

The Mathematics of Causal Relations

Judea Pearl
Cognitive Systems Laboratory
Computer Science Department
University of California, Los Angeles, CA 90024 USA

March 24, 2008

Abstract

This paper introduces empirical researchers to recent advances in causal inference and stresses the paradigmatic shifts that must be undertaken in moving from traditional statistical analysis to causal analysis of multivariate data. Special emphasis is placed on the assumptions that underly all causal inferences, the languages used in formulating those assumptions, and the conditional nature of causal claims inferred from nonexperimental studies.

In particular, the paper advocates a formalism based on nonparametric structural equations [Pearl, 2000a] which provides both a mathematical foundation for the analysis of counterfactuals and a conceptually transparent language for expressing causal knowledge. This framework gives rise to a friendly calculus of causation that unifies the graphical, potential outcome (Neyman-Rubin) and structural equation approaches and resolves long-standing problems in several of the sciences. These include questions of confounding, causal effect estimation, policy analysis, legal responsibility, direct and indirect effects, instrumental variables, surrogate designs, and the integration of data from experimental and observational studies.

KEY WORDS: Structural equation models, confounding, Rubin causal model, graphical methods, counterfactuals, causal effects.

1 Introduction

Almost two decades have passed since Paul Holland published his highly cited review paper on the Neyman-Rubin (NR) approach to causal inference [Holland, 1986]. Our understanding of causal inference has since increased several folds, due primarily to advances in three areas:

1. Nonparametric structural equations
2. Graphical models
3. Symbiosis between counterfactual and graphical methods.

These advances are central to the empirical sciences because the research questions that motivate most studies in the health, social and behavioral sciences are not statistical but causal in nature. For example, what is the efficacy of a given drug in a given population? Whether data can prove an employer guilty of hiring discrimination? What fraction of past crimes could have been avoided by a given policy? What was the cause of death of a given individual, in a specific incident?

Remarkably, although much of the conceptual framework and algorithmic tools needed for tackling such problems are now well established, they are hardly known to researchers in the field who could put them into practical use. Why?

Solving causal problems mathematically requires certain extensions in the standard mathematical language of statistics, and these extensions are not generally emphasized in the mainstream literature and

education. As a result, large segments of the statistical research community find it hard to appreciate and benefit from the many results that causal analysis has produced in the past two decades.

This paper aims at making these advances more accessible to the general research community by, first, contrasting causal analysis with standard statistical analysis and, second, by comparing and unifying various approaches to causal analysis.

2 From Associational to Causal Analysis: Distinctions and Barriers

2.1 The Basic Distinction: Coping With Change

The aim of standard statistical analysis, typified by regression, estimation, and hypothesis testing techniques, is to assess parameters of a distribution from samples drawn of that distribution. With the help of such parameters, one can infer associations among variables, estimate the likelihood of past and future events, as well as update the likelihood of events in light of new evidence or new measurements. These tasks are managed well by standard statistical analysis so long as experimental conditions remain the same. Causal analysis goes one step further; its aim is to infer not only the likelihood of events under static conditions, but also the dynamics of events under *changing conditions*, for example, changes induced by treatments or external interventions.

This distinction implies that causal and associational concepts do not mix. There is nothing in the joint distribution of symptoms and diseases to tell us that curing the former would or would not cure the latter. More generally, there is nothing in a distribution function to tell us how that distribution would differ if external conditions were to change—say from observational to experimental setup—because the laws of probability theory do not dictate how one property of a distribution ought to change when another property is modified. This information must be provided by causal assumptions which identify relationships that remain invariant when external conditions change.

These considerations imply that the slogan “correlation does not imply causation” can be translated into a useful principle: one cannot substantiate causal claims from associations alone, even at the population level—behind every causal conclusion there must lie some causal assumption that is not testable in observational studies.

2.2 Formulating the Basic Distinction

A useful demarcation line that makes the distinction between associational and causal concepts crisp and easy to apply, can be formulated as follows. An associational concept is any relationship that can be defined in terms of a joint distribution of observed variables, and a causal concept is any relationship that cannot be defined from the distribution alone. Examples of associational concepts are: correlation, regression, dependence, conditional independence, likelihood, collapsibility, risk ratio, odd ratio, marginalization, conditionalization, “controlling for,” and so on. Examples of causal concepts are: randomization, influence, effect, confounding, “holding constant,” disturbance, spurious correlation, instrumental variables, intervention, explanation, attribution, and so on. The former can, while the latter cannot be defined in term of distribution functions.

This demarcation line is extremely useful in causal analysis for it helps investigators to trace the assumptions that are needed for substantiating various types of scientific claims. Every claim invoking causal concepts must rely on some premises that invoke such concepts; it cannot be inferred from, or even defined in terms statistical associations alone.

