

Causal diagrams for empirical research

BY JUDEA PEARL

*Cognitive Systems Laboratory, Computer Science Department, University of California,
Los Angeles, California 90024, U.S.A.*

SUMMARY

The primary aim of this paper is to show how graphical models can be used as a mathematical language for integrating statistical and subject-matter information. In particular, the paper develops a principled, nonparametric framework for causal inference, in which diagrams are queried to determine if the assumptions available are sufficient for identifying causal effects from nonexperimental data. If so the diagrams can be queried to produce mathematical expressions for causal effects in terms of observed distributions; otherwise, the diagrams can be queried to suggest additional observations or auxiliary experiments from which the desired inferences can be obtained.

Some key words: Causal inference; Graph model; Structural equations; Treatment effect.

1. INTRODUCTION

The tools introduced in this paper are aimed at helping researchers communicate qualitative assumptions about cause-effect relationships, elucidate the ramifications of such assumptions, and derive causal inferences from a combination of assumptions, experiments, and data.

The basic philosophy of the proposed method can best be illustrated through the classical example due to Cochran (Wainer, 1989). Consider an experiment in which soil fumigants, X , are used to increase oat crop yields, Y , by controlling the eelworm population, Z , but may also have direct effects, both beneficial and adverse, on yields beside the control of eelworms. We wish to assess the total effect of the fumigants on yields when this study is complicated by several factors. First, controlled randomised experiments are infeasible: farmers insist on deciding for themselves which plots are to be fumigated. Secondly, farmers' choice of treatment depends on last year's eelworm population, Z_0 , an unknown quantity strongly correlated with this year's population. Thus we have a classical case of confounding bias, which interferes with the assessment of treatment effects, regardless of sample size. Fortunately, through laboratory analysis of soil samples, we can determine the eelworm populations before and after the treatment and, furthermore, because the fumigants are known to be active for a short period only, we can safely assume that they do not affect the growth of eelworms surviving the treatment. Instead, eelworm growth depends on the population of birds and other predators, which is correlated, in turn, with last year's eelworm population and hence with the treatment itself.

The method proposed in this paper permits the investigator to translate complex considerations of this sort into a formal language, thus facilitating the following tasks.

- (i) Explicate the assumptions underlying the model.

- (ii) Decide whether the assumptions are sufficient for obtaining consistent estimates of the target quantity: the total effect of the fumigants on yields.
- (iii) If the answer to (ii) is affirmative, provide a closed-form expression for the target quantity, in terms of distributions of observed quantities.
- (iv) If the answer to (ii) is negative, suggest a set of observations and experiments which, if performed, would render a consistent estimate feasible.

The first step in this analysis is to construct a causal diagram such as the one given in Fig. 1, which represents the investigator's understanding of the major causal influences among measurable quantities in the domain. The quantities Z_1 , Z_2 and Z_3 denote, respectively, the eelworm population, both size and type, before treatment, after treatment, and at the end of the season. Quantity Z_0 represents last year's eelworm population; because it is an unknown quantity, it is represented by a hollow circle, as is B , the population of birds and other predators. Links in the diagram are of two kinds: those that connect unmeasured quantities are designated by dashed arrows, those connecting measured quantities by solid arrows. The substantive assumptions embodied in the diagram are negative causal assertions, which are conveyed through the links missing from the diagram. For example, the missing arrow between Z_1 and Y signifies the investigator's understanding that pre-treatment eelworms cannot affect oat plants directly; their entire influence on oat yields is mediated by post-treatment conditions, namely Z_2 and Z_3 . The purpose of the paper is not to validate or repudiate such domain-specific assumptions but, rather, to test whether a given set of assumptions is sufficient for quantifying causal effects from non-experimental data, for example, estimating the total effect of fumigants on yields.

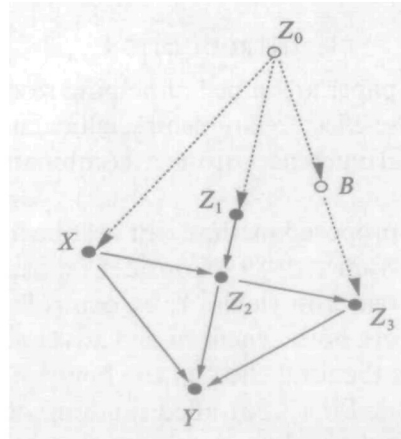


Fig. 1. A causal diagram representing the effect of fumigants, X , on yields, Y .

