} &= \frac{\sigma_{xy} - \sigma_{yz}\sigma_{zx}}{1 - \sigma_{xz}^2} \\ &= \frac{\gamma_0 + cd + ab - (d + c\gamma_0)c}{1 - c^2} \\ &= \frac{\gamma_0(1 - c^2) + ab}{1 - c^2} \\ &= \gamma_0 + \frac{ab}{1 - c^2} \end{aligned} \tag{23}$$

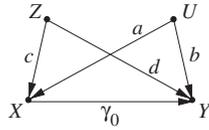


Figure 11 A diagram where Z acts both as an instrument and as a confounder.

We see that the baseline bias $ab + cd$ is first reduced to ab and then magnified by the factor $[(1 - c^2)^{-1}]$. For Z to be a bias-reducer, its effect on Y (i.e., d) must exceed its effect on X (i.e., c) by a factor $ab/(1 - c^2)$. This trade-off was assessed by simulations in Myers et al. [28] and analytically in Pearl [29], including an analysis of multi-confounders, and nonlinear models.

3.10 The butterfly

Another model in which conditioning on Z may have both harmful and beneficial effects is seen in Figure 12.

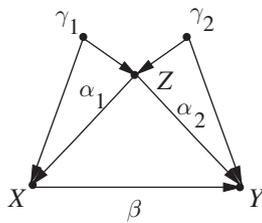


Figure 12 Adjusting for Z may be harmful or beneficial depending on the model's parameters.

Here, Z is both a collider and a confounder. Conditioning on Z blocks the confounding path through α_1 and α_2 and at the same time induces a virtual confounding path through the latent variables that create the covariances $C_{XZ} = \gamma_1$ and $C_{ZY} = \gamma_2$.

This trade-off can be evaluated from our path-tracing formula eq. [10] which yields

$$\begin{aligned} \beta_{yx \cdot z} &= \frac{\beta_{yx} - \beta_{yz}\beta_{zx}}{1 - \beta_{zx}^2} = \frac{[\beta + (\alpha_1 + \gamma_1)\alpha_2 + \alpha_1\gamma_2] - [\alpha_2 + \gamma_2 + \beta(\gamma_1 + \alpha_1)][\gamma_1 + \alpha_1]}{1 - (\alpha_1 + \gamma_1)^2} \\ &= \frac{\beta - \gamma_2\gamma_1 - \beta(\gamma_1 + \alpha_1)^2}{1 - (\alpha_1 + \gamma_1)^2} \end{aligned} \tag{24}$$

We first note that the pre-conditioning bias

$$\beta_{xy} - \beta = \alpha_2(\alpha_1 + \gamma_1) + \alpha_1\gamma_2 \tag{25}$$

may have positive or negative values even when both $\sigma_{xz} = 0$ and $\sigma_{zy} = 0$. This refutes folklore wisdom, according to which a variable Z can be exonerated from confounding considerations if it is uncorrelated with both treatment (X) and outcome (Y).

Second, we notice that conditioning on Z may either increase or decrease bias, depending on the structural parameters. This can be seen by comparing eq. [25] with the post-conditioning bias:

$$\beta_{xy \cdot z} - \beta = -\gamma_1\gamma_2/[1 - (\alpha_1 + \gamma_1)^2] \tag{26}$$

In particular, since eq. [26] is independent on α_2 , it is easy to choose values of α_2 that make eq. [25] either higher or lower than eq. [26].

3.11 Measurement error

Assume the confounder U in Figure 13(a) is unobserved but we can measure a proxy Z of U . Can we assess the amount of bias introduced by adjusting for Z instead of U ? The answer, again, can be extracted from our path-tracing formula, which yields

$$\begin{aligned} \beta_{yx \cdot z} &= \frac{\sigma_{yx} - \sigma_{yz}\sigma_{zx}}{1 - \sigma_{zx}^2} = \frac{(\alpha + \beta\gamma) - (\gamma\delta + \alpha\beta\delta)\beta\delta}{1 - \beta^2\delta^2} \\ &= \frac{\alpha + \beta\delta - \beta\delta^2(\gamma + \alpha\beta)}{1 - \beta^2\delta^2} = \frac{\alpha(1 - \beta^2\delta^2) + \gamma\beta(1 - \delta^2)}{1 - \beta^2\delta^2} \\ &= \alpha + \frac{\gamma\beta(1 - \delta^2)}{1 - \beta^2\delta^2} \end{aligned} \tag{27}$$

As expected, the bias vanishes when δ approaches unity, indicating a faithful proxy. Moreover, if δ can be estimated from an external pilot study, the causal effect α can be identified [See 30, 31] Remarkably, identical behavior emerges in the model of Figure 13(b) in which Z is a driver of U , rather than a proxy.