2.3 Ramifications of the Basic Distinction

This principle has far reaching consequences that are not generally recognized in the standard statistical literature. Many researchers, for example, are still convinced that confounding is solidly founded in standard, frequentist statistics, and that it can be given an associational definition saying (roughly): “ U is a potential confounder for examining the effect of treatment X on outcome Y when both U and X and

U and Y are not independent.” That this definition and all its many variants must fail, is obvious from the demarcation line above; “independence” is an associational concept while confounding is needed for establishing causal relations. The two do not mix, hence, the definition must be false. Therefore, to the bitter disappointment of generations of epidemiology researchers, confounding bias cannot be detected or corrected by statistical methods alone; one must make some judgmental assumptions regarding causal relationships in the problem before an adjustment (e.g., by stratification) can safely correct for confounding bias.

Another ramification of the sharp distinction between associational and causal concepts is that any mathematical approach to causal analysis must acquire new notation for expressing causal relations – probability calculus is insufficient. To illustrate, the syntax of probability calculus does not permit us to express the simple fact that “symptoms do not cause diseases”, let alone draw mathematical conclusions from such facts. All we can say is that two events are dependent—meaning that if we find one, we can expect to encounter the other, but we cannot distinguish statistical dependence, quantified by the conditional probability $P(\text{disease}|\text{symptom})$ from causal dependence, for which we have no expression in standard probability calculus. Scientists seeking to express causal relationships must therefore supplement the language of probability with a vocabulary for causality, one in which the symbolic representation for the relation “symptoms cause disease” is distinct from the symbolic representation of “symptoms are associated with disease.”

2.4 Two Mental Barriers: Untested Assumptions and New Notation

The preceding two requirements: (1) to commence causal analysis with untested,¹ theoretically or judgmentally based assumptions, and (2) to extend the syntax of probability calculus, constitute the two main obstacles to the acceptance of causal analysis among statisticians and among professionals with traditional training in statistics.

Associational assumptions, even untested, are testable in principle, given sufficiently large sample and sufficiently fine measurements. Causal assumptions, in contrast, cannot be verified even in principle, unless one resorts to experimental control. This difference stands out in Bayesian analysis. Though the priors that Bayesians commonly assign to statistical parameters are untested quantities, the sensitivity to these priors tends to diminish with increasing sample size. In contrast, sensitivity to prior causal assumptions, say that treatment does not change gender, remains substantial regardless of sample size.

This makes it doubly important that the notation we use for expressing causal assumptions be meaningful and unambiguous so that one can clearly judge the plausibility or inevitability of the assumptions articulated. Statisticians can no longer ignore the mental representation in which scientists store experiential knowledge, since it is this representation, and the language used to access that representation that determine the reliability of the judgments upon which the analysis so crucially depends.

How does one recognize causal expressions in the statistical literature? Those versed in the potential-outcome notation [Neyman, 1923; Rubin, 1974; Holland, 1988], can recognize such expressions through the subscripts that are attached to counterfactual events and variables, e.g. $Y_x(u)$ or Z_{xy} . (Some authors use parenthetical expressions, e.g. $Y(x, u)$ or $Z(x, y)$.) The expression $Y_x(u)$, for example, stands for the value that outcome Y would take in individual u , had treatment X been at level x . If u is chosen at random, Y_x is a random variable, and one can talk about the probability that Y_x would attain a value y in the population, written $P(Y_x = y)$. Alternatively, Pearl [1995] used expressions of the form $P(Y = y|set(X = x))$ or $P(Y = y|do(X = x))$ to denote the probability (or frequency) that event ($Y = y$) would occur if treatment condition $X = x$ were enforced uniformly over the population.² Still a third notation that distinguishes causal expressions is provided by graphical models, where the arrows convey causal directionality.³

However, few have taken seriously the textbook requirement that any introduction of new notation must entail a systematic definition of the syntax and semantics that governs the notation. Moreover, in

¹By “untested” I mean untested using frequency data in nonexperimental studies.

²Clearly, $P(Y = y|do(X = x))$ is equivalent to $P(Y_x = y)$, This is what we normally assess in a controlled experiment, with X randomized, in which the distribution of Y is estimated for each level x of X .

³These notational clues should be useful for detecting inadequate definitions of causal concepts; any definition of confounding, randomization or instrumental variables that is cast in standard probability expressions, void of graphs, counterfactual subscripts or $do(*)$ operators, can safely be discarded as inadequate.

the bulk of the statistical literature before 2000, causal claims rarely appear in the mathematics. They surface only in the verbal interpretation that investigators occasionally attach to certain associations, and in the verbal description with which investigators justify assumptions. For example, the assumption that a covariate is not affected by a treatment, a necessary assumption for the control of confounding [Cox, 1958], is expressed in plain English, not in a mathematical expression.