The proposed method allows an investigator to inspect the diagram of Fig. 1 and conclude immediately the following.

- (a) The total effect of X on Y can be estimated consistently from the observed distribution of X , Z_1 , Z_2 , Z_3 and Y .
- (b) The total effect of X on Y , assuming discrete variables throughout, is given by the formula

$$\text{pr}(y|\check{x}) = \sum_{z_1} \sum_{z_2} \sum_{z_3} \text{pr}(y|z_2, z_3, x) \text{pr}(z_2|z_1, x) \sum_{x'} \text{pr}(z_3|z_1, z_2, x') \text{pr}(z_1, x'), \quad (1)$$

where the symbol \check{x} , read 'x check', denotes that the treatment is set to level $X = x$ by external intervention.

- (c) Consistent estimation of the total effect of X on Y would not be feasible if Y were confounded with Z_3 ; however, confounding Z_2 and Y will not invalidate the formula for $\text{pr}(y|\bar{x})$.

These conclusions can be obtained either by analysing the graphical properties of the diagram, or by performing a sequence of symbolic derivations, governed by the diagram, which gives rise to causal effect formulae such as (1).

The formal semantics of the causal diagrams used in this paper will be defined in § 2, following a review of directed acyclic graphs as a language for communicating conditional independence assumptions. Section 2.2 introduces a causal interpretation of directed graphs based on nonparametric structural equations and demonstrates their use in predicting the effect of interventions. Section 3 demonstrates the use of causal diagrams to control confounding bias in observational studies. We establish two graphical conditions ensuring that causal effects can be estimated consistently from nonexperimental data. The first condition, named the back-door criterion, is equivalent to the ignorability condition of Rosenbaum & Rubin (1983). The second condition, named the front-door criterion, involves covariates that are affected by the treatment, and thus introduces new opportunities for causal inference. In § 4, we introduce a symbolic calculus that permits the stepwise derivation of causal effect formulae of the type shown in (1). Using this calculus, § 5 characterises the class of graphs that permit the quantification of causal effects from nonexperimental data, or from surrogate experimental designs.

2. GRAPHICAL MODELS AND THE MANIPULATIVE ACCOUNT OF CAUSATION

2.1. Graphs and conditional independence

The usefulness of directed acyclic graphs as economical schemes for representing conditional independence assumptions is well evidenced in the literature (Pearl, 1988; Whittaker, 1990). It stems from the existence of graphical methods for identifying the conditional independence relationships implied by recursive product decompositions

$$\text{pr}(x_1, \dots, x_n) = \prod_i \text{pr}(x_i | pa_i), \quad (2)$$

where pa_i stands for the realisation of some subset of the variables that precede X_i in the order (X_1, X_2, \dots, X_n) . If we construct a directed acyclic graph in which the variables corresponding to pa_i are represented as the parents of X_i , also called adjacent predecessors or direct influences of X_i , then the independencies implied by the decomposition (2) can be read off the graph using the following test.

DEFINITION 1 (*d*-separation). Let X , Y and Z be three disjoint subsets of nodes in a directed acyclic graph G , and let p be any path between a node in X and a node in Y , where by 'path' we mean any succession of arcs, regardless of their directions. Then Z is said to block p if there is a node w on p satisfying one of the following two conditions: (i) w has converging arrows along p , and neither w nor any of its descendants are in Z , or, (ii) w does not have converging arrows along p , and w is in Z . Further, Z is said to *d*-separate X from Y , in G , written $(X \perp\!\!\!\perp Y | Z)_G$, if and only if Z blocks every path from a node in X to a node in Y .

It can be shown that there is a one-to-one correspondence between the set of conditional independencies $X \perp\!\!\!\perp Y | Z$ (Dawid, 1979) implied by the recursive decomposition (2), and the set of triples (X, Z, Y) that satisfy the *d*-separation criterion in G (Geiger, Verma & Pearl, 1990).

An alternative test for d -separation has been given by Lauritzen et al. (1990). To test for $(X \perp\!\!\!\perp Y|Z)_G$, delete from G all nodes except those in $X \cup Y \cup Z$ and their ancestors, connect by an edge every pair of nodes that share a common child, and remove all arrows from the arcs. Then $(X \perp\!\!\!\perp Y|Z)_G$ holds if and only if Z is a cut-set of the resulting undirected graph, separating nodes of X from those of Y . Additional properties of directed acyclic graphs and their applications to evidential reasoning in expert systems are discussed by Pearl (1988), Lauritzen & Spiegelhalter (1988), Spiegelhalter et al. (1993) and Pearl (1993a).