The same treatment can be applied to errors in measurements of X or of Y and, in each case, the formula of $\sigma_{xy \cdot z}$ reveals what model parameters are the ones affecting the resulting bias.

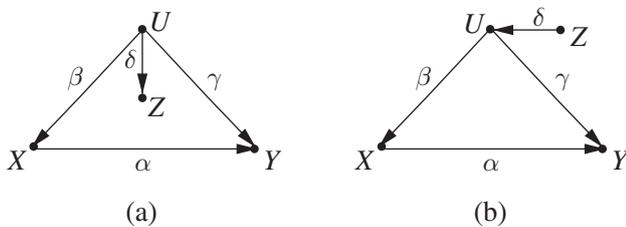


Figure 13 Conditioning on Z , a proxy for the unobserved confounder U , does not remove the bias ($\beta\gamma$).

4 Conclusions

We have demonstrated how path-analytic techniques can illuminate the emergence of several phenomena in causal analysis and how these phenomena depend on the structural features of the model. Although the techniques are limited to linear analysis, hence, restricted to homogeneous populations with no interactions, they can be superior to simulation studies whenever conceptual understanding is of essence, and problem size is manageable.

Acknowledgment: This research was supported in parts by grants from NSF #IIS-1249822 and ONR #N00014-13-1-0153.

Appendix

In linear systems, the explanation for the equality $\sigma_{yx} = \sigma_{yx \cdot w}$ in Figure 3 is simple. Conditioning on W does not physically constrain Z , it merely limits the variance of Z in the subpopulation satisfying $W = w$ which was chosen for observations. Given that effect-homogeneity prevails of linear models, we know that the effect of X on Z remains invariant to the level w chosen for observation and therefore this w -specific effect reflects the effect of X on the entire population. This dictates (in a confounding-free model) $\beta_{xy \cdot w} = \beta_{xy}$.

But how can we explain the persistence of this phenomenon in nonparametric models, where we know (e.g., using *do*-calculus [7]) that adjustment for W does not have any effect on the resulting estimand? In other words, the equality

$$E[Y | X = x] = E_w E[Y | X = x, W = w]$$

will hold in the model of Figure 3 even when the structural equations are nonlinear. Indeed, the independence of W and X , implies

$$\begin{aligned} E[Y | X = x] &= \sum_w E[Y | X = x, W = w] P(W = w | X = x) \\ &= \sum_w E[Y | X = x, W = w] P(W = w) \\ &= E_w E[Y | X = x, W = w] \end{aligned}$$

The answer is that adjustment for W involves averaging over W ; conditioning on W does not. In other words, whereas the effect of X on Z may vary across strata of W , the average of this effect is none other but the effect over the entire population, that is, $E[Y | do(X = x)]$, which equals $E[Y | X = x]$ in the non-confounding case.

Symbolically, we have

$$\begin{aligned} E[Y | do(X = x)] &= \sum_w E[Y | do(X = x), W = w] P[W = w | do(X = x)] \\ &= \sum_w E[Y | do(X = x), W = w] P(W = w) \\ &= \sum_w E[Y | X = x, W = w] P(W = w) \\ &= E(Y | X = x) \end{aligned}$$

The first reduction is licensed by the fact that X has no effect on W and the second by the back-door condition.

References

1. Wright S. Correlation and causation. *J Agric Res* 1921;20:557–85.
2. Duncan O. Introduction to structural equation models. New York: Academic Press, 1975.
3. Kenny D. Correlation and Causality. New York: Wiley, 1979.
4. Heise D. Causal analysis. New York: John Wiley and Sons, 1975.
5. Crámer H. Mathematical methods of statistics. Princeton, NJ: Princeton University Press, 1946.
6. Pearl J. [Causal diagrams for empirical research](#). *Biometrika* 1995;82:669–710.
7. Pearl J. Causality: models, reasoning, and inference, 2nd ed. New York: Cambridge University Press, 2009.