Remarkably, though the necessity of explicit causal notation is now recognized by most leaders in the field, the use of such notation has remained enigmatic to most rank and file researchers, and its potentials still lay grossly underutilized in the statistics based sciences. The reason for this, I am firmly convinced, can be traced to the unfriendly and ad-hoc way in which causal analysis, has been presented to the research community, relying primarily on the NR and structural equation models.

The next section provides a conceptualization that overcomes these mental barriers; it offers both a friendly mathematical machinery for cause-effect analysis and a formal foundation for counterfactual analysis.

3 The Language of Diagrams and Structural Equations

3.1 Semantics: Causal Effects and Counterfactuals

How can one express mathematically the common understanding that symptoms do not cause diseases? The earliest attempt to formulate such relationship mathematically was made in the 1920's by the geneticist Sewall Wright [1921], who used a combination of equations and graphs. For example, if X stands for a disease variable and Y stands for a certain symptom of the disease, Wright would write a linear equation:

$$y = \beta x + u \tag{1}$$

where x stands for the level (or severity) of the disease, y stands for the level (or severity) of the symptom, and u stands for all factors, other than the disease in question, that could possibly affect Y . In interpreting this equation one should think of a physical process whereby Nature *examines* the values of x and u and, accordingly, *assigns* variable Y the value $y = \beta x + u$.

To express the directionality inherent in this process, Wright augmented the equation with a diagram, later called “path diagram,” in which arrows are drawn from perceived causes to their (perceived) effects and, more importantly, the absence of an arrow makes the empirical claim that the value Nature assigns to one variable is not determined by the value taken by another.

The variables V and U are called “exogenous” ; they represent observed or unobserved background factors that the modeler decides to keep unexplained, that is, factors that influence but are not influenced by the other variables (called “endogenous”) in the model.

If correlation is judged possible between two exogenous variables, U and V , it is customary to connect them by a dashed double arrow, as shown in Fig. 1(b).

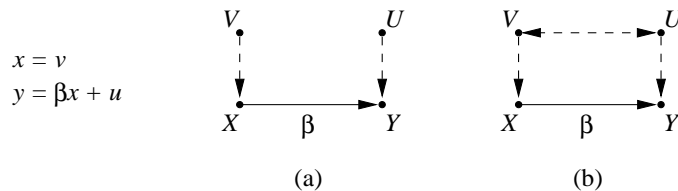


Figure 1: A simple structural equation model, and its associated diagrams. Unobserved exogenous variables are connected by dashed arrows.

To summarize, path diagrams encode causal assumptions via missing arrows, representing claims of zero influence, and missing double arrows (e.g., between V and U), representing the (causal) assumption $Cov(U, V)=0$.

The generalization to nonlinear system of equations is straightforward. For example, the non-parametric interpretation of the diagram of Fig. 2(a) corresponds to a set of three functions, each

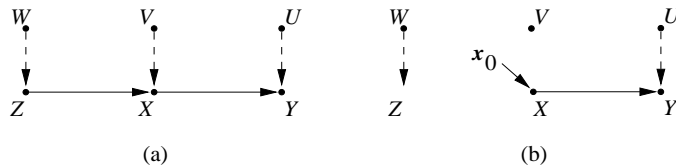


Figure 2: (a) The diagram associated with the structural model of Eq. (2). (b) The diagram associated with the modified model of Eq. (3), representing the intervention $do(X = x_0)$.

corresponding to one of the observed variables:

$$\begin{aligned}
 z &= f_Z(w) \\
 x &= f_X(z, v) \\
 y &= f_Y(x, u)
 \end{aligned} \tag{2}$$

where W, V and U are assumed to be jointly independent but, otherwise, arbitrarily distributed.

Remarkably, unknown to most economists and philosophers, structural equation models provide a formal interpretation and symbolic machinery for analyzing counterfactual relationships of the type: “ Y would be y had X been x in situation $U = u$,” denoted $Y_x(u) = y$. Here U represents the vector of all exogenous variables.

The key idea is to interpret the phrase “had X been x_0 ” as an instruction to modify the original model and replace the equation for X by a constant x_0 , yielding

$$\begin{aligned}
 z &= f_Z(w) \\
 x &= x_0 \\
 y &= f_Y(x, u)
 \end{aligned} \tag{3}$$

the graphical description of which is shown in Fig. 2(b).