2.2. Graphs as models of interventions

The use of directed acyclic graphs as carriers of independence assumptions has also been instrumental in predicting the effect of interventions when these graphs are given a causal interpretation (Spirtes, Glymour & Scheines, 1993, p. 78; Pearl, 1993b). Pearl (1993b), for example, treated interventions as variables in an augmented probability space, and their effects were obtained by ordinary conditioning.

In this paper we pursue a different, though equivalent, causal interpretation of directed graphs, based on nonparametric structural equations, which owes its roots to early works in econometrics (Frisch, 1938; Haavelmo, 1943; Simon, 1953). In this account, assertions about causal influences, such as those specified by the links in Fig. 1, stand for autonomous physical mechanisms among the corresponding quantities, and these mechanisms are represented as functional relationships perturbed by random disturbances. In other words, each child-parent family in a directed graph G represents a deterministic function

$$X_i = f_i(pa_i, \varepsilon_i) \quad (i = 1, \dots, n), \quad (3)$$

where pa_i denote the parents of variable X_i in G , and ε_i ($1 \leq i \leq n$) are mutually independent, arbitrarily distributed random disturbances (Pearl & Verma, 1991). These disturbance terms represent exogenous factors that the investigator chooses not to include in the analysis. If any of these factors is judged to be influencing two or more variables, thus violating the independence assumption, then that factor must enter the analysis as an unmeasured, or latent, variable, to be represented in the graph by a hollow node, such as Z_0 or B in Fig. 1. For example, the causal assumptions conveyed by the model in Fig. 1 correspond to the following set of equations:

$$\begin{aligned} Z_0 &= f_0(\varepsilon_0), & Z_2 &= f_2(X, Z_1, \varepsilon_2), & B &= f_B(Z_0, \varepsilon_B), & Z_3 &= f_3(B, Z_2, \varepsilon_3), \\ Z_1 &= f_1(Z_0, \varepsilon_1), & Y &= f_Y(X, Z_2, Z_3, \varepsilon_Y), & X &= f_X(Z_0, \varepsilon_X). \end{aligned} \quad (4)$$

The equational model (3) is the nonparametric analogue of a structural equations model (Wright, 1921; Goldberger, 1972), with one exception: the functional form of the equations, as well as the distribution of the disturbance terms, will remain unspecified. The equality signs in such equations convey the asymmetrical counterfactual relation 'is determined by', forming a clear correspondence between causal diagrams and Rubin's model of potential outcome (Rubin, 1974; Holland, 1988; Pratt & Schlaifer, 1988; Rubin, 1990). For example, the equation for Y states that, regardless of what we currently observe about Y , and regardless of any changes that might occur in other equations, if $(X, Z_2, Z_3, \varepsilon_Y)$ were to assume the values $(x, z_2, z_3, \varepsilon_Y)$, respectively, Y would take on the value dictated by the function f_Y . Thus, the corresponding potential response variable in Rubin's model $Y_{(x)}$, the value that Y would take if X were x , becomes a deterministic function of Z_2, Z_3 and

ε_Y , whose distribution is thus determined by those of Z_2 , Z_3 and ε_Y . The relation between graphical and counterfactual models is further analysed by Pearl (1994a).

Characterising each child-parent relationship as a deterministic function, instead of by the usual conditional probability $\text{pr}(x_i|pa_i)$, imposes equivalent independence constraints on the resulting distributions, and leads to the same recursive decomposition (2) that characterises directed acyclic graph models. This occurs because each ε_i is independent of all nondescendants of X_i . However, the functional characterisation $X_i = f_i(pa_i, \varepsilon_i)$ also provides a convenient language for specifying how the resulting distribution would change in response to external interventions. This is accomplished by encoding each intervention as an alteration to a selected subset of functions, while keeping the others intact. Once we know the identity of the mechanisms altered by the intervention, and the nature of the alteration, the overall effect can be predicted by modifying the corresponding equations in the model, and using the modified model to compute a new probability function.