8. Cox D. The planning of experiments. New York: John Wiley and Sons, 1958.
9. Weinberg C. [Toward a clearer definition of confounding](#). Am J Epidemiol 1993;137:1–8.
10. Heckman JJ. [Sample selection bias as a specification error](#). Econometrica 1979;47:153–61.
11. Robins J. Data, design, and background knowledge in etiologic inference. Epidemiology 2001;12:313–20.
12. Bareinboim E, Pearl J. Controlling selection bias in causal inference. In Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS). La Palma, Canary Islands, 2012;100–8.
13. Daniel RM, Kenward MG, Cousens SN, Stavola BLD. [Using causal diagrams to guide analysis in missing data problems](#). Stat Methods Med Res 2011;21:243–56.
14. Geneletti S, Richardson S, Best N. Adjusting for selection bias in retrospective, case-control studies. Biostatistics 2009;10:17–31.
15. Pearl J. A solution to a class of selection-bias problems. Technical Report, R-405, http://ftp.cs.ucla.edu/pub/stat_ser/r405.pdf, Department of Computer Science, University of California, Los Angeles, CA, 2012.
16. Berkson J. [Limitations of the application of fourfold table analysis to hospital data](#). Biometrics Bull 1946;2:47–53.
17. Kim J, Pearl J. A computational model for combined causal and diagnostic reasoning in inference systems. In Proceedings of the Eighth International Joint Conference on Artificial Intelligence (IJCAI-83). Karlsruhe, Germany, 1983.
18. Little RJ, Rubin DB. Statistical analysis with missing data, Vol. 4. New York: Wiley, 1987.
19. Pearl J. Myth, confusion, and science in causal analysis. Technical Report, R-348, University of California, Los Angeles, CA. http://ftp.cs.ucla.edu/pub/stat_ser/r348.pdf, 2009.
20. Rubin D. Author’s reply: should observational studies be designed to allow lack of balance in covariate distributions across treatment group? Stat Med 2009;28:1420–23.
21. Mohan K, Pearl J, Tian J. Missing data as a causal inference problem. Technical Report R-410, http://ftp.cs.ucla.edu/pub/stat_ser/r410.pdf, University of California Los Angeles, Computer Science Department, Los Angeles, CA, 2013.
22. Rosenbaum P. Observational studies, 2nd ed. New York: Springer-Verlag, 2002.
23. Goldberger A. [Reverse regression and salary discrimination](#). J Hum Resour 1984;19:293–318.
24. Hirano K, Imbens G. Estimation of causal effects using propensity score weighting: an application to data on right heart catheterization. Health Serv Outcomes Res Methodol 2001;2:259–78.
25. Bhattacharya J, Vogt W. Do instrumental variables belong in propensity scores? Tech. Rep. NBER Technical Working Paper 343, National Bureau of Economic Research, MA, 2007.
26. Pearl J. On a class of bias-amplifying variables that endanger effect estimates. In Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence. AUAI, Corvallis, OR, 417–24. http://ftp.cs.ucla.edu/pub/stat_ser/r356.pdf, 2010.
27. Wooldridge J. Should instrumental variables be used as matching variables? Technical Report, <https://www.msu.edu/ec/faculty/wooldridge/current%20research/treat1r6.pdf>, Michigan State University, MI, 2009.
28. Myers JA, Rassen JA, Gagne JJ, Huybrechts KF, Schneeweiss S, Rothman KJ, Joffe MM, Glynn RJ. [Effects of adjusting for instrumental variables on bias and precision of effect estimates](#). Am J Epidemiol 2011;174:1213–22.
29. Pearl J. Invited commentary: understanding bias amplification. Am J Epidemiol 2011 [online]. DOI: 10.1093/aje/kwr352.
30. Pearl J. On measurement bias in causal inferences. In Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence. AUAI, Corvallis, OR, 425–32. http://ftp.cs.ucla.edu/pub/stat_ser/r357.pdf, 2010.
31. Kuroki M, Pearl J. Measurement bias and effect restoration in causal inference. Technical Report, R-366, http://ftp.cs.ucla.edu/pub/stat_ser/r366.pdf, University of California Los Angeles, Computer Science Department, Los Angeles, CA, 2013.