This replacement permits the constant x_0 to differ from the actual value of X (namely $f_X(z, v)$) without rendering the system of equations inconsistent, thus yielding a formal interpretation of counterfactuals in multi-stage models, where the dependent variable in one equation may be an independent variable in another [Balke and Pearl, 1994ab; Pearl, 2000b]. For example, to compute the average causal effect of X on Y , i.e., $E(Y_{x_0})$ we solve Eq. (3) for Y in terms of the exogenous variables, yielding $Y_{x_0} = f_Y(x_0, u)$, and average over U and V . To answer more sophisticated questions such as whether Y would be y_1 if X were x_1 , given that in fact Y is y_0 and X is x_0 , we need to compute the conditional probability $P(Y_{x_1} = y_1 | Y = y_0, X = x_0)$ which is well defined once we know the forms of the structural equations and the distribution of the exogenous variables in the model.

This interpretation of counterfactuals, cast as solutions to modified systems of equations, provides the conceptual and formal link between structural equation models, used in economics and social science and the Neyman-Rubin potential-outcome framework to be discussed in Section 4. But first we discuss two long-standing problems that have been completely resolved in purely graphical terms, without delving into algebraic techniques.

3.2 Confounding and Causal Effect Estimation

The central target of most studies in the social and health sciences is the elucidation of cause-effect relationships among variables of interests, for example, treatments, policies, preconditions and outcomes. While good statisticians have always known that the elucidation of causal relationships from observational studies must be shaped by assumptions about how the data were generated, the relative roles of assumptions and data, and ways of using those assumptions to eliminate confounding bias have been a subject of much controversy. The structural framework of Section 3.1 puts these controversies to rest.

Covariate Selection: The back-door criterion

Consider an observational study where we wish to find the effect of X on Y , for example, treatment on response, and assume that the factors deemed relevant to the problem are structured as in Fig. 3; some are affecting the response, some are affecting the treatment and some are affecting both treatment

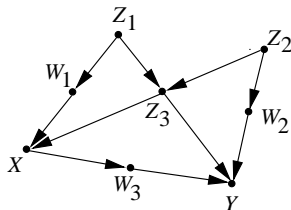


Figure 3: Graphical model illustrating the back-door criterion. Error terms are not shown explicitly.

and response. Some of these factors may be unmeasurable, such as genetic trait or life style, others are measurable, such as gender, age, and salary level. Our problem is to select a subset of these factors for measurement and adjustment, namely, that if we compare treated vs. untreated subjects having the same values of the selected factors, we get the correct treatment effect in that subpopulation of subjects. Such a set of factors is called a “sufficient set” or a set “appropriate for adjustment”. The problem of defining a sufficient set, let alone finding one, has baffled epidemiologists and social science for decades (see Greenland *et al.*, [1999], Pearl [2000a] and [2003] for review).

The following criterion, named “back-door” in [Pearl, 1993a], provides a graphical method of selecting such a set of factors for adjustment. It states that a set S is appropriate for adjustment if two conditions hold:

1. No element of S is a descendant of X
2. The elements of S “block” all “back-door” paths from X to Y , namely all paths that end with an arrow pointing to X .⁴

Based on this criterion we see, for example, that the sets $\{Z_1, Z_2, Z_3\}$, $\{Z_1, Z_3\}$, and $\{W_2, Z_3\}$, each is sufficient for adjustment, because each blocks all back-door paths between X and Y . The set $\{Z_3\}$, however, is not sufficient for adjustment because, as explained above, it does not block the path $X \leftarrow W_1 \leftarrow Z_1 \rightarrow Z_3 \leftarrow Z_2 \rightarrow W_2 \rightarrow Y$.

The implication of finding a sufficient set S is that, stratifying on S is guaranteed to remove all confounding bias relative the causal effect of X on Y . In other words, it renders the causal effect of X on Y identifiable, via

$$\begin{aligned} P(Y = y|do(X = x)) \\ = \sum_s P(Y = y|X = x, S = s)P(S = s) \end{aligned} \tag{4}$$

Since all factors on the right hand side of the equation are estimable (e.g., by regression) from the pre-interventional data, the causal effect can likewise be estimated from such data without bias.

The back-door criterion allows us to write Eq. (4) directly, after selecting a sufficient set S from the diagram, without resorting to any algebraic manipulation. The selection criterion can be applied systematically to diagrams of any size and shape, thus freeing analysts from judging whether “ X is conditionally ignorable given S ,” a formidable mental task required in the potential-response framework [Rosenbaum and Rubin, 1983]. The criterion also enables the analyst to search for an optimal set of covariate—namely, a set S that minimizes measurement cost or sampling variability [Tian et al., 1998].