The simplest type of external intervention is one in which a single variable, say X_i , is forced to take on some fixed value x_i . Such an intervention, which we call atomic, amounts to lifting X_i from the influence of the old functional mechanism $X_i = f_i(pa_i, \varepsilon_i)$ and placing it under the influence of a new mechanism that sets its value to x_i while keeping all other mechanisms unperturbed. Formally, this atomic intervention, which we denote by $\text{set}(X_i = x_i)$, or $\text{set}(x_i)$ for short, amounts to removing the equation $X_i = f_i(pa_i, \varepsilon_i)$ from the model, and substituting x_i for X_i in the remaining equations. The model thus created represents the system's behaviour under the intervention $\text{set}(X_i = x_i)$ and, when solved for the distribution of X_j , yields the causal effect of X_i on X_j , denoted by $\text{pr}(x_j|\check{x}_i)$. More generally, when an intervention forces a subset X of variables to fixed values x , a subset of equations is to be pruned from the model given in (3), one for each member of X , thus defining a new distribution over the remaining variables, which completely characterises the effect of the intervention. We thus introduce the following.

DEFINITION 2 (causal effect). *Given two disjoint sets of variables, X and Y , the causal effect of X on Y , denoted $\text{pr}(y|\check{x})$, is a function from X to the space of probability distributions on Y . For each realisation x of X , $\text{pr}(y|\check{x})$ gives the probability of $Y = y$ induced on deleting from the model (3) all equations corresponding to variables in X and substituting x for X in the remainder.*

An explicit translation of intervention into 'wiping out' equations from the model was first proposed by Strotz & Wold (1960), and used by Fisher (1970) and Sobel (1990). Graphical ramifications were explicated by Spirtes et al. (1993) and Pearl (1993b). A related mathematical model using event trees has been introduced by Robins (1986, pp. 1422–5).

Regardless of whether we represent interventions as a modification of an existing model as in Definition 2, or as a conditionalisation in an augmented model (Pearl, 1993b), the result is a well-defined transformation between the pre-intervention and the post-intervention distributions. In the case of an atomic intervention $\text{set}(X_i = x'_i)$, this transformation can be expressed in a simple algebraic formula that follows immediately from (3) and Definition 2:

$$\text{pr}(x_1, \dots, x_n|\check{x}'_i) = \begin{cases} \{\text{pr}(x_1, \dots, x_n)\} / \{\text{pr}(x_i|pa_i)\} & \text{if } x_i = x'_i, \\ 0 & \text{if } x_i \neq x'_i. \end{cases} \quad (5)$$

This formula reflects the removal of the terms $\text{pr}(x_i|pa_i)$ from the product in (2), since pa_i no longer influence X_i . Graphically, this is equivalent to removing the links between pa_i

and X_i while keeping the rest of the network intact. Equation (5) can also be obtained from the G -computation formula of Robins (1986, p. 1423) and the Manipulation Theorem of Spirtes et al. (1993), who state that this formula was 'independently conjectured by Fienberg in a seminar in 1991'. Clearly, an intervention set (x_i) can affect only the descendants of X_i in G . Additional properties of the transformation defined in (5) are given by Pearl (1993b).

The immediate implication of (5) is that, given a causal diagram in which all parents of manipulated variables are observable, one can infer post-intervention distributions from pre-intervention distributions; hence, under such assumptions we can estimate the effects of interventions from passive, i.e. nonexperimental observations. The aim of this paper, however, is to derive causal effects in situations such as Fig. 1, where some members of pa_i may be unobservable, thus preventing estimation of $\text{pr}(x_i|pa_i)$. The next two sections provide simple graphical tests for deciding when $\text{pr}(x_j|\tilde{x}_i)$ is estimable in a given model.

3. CONTROLLING CONFOUNDING BIAS

3.1. *The back-door criterion*

Assume we are given a causal diagram G together with nonexperimental data on a subset V_0 of observed variables in G , and we wish to estimate what effect the intervention set $(X_i = x_i)$ would have on some response variable X_j . In other words, we seek to estimate $\text{pr}(x_j|\tilde{x}_i)$ from a sample estimate of $\text{pr}(V_0)$.

The variables in $V_0 \setminus \{X_i, X_j\}$, are commonly known as concomitants (Cox, 1958, p. 48). In observational studies, concomitants are used to reduce confounding bias due to spurious correlations between treatment and response. The condition that renders a set Z of concomitants sufficient for identifying causal effect, also known as ignorability, has been given a variety of formulations, all requiring conditional independence judgments involving counterfactual variables (Rosenbaum & Rubin, 1983; Pratt & Schlaifer, 1988). Pearl (1993b) shows that such judgments are equivalent to a simple graphical test, named the 'back-door criterion', which can be applied directly to the causal diagram.