⁴In this criterion, a set S of nodes is said to block a path p if either (i) p contains at least one arrow-emitting node that is in S , or (ii) p contains at least one collision node that is outside S and has no descendant in S . See Pearl [2000a, pp. 16-17]

General control of confounding

Adjusting for covariates is only one of many methods that permits us to estimate causal effects in nonexperimental studies. A much more general identification criterion is provided by the following theorem:

Theorem 1 [*Tian and Pearl, 2002*]

*A sufficient condition for identifying the causal effect $P(y|do(x))$ is that every path between X and any of its children traces at least one arrow emanating from a measured variable.*⁵

For example, if W_3 is the only observed covariate in the model of Fig. 3, then there exists no sufficient set for adjustment (because no set of observed covariates can block the paths from X to Y through Z_3), yet $P(y|do(x))$ can nevertheless be estimated since the one path from X to W_3 (the only child of X) traces the arrow $X \rightarrow W_3$, which emanates from a measured variable, X . In this example, the variable W_3 acts as a “mediating instrumental variable” [Pearl 1993b; Chalak and White, 2006] and yields the estimand:

$$\begin{aligned} P(Y = y|do(X = x)) &= \sum_{w_3} P(W_3 = w|do(X = x))P(Y = y|do(W_3 = w)) \\ &= \sum_w P(w|x) \sum_{x'} P(y|w, x')P(x') \end{aligned} \tag{5}$$

More recent results extend this theorem by (1) presenting a necessary and sufficient condition for identification [Shpitser and Pearl, 2006], and (2) extending the condition from causal effects to any counterfactual expression [Shpitser and Pearl, 2007]. The corresponding unbiased estimands for these causal quantities, are readable directly from the diagram.

4 The Language of Potential Outcomes

The elementary object of analysis in the potential-outcome framework is the unit-based response variable, denoted $Y_x(u)$, read: “the value that Y would obtain in unit u , had treatment X been x ” [Neyman, 1923; Rubin, 1974]. These subscripted variables are treated as undefined quantities, useful for expressing the causal quantities we seek, but are not derived from other quantities in the model. In contrast, in the previous section counterfactual entities were derived from a set of meaningful physical processes, each represented by an equation, and *unit* was interpreted a vector u of background factors that characterize an experimental unit. Each structural equation model thus provides a compact representation for a huge number of counterfactual claims, guaranteed to be consistent. The potential outcome framework lacks such compactness, nor does it provide guarantees that any given set of claims is consistent.

In view of these features, the structural definition of $Y_x(u)$ can be regarded as the formal basis for the potential outcome approach. It interprets the opaque English phrase “the value that Y would obtain in unit u , had X been x ” in terms of a meaningful mathematical model that allows such values to be computed unambiguously. Consequently, important concepts in potential response analysis that researchers find ill-defined or overly esoteric often obtain meaningful and natural interpretation in the structural semantics. Examples are: “unit” (“exogenous variables” in structural semantics), “principal stratification” (“equivalence classes” in structural semantics [Balke and Pearl, 1994a] and [Pearl 2000b] “conditional ignorability” (“back-door condition” in [Pearl 1993a] “assignment mechanism” ($P(x|direct-causes\ of\ X)$ in structural semantics) and so on. The next two subsections examine how assumptions and inferences are handled in the potential outcome approach vis a vis the graphical-structural approach.

4.1 Formulating Assumptions

The distinct characteristic of the potential outcome approach is that, although its primitive objects are undefined, hypothetical quantities, the analysis itself is conducted almost entirely within the axiomatic framework of probability theory. This is accomplished, by postulating a “super” probability function

⁵Before applying this criterion, one may delete from the causal graph all nodes that are not ancestors of Y .

on both hypothetical and real events, treating the former as "missing data". In other words, if U is treated as a random variable then the value of the counterfactual $Y_x(u)$ becomes a random variable as well, denoted as Y_x . The potential-outcome analysis proceeds by treating the observed distribution $P(x_1, \dots, x_n)$ as the marginal distribution of an augmented probability function P^* defined over both observed and counterfactual variables. Queries about causal effects are phrased as queries about the probability distribution of the counterfactual variable of interest, written $P^*(Y_x = y)$. The new hypothetical entities Y_x are treated as ordinary random variables; for example, they are assumed to obey the axioms of probability calculus, the laws of conditioning, and the axioms of conditional independence. Moreover, these hypothetical entities are not entirely whimsy, but are assumed to be connected to observed variables via consistency constraints [Robins, 1986] such as

$$X = x \implies Y_x = Y, \tag{6}$$

which states that, for every u , if the actual value of X turns out to be x , then the value that Y would take on if X were x is equal to the actual value of Y . For example, a person who chose treatment x and recovered, would also have recovered if given treatment x by design.

The main conceptual difference between the two approaches is that, whereas the structural approach views the subscript x as an operation that changes the distribution but keeps the variables the same, the potential-outcome approach views the variable Y_x , to be a different variable, loosely connected to Y through relations such as (6).

Pearl [2000a, Chapter 7] shows, using the structural interpretation of $Y_x(u)$, that it is indeed legitimate to treat counterfactuals as jointly distributed random variables in all respects, that consistency constraints like (6) are automatically satisfied in the structural interpretation and, moreover, that investigators need not be concerned about any additional constraints except the following two⁶:

$$Y_{yz} = y \quad \text{for all } y \text{ and } z \tag{7}$$

$$X_z = x \implies Y_{xz} = Y_z \quad \text{for all } x \text{ and } z \tag{8}$$

Eq. (7) ensures that the interventions $do(Y = y)$ results in the condition $Y = y$, regardless of concurrent interventions, say $do(Z = z)$, that are applied to variables other than Y . Equation (8) generalizes (6) to cases where Z is held fixed, at z .

To communicate substantive causal knowledge, the potential-outcome analyst must express causal assumptions as constraints on P^* , usually in the form of conditional independence assertions involving counterfactual variables. In Fig. 2(a) for instance, to communicate the understanding that a treatment assignment (Z) is randomized (hence independent of both U and V), the potential-outcome analyst needs to use the independence constraint $Z \perp\!\!\!\perp \{X_z, Y_x\}$. To further formulate the understanding that Z does not affect Y directly, except through X , the analyst would write a, so called, "exclusion restriction": $Y_{xz} = Y_x$. Clearly, no mortal can judge the validity of such assumptions in any real life problem without resorting to graphs.⁷

4.2 Performing Inferences

A collection of assumptions of this type might sometimes be sufficient to permit a unique solution to the query of interest; in other cases, only bounds on the solution can be obtained. For example, if one can plausibly assume that a set Z of covariates satisfies the conditional independence

$$Y_x \perp\!\!\!\perp X | Z \tag{9}$$

⁶This *completeness* result is due to Halpern [1998], who noted that an additional axiom

$$\{Y_{xz} = y\} \& \{Z_{xy} = z\} \implies Y_x = y$$

must hold in non-recursive models. This fundamental axiom may come to haunt economists and social scientists who blindly apply NR analysis in their fields.

⁷Even with the use of graphs the task is not easy, for example, the reader should try to verify whether $\{Z \perp\!\!\!\perp X_z | Y\}$ holds in the simple model of Fig. 2(a). The answer is given in Pearl [2000a, page 214]

(an assumption that was termed “conditional ignorability” by [Rosenbaum and Rubin, 1983] then the causal effect $P^*(Y_x = y)$ can readily be evaluated to yield

$$\begin{aligned}
 P^*(Y_x = y) &= \sum_z P^*(Y_x = y|z)P(z) \\
 &= \sum_z P^*(Y_x = y|x, z)P(z) \quad (\text{using (9)}) \\
 &= \sum_z P^*(Y = y|x, z)P(z) \quad (\text{using (6)}) \\
 &= \sum_z P(y|x, z)P(z). \tag{10}
 \end{aligned}$$

which is the usual covariate-adjustment formula, as in Eq. (4).

Note that almost all mathematical operations in this derivation are conducted within the safe confines of probability calculus. Save for an occasional application of rule (8) or (6)), the analyst may forget that Y_x stands for a counterfactual quantity—it is treated as any other random variable, and the entire derivation follows the course of routine probability exercises.

However, this mathematical illusion comes at the expense of conceptual clarity, especially at a stage where causal assumptions need be formulated. The reader may appreciate this aspect by attempting to judge whether the assumption of conditional ignorability Eq. (9), the key to the derivation of Eq. (10), holds in any familiar situation, say in the experimental setup of Fig. 2(a). This assumption reads: “the value that Y would obtain had X been x , is independent of X , given Z ”. Such assumptions of conditional independence among counterfactual variables are not straightforward to comprehend or ascertain, for they are cast in a language far removed from ordinary understanding of cause and effect. When counterfactual variables are not viewed as byproducts of a deeper, process-based model, it is also hard to ascertain whether *all* relevant counterfactual independence judgments have been articulated, whether the judgments articulated are redundant, or whether those judgments are self-consistent.

The need to express, defend, and manage formidable counterfactual relationships of this type explains the slow acceptance of causal analysis among epidemiologists and statisticians, and why economists and social scientists continue to use structural equation models instead of the potential-outcome alternatives advocated in Holland [1988], Angrist *et al.* [1996], and Sobel [1998].