DEFINITION 3 (*Back-door criterion*). *A set of variables Z satisfies the back-door criterion relative to an ordered pair of variables (X_i, X_j) in a directed acyclic graph G if: (i) no node in Z is a descendant of X_i , and (ii) Z blocks every path between X_i and X_j which contains an arrow into X_i . If X and Y are two disjoint sets of nodes in G , Z is said to satisfy the back-door criterion relative to (X, Y) if it satisfies it relative to any pair (X_i, X_j) such that $X_i \in X$ and $X_j \in Y$.*

The name 'back-door' echoes condition (ii), which requires that only paths with arrows pointing at X_i be blocked; these paths can be viewed as entering X_i through the back door. In Fig. 2, for example, the sets $Z_1 = \{X_3, X_4\}$ and $Z_2 = \{X_4, X_5\}$ meet the back-door criterion, but $Z_3 = \{X_4\}$ does not, because X_4 does not block the path $(X_i, X_3, X_1, X_4, X_2, X_5, X_j)$. An equivalent, though more complicated, graphical criterion is given in Theorem 7.1 of Spirtes et al. (1993). An alternative criterion using a single d -separation test will be established in § 4.4.

We summarise this finding in a theorem, after formally defining 'identifiability'.

DEFINITION 4 (*Identifiability*). *The causal effect of X on Y is said to be identifiable if the quantity $\text{pr}(y|\tilde{x})$ can be computed uniquely from any positive distribution of the observed variables that is compatible with G .*

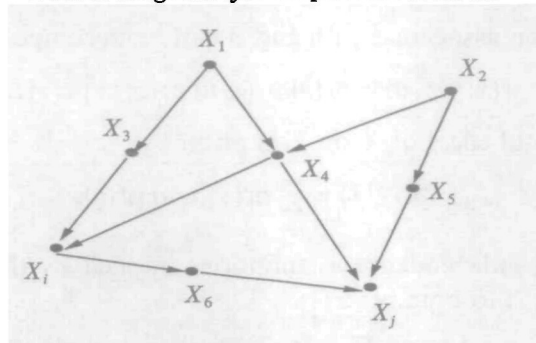


Fig. 2. A diagram representing the back-door criterion; adjusting for variables $\{X_3, X_4\}$ or $\{X_4, X_5\}$ yields a consistent estimate of $\text{pr}(x_j|\tilde{x}_i)$.

Identifiability means that $\text{pr}(y|\tilde{x})$ can be estimated consistently from an arbitrarily large sample randomly drawn from the joint distribution. To prove nonidentifiability, it is sufficient to present two sets of structural equations, both complying with (3), that induce identical distributions over observed variables but different causal effects.

THEOREM 1. *If a set of variables Z satisfies the back-door criterion relative to (X, Y) , then the causal effect of X on Y is identifiable and is given by the formula*

$$\text{pr}(y|\tilde{x}) = \sum_z \text{pr}(y|x, z) \text{pr}(z). \tag{6}$$

Equation (6) represents the standard adjustment for concomitants Z when X is conditionally ignorable given Z (Rosenbaum & Rubin, 1983). Reducing ignorability conditions to the graphical criterion of Definition 3 replaces judgments about counterfactual dependencies with systematic procedures that can be applied to causal diagrams of any size and shape. The graphical criterion also enables the analyst to search for an optimal set of concomitants, to minimise measurement cost or sampling variability.

3.2. The front-door criteria

An alternative criterion, ‘the front-door criterion’, may be applied in cases where we cannot find observed covariates Z satisfying the back-door conditions. Consider the diagram in Fig. 3. Although Z does not satisfy any of the back-door conditions, measurements of Z nevertheless enable consistent estimation of $\text{pr}(y|\tilde{x})$. This can be shown by reducing the expression for $\text{pr}(y|\tilde{x})$ to formulae computable from the observed distribution function $\text{pr}(x, y, z)$.

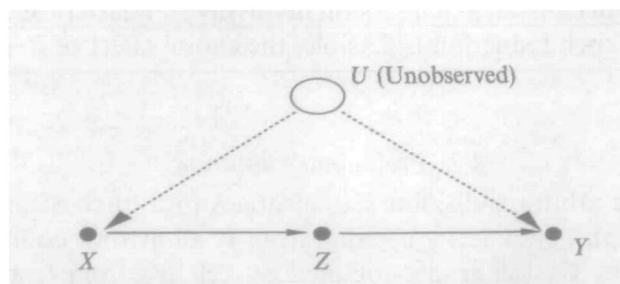


Fig. 3. A diagram representing the front-door criterion.

The joint distribution associated with Fig. 3 can be decomposed into

$$\text{pr}(x, y, z, u) = \text{pr}(u) \text{pr}(x|u) \text{pr}(z|x) \text{pr}(y|z, u) \quad (7)$$

and, from (5), the causal effect of X on Y is given by

$$\text{pr}(y|\tilde{x}) = \sum_u \text{pr}(y|x, u) \text{pr}(u). \quad (8)$$

Using the conditional independence assumptions implied by the decomposition (7), we can eliminate u from (8) to obtain

$$\text{pr}(y|\tilde{x}) = \sum_x \text{pr}(z|x) \sum_{x'} \text{pr}(y|x', z) \text{pr}(x'). \quad (9)$$

We summarise this result by a theorem.

THEOREM 2. *Suppose a set of variables Z satisfies the following conditions relative to an ordered pair of variables (X, Y) : (i) Z intercepts all directed paths from X to Y , (ii) there is no back-door path between X and Z , and (iii) every back-door path between Z and Y is blocked by X . Then the causal effect of X on Y is identifiable and is given by (9).*

The graphical criterion of Theorem 2 uncovers many new structures that permit the identification of causal effects from measurements of variables that are affected by treatments: see § 5.2. The relevance of such structures in practical situations can be seen, for instance, if we identify X with smoking, Y with lung cancer, Z with the amount of tar deposited in a subject's lungs, and U with an unobserved carcinogenic genotype that, according to some, also induces an inborn craving for nicotine. In this case, (9) would provide us with the means to quantify, from nonexperimental data, the causal effect of smoking on cancer, assuming, of course, that $\text{pr}(x, y, z)$ is available and that we believe that smoking does not have any direct effect on lung cancer except that mediated by tar deposits.

4. A CALCULUS OF INTERVENTION

4.1. General

This section establishes a set of inference rules by which probabilistic sentences involving interventions and observations can be transformed into other such sentences, thus providing a syntactic method of deriving or verifying claims about interventions. We shall assume that we are given the structure of a causal diagram G in which some of the nodes are observable while the others remain unobserved. Our main problem will be to facilitate the syntactic derivation of causal effect expressions of the form $\text{pr}(y|\tilde{x})$, where X and Y denote sets of observed variables. By derivation we mean step-wise reduction of the expression $\text{pr}(y|\tilde{x})$ to an equivalent expression involving standard probabilities of observed quantities. Whenever such reduction is feasible, the causal effect of X on Y is identifiable: see Definition 4.

4.2. Preliminary notation

Let X , Y and Z be arbitrary disjoint sets of nodes in a directed acyclic graph G . We denote by $G_{\tilde{X}}$ the graph obtained by deleting from G all arrows pointing to nodes in X . Likewise, we denote by $G_{\tilde{X}}$ the graph obtained by deleting from G all arrows emerging from nodes in X . To represent the deletion of both incoming and outgoing arrows, we

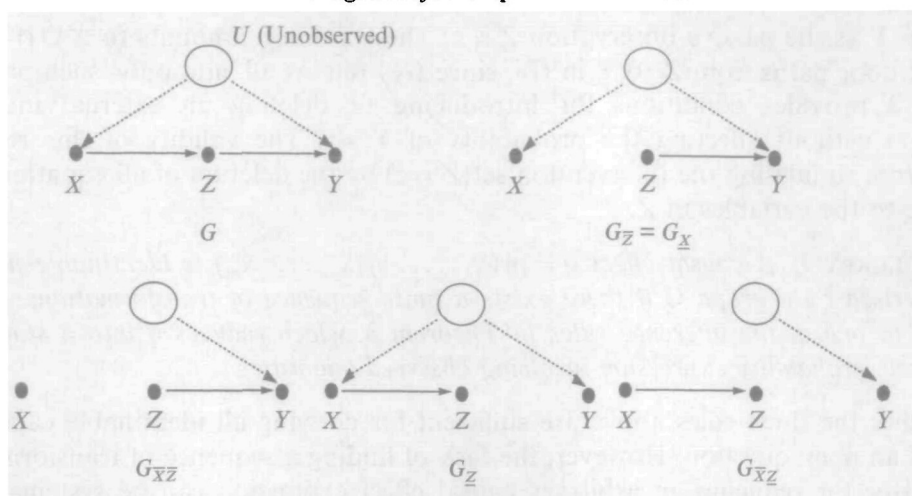


Fig. 4. Subgraphs of G used in the derivation of causal effects.