On the other hand, the algebraic machinery offered by the potential-outcome notation, once a problem is properly formalized, can be powerful in refining assumptions [Angrist et al., 1996], deriving consistent estimands [Robins, 1986], bounding probabilities of causation [Tian and Pearl, 2000], and combining data from experimental and nonexperimental studies [Pearl, 2000a, pages 302-303]

4.3 Combining Graphs and Algebra – Methods and Accomplishments.

Pearl [2000a, page 232] presents a way of combining the best features of the two approaches. It is based on encoding causal assumptions in the language of diagrams, translating these assumptions into potential outcome notation, performing the mathematics in the algebraic language of counterfactuals and, finally, interpreting the result in plain causal language. Often, the answer desired can be obtained directly from the diagram, and no translation is necessary (as demonstrated in Section 3.2).

This method has scored an impressive list of accomplishments, including solutions to the long-standing problems of legal responsibility [Tian and Pearl, 2000; Pearl, 2001], non-compliance [Balke and Pearl, 1997; Chickering and Pearl, 1997], direct and indirect effects [Pearl, 2001], mediating instrumental variables [Pearl, 1993b; Brito and Pearl, 2006], robustness analysis [Pearl, 2004], and the integration of data from experimental and observational studies [Tian and Pearl, 2000; Pearl, 2000a]. Detailed descriptions of these results are given in the corresponding articles which are available on bayes.cs.ucla.edu/jp.home.html.

5 Conclusions

Statistics is strong in devising ways of describing data and inferring distributional parameters from sample. Causal inference require two addition ingredients: a science-friendly language for articulating

causal knowledge, and a mathematical machinery for processing that knowledge, combining it with data and drawing new causal conclusions about a phenomena. This paper introduces nonparametric structural equations models as a formal and meaningful language for formulating causal knowledge and for explicating causal concepts used in scientific discourse. These include: randomization, intervention, direct and indirect effects, confounding, counterfactuals, and attribution. The algebraic component of the structural language coincides with the potential-outcome framework, and its graphical component embraces Wright's method of path diagrams (in its nonparametric version.) When unified and synthesized, the two components offer investigators a powerful methodology for empirical research. The merits of this methodology have quickly been recognized by several research communities [e.g., Morgan and Winship, 2007; Greenland *et al.*, 1999; Petersen *et al.*, 2006; Chalak and White, 2006] and are making their way, past obvious pockets of resistance, to statistical education as well.

Perhaps the most important message of the discussion and methods presented in this paper would be a widespread awareness that (1) all studies concerning causal relations must begin with causal assumptions of some sort and (2) that a friendly and formal language is currently available for articulating such assumptions. This means that scientific articles concerning questions of causation must contain a section in which causal assumptions are articulated using either graphs or subscripted formulas. Authors who wish their assumptions to be understood, scrutinized and discussed by readers and colleagues would do well to use graphs. Authors who refrain from using graphs would be risking a suspicion of attempting to avoid transparency of their working assumptions.

Another important implication of this paper is that every causal inquiry can be mathematized. In other words, mechanical procedures can now be invoked to determine what assumptions investigators must be willing to make in order for desired quantities to be estimable consistently from the data. This is not to say that the needed assumptions would be reasonable, or that the resulting estimation method would be easy. It means that the needed causal assumptions can be made transparent, brought up for discussion and refinement and, once consistency is assured, causal quantities can be estimated from data through ordinary statistical methods, free of the mystical aura that has shrouded causal analysis in the past.

Acknowledgments

This research was supported in parts by grants from NSF IIS-0535223 and NLM T15 LM07356.

References

- [Angrist et al., 1996] J.D. Angrist, G.W. Imbens, and Rubin D.B. Identification of causal effects using instrumental variables (with comments). *Journal of the American Statistical Association*, 91(434):444–472, June 1996.
- [Balke and Pearl, 1994a] A. Balke and J. Pearl. Counterfactual probabilities: Computational methods, bounds, and applications. In R. Lopez de Mantaras and D. Poole, editors, *Uncertainty in Artificial Intelligence 10*, pages 46–54. Morgan Kaufmann, San Mateo, CA, 1994.
- [Balke and Pearl, 1994b] A. Balke and J. Pearl. Probabilistic evaluation of counterfactual queries. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, volume I, pages 230–237. MIT Press, Menlo Park, CA, 1994.
- [Balke and Pearl, 1997] A. Balke and J. Pearl. Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association*, 92(439):1172–1176, September 1997.
- [Brito and Pearl, 2006] C. Brito and J Pearl. Graphical condition for identification in recursive SEM. In *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence*, pages 47–54. AUAI Press, Corvallis, OR, 2006.
- [Chalak and White, 2006] K. Chalak and H. White. An extended class of instrumental variables for the estimation of causal effects. Technical Report Discussion Paper, UCSD, Department of Economics, July 2006.