use the notation G_{XZ} : see Fig. 4 for illustration. Finally, $\text{pr}(y|\check{x}, z) := \text{pr}(y, z|\check{x})/\text{pr}(z|\check{x})$ denotes the probability of $Y = y$ given that $Z = z$ is observed and X is held constant at x .

4.3. Inference rules

The following theorem states the three basic inference rules of the proposed calculus. Proofs are provided in the Appendix.

THEOREM 3. *Let G be the directed graph associated with a causal model as defined in (3), and let $\text{pr}(\cdot)$ stand for the probability distribution induced by that model. For any disjoint subsets of variables X, Y, Z and W we have the following.*

Rule 1 (insertion/deletion of observations):

$$\text{pr}(y|\check{x}, z, w) = \text{pr}(y|\check{x}, w) \quad \text{if } (Y \perp\!\!\!\perp Z | X, W)_{G_{\check{X}}}. \tag{10}$$

Rule 2 (action/observation exchange):

$$\text{pr}(y|\check{x}, \check{z}, w) = \text{pr}(y|\check{x}, z, w) \quad \text{if } (Y \perp\!\!\!\perp Z | X, W)_{G_{\check{X}\check{Z}}}. \tag{11}$$

Rule 3 (insertion/deletion of actions):

$$\text{pr}(y|\check{x}, \check{z}, w) = \text{pr}(y|\check{x}, w) \quad \text{if } (Y \perp\!\!\!\perp Z | X, W)_{G_{\check{Z}, \overline{W}}}, \tag{12}$$

where $Z(W)$ is the set of Z -nodes that are not ancestors of any W -node in $G_{\check{X}}$.

Each of the inference rules above follows from the basic interpretation of the ‘ \check{x} ’ operator as a replacement of the causal mechanism that connects X to its pre-intervention parents by a new mechanism $X = x$ introduced by intervening force. The result is a submodel characterised by the subgraph $G_{\check{X}}$, called the ‘manipulated graph’ by Spirtes et al. (1993), which supports all three rules.

Rule 1 reaffirms d -separation as a valid test for conditional independence in the distribution resulting from the intervention set ($X = x$), hence the graph $G_{\check{X}}$. This rule follows from the fact that deleting equations from the system does not introduce any dependencies among the remaining disturbance terms: see (3).

Rule 2 provides a condition for an external intervention set ($Z = z$) to have the same

effect on Y as the passive observation $Z = z$. The condition amounts to $X \cup W$ blocking all back-door paths from Z to Y in $G_{\underline{X}}$, since $G_{\underline{XZ}}$ retains all, and only, such paths.

Rule 3 provides conditions for introducing or deleting an external intervention set($Z = z$) without affecting the probability of $Y = y$. The validity of this rule stems, again, from simulating the intervention set($Z = z$) by the deletion of all equations corresponding to the variables in Z .

COROLLARY 1. *A causal effect $q = \text{pr}(y_1, \dots, y_k | \check{x}_1, \dots, \check{x}_m)$ is identifiable in a model characterised by a graph G if there exists a finite sequence of transformations, each conforming to one of the inference rules in Theorem 3, which reduces q into a standard, i.e. check-free, probability expression involving observed quantities.*

Whether the three rules above are sufficient for deriving all identifiable causal effects remains an open question. However, the task of finding a sequence of transformations, if such exists, for reducing an arbitrary causal effect expression can be systematised and executed by efficient algorithms as described by Galles & Pearl (1995). As § 4.4 illustrates, symbolic derivations using the check notation are much more convenient than algebraic derivations that aim at eliminating latent variables from standard probability expressions, as in § 3.2.

4.4. Symbolic derivation of causal effects: An example

We now demonstrate how Rules 1–3 can be used to derive causal effect estimands in the structure of Fig. 3 above. Figure 4 displays the subgraphs that will be needed for the derivations that follow.