- [Chickering and Pearl, 1997] D.M. Chickering and J. Pearl. A clinician’s tool for analyzing non-compliance. *Computing Science and Statistics*, 29(2):424–431, 1997.
- [Cox, 1958] D.R. Cox. *The Planning of Experiments*. John Wiley and Sons, NY, 1958.
- [Greenland et al., 1999] S. Greenland, J. Pearl, and J.M Robins. Causal diagrams for epidemiologic research. *Epidemiology*, 10(1):37–48, 1999.
- [Halpern, 1998] J.Y. Halpern. Axiomatizing causal reasoning. In G.F. Cooper and S. Moral, editors, *Uncertainty in Artificial Intelligence*, pages 202–210. Morgan Kaufmann, San Francisco, CA, 1998.
- [Holland, 1986] P.W. Holland. Statistics and causal inference. *Journal of the American Statistical Association*, 81(396):945–960, December 1986.
- [Holland, 1988] P.W. Holland. Causal inference, path analysis, and recursive structural equations models. In C. Clogg, editor, *Sociological Methodology*, pages 449–484. American Sociological Association, Washington, D.C., 1988.
- [Morgan and Winship, 2007] S.L. Morgan and C. Winship. *Counterfactuals and Causal Inference: Methods and Principles for Social Research (Analytical Methods for Social Research)*. Cambridge University Press, New York, NY, 2007.
- [Neyman, 1923] J. Neyman. On the application of probability theory to agricultural experiments. Essay on principles. Section 9. *Statistical Science*, 5(4):465–480, 1923.
- [Pearl, 1993a] J. Pearl. Comment: Graphical models, causality, and intervention. *Statistical Science*, 8:266–269, 1993.
- [Pearl, 1993b] J. Pearl. Mediating instrumental variables. Technical Report Technical Report R-210, Computer Science Department, UCLA, 1993.
- [Pearl, 1995] J. Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–710, December 1995.
- [Pearl, 2000a] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000.
- [Pearl, 2000b] J. Pearl. Comment on A.P. Dawid’s, causal inference without counterfactuals. *Journal of the American Statistical Association*, 95(450):428–431, June 2000.
- [Pearl, 2001] J. Pearl. Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, pages 411–420. Morgan Kaufmann, San Francisco, CA, 2001.
- [Pearl, 2003] J. Pearl. Statistics and causal inference: A review. *Test Journal*, 12(2):281–345, December 2003.
- [Pearl, 2004] J. Pearl. Robustness of causal claims. *Proceedings of the Twentieth Conference Uncertainty in Artificial Intelligence*, pages 446–453, 2004.
- [Petersen et al., 2006] M.L. Petersen, S.E. Sinisi, and M.J. van der Laan. Estimation of direct causal effects. *Epidemiology*, 17(3):276–284, 2006.
- [Robins, 1986] J.M. Robins. A new approach to causal inference in mortality studies with a sustained exposure period – applications to control of the healthy workers survivor effect. *Mathematical Modeling*, 7:1393–1512, 1986.
- [Rosenbaum and Rubin, 1983] P. Rosenbaum and D. Rubin. The central role of propensity score in observational studies for causal effects. *Biometrika*, 70:41–55, 1983.
- [Rubin, 1974] D.B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66:688–701, 1974.

- [Shpitser and Pearl, 2006] I. Shpitser and J Pearl. Identification of conditional interventional distributions. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, pages 437–444. AUAI Press, Corvallis, OR, 2006.
- [Shpitser and Pearl, 2007] I. Shpitser and J Pearl. What counterfactuals can be tested. In *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence*, pages 352–359. AUAI Press, Vancouver, BC Canada, 2007.
- [Sobel, 1998] M.E. Sobel. Causal inference in statistical models of the process of socioeconomic achievement. *Sociological Methods & Research*, 27(2):318–348, November 1998. cited in jp-98 book.
- [Tian and Pearl, 2000] J. Tian and J. Pearl. Probabilities of causation: Bounds and identification. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, pages 589–598. Morgan Kaufmann, San Francisco, CA, 2000.
- [Tian and Pearl, 2002] J. Tian and J. Pearl. A general identification condition for causal effects. In *Proceedings of the Eighteenth National Conference on Artificial Intelligence*, pages 567–573. AAAI Press/The MIT Press, Menlo Park, CA, 2002.
- [Tian et al., 1998] J. Tian, A. Paz, and J. Pearl. Finding minimal separating sets. Technical Report R-254, University of California, Los Angeles, CA, 1998.
- [Wright, 1921] S. Wright. Correlation and causation. *Journal of Agricultural Research*, 20:557–585, 1921.