Task 1: compute $\text{pr}(z | \check{x})$. This task can be accomplished in one step, since G satisfies the applicability condition for Rule 2, namely, $X \perp\!\!\!\perp Z$ in $G_{\underline{X}}$, because the path $X \leftarrow U \rightarrow Y \leftarrow Z$ is blocked by the converging arrows at Y , and we can write

$$\text{pr}(z | \check{x}) = \text{pr}(z | x). \quad (13)$$

Task 2: compute $\text{pr}(y | \check{z})$. Here we cannot apply Rule 2 to exchange \check{z} with z because $G_{\underline{Z}}$ contains a back-door path from Z to Y : $Z \leftarrow X \leftarrow U \rightarrow Y$. Naturally, we would like to block this path by measuring variables, such as X , that reside on that path. This involves conditioning and summing over all values of X :

$$\text{pr}(y | \check{z}) = \sum_x \text{pr}(y | x, \check{z}) \text{pr}(x | \check{z}). \quad (14)$$

We now have to deal with two expressions involving \check{z} , $\text{pr}(y | x, \check{z})$ and $\text{pr}(x | \check{z})$. The latter can be readily computed by applying Rule 3 for action deletion:

$$\text{pr}(x | \check{z}) = \text{pr}(x) \quad \text{if } (Z \perp\!\!\!\perp X)_{G_{\underline{Z}}}, \quad (15)$$

since X and Z are d -separated in $G_{\underline{Z}}$. Intuitively, manipulating Z should have no effect on X , because Z is a descendant of X in G . To reduce $\text{pr}(y | x, \check{z})$, we consult Rule 2:

$$\text{pr}(y | x, \check{z}) = \text{pr}(y | x, z) \quad \text{if } (Z \perp\!\!\!\perp Y | X)_{G_{\underline{Z}}}, \quad (16)$$

noting that X d -separates Z from Y in $G_{\underline{Z}}$. This allows us to write (14) as

$$\text{pr}(y | \check{z}) = \sum_x \text{pr}(y | x, z) \text{pr}(x) = E_x \text{pr}(y | x, z), \quad (17)$$

which is a special case of the back-door formula (6). The legitimising condition, $(Z \perp\!\!\!\perp Y|X)_{G_{\underline{Z}}}$, offers yet another graphical test for the ignorability condition of Rosenbaum & Rubin (1983).

Task 3: compute $\text{pr}(y|\check{x})$. Writing

$$\text{pr}(y|\check{x}) = \sum_z \text{pr}(y|z, \check{x}) \text{pr}(z|\check{x}), \quad (18)$$

we see that the term $\text{pr}(z|\check{x})$ was reduced in (13) but that no rule can be applied to eliminate the ‘check’ symbol from the term $\text{pr}(y|z, \check{x})$. However, we can add a ‘check’ symbol to this term via Rule 2:

$$\text{pr}(y|z, \check{x}) = \text{pr}(y|z, \check{x}), \quad (19)$$

since the applicability condition $(Y \perp\!\!\!\perp Z|X)_{G_{\underline{XZ}}}$, holds true. We can now delete the action \check{x} from $\text{pr}(y|z, \check{x})$ using Rule 3, since $Y \perp\!\!\!\perp X|\check{Z}$ holds in $G_{\overline{XZ}}$. Thus, we have

$$\text{pr}(y|z, \check{x}) = \text{pr}(y|z), \quad (20)$$

which was calculated in (17). Substituting (17), (20) and (13) back into (18) finally yields

$$\text{pr}(y|\check{x}) = \sum_z \text{pr}(z|x) \sum_{x'} \text{pr}(y|x', z) \text{pr}(x'), \quad (21)$$

which is identical to the front-door formula (9).

The reader may verify that all other causal effects, for example, $\text{pr}(y, z|\check{x})$ and $\text{pr}(x, z|\check{y})$, can likewise be derived through the rules of Theorem 3. Note that in all the derivations the graph G provides both the license for applying the inference rules and the guidance for choosing the right rule to apply.

4.5. Causal inference by surrogate experiments

Suppose we wish to learn the causal effect of X on Y when $\text{pr}(y|\check{x})$ is not identifiable and, for practical reasons of cost or ethics, we cannot control X by randomised experiment. The question arises whether $\text{pr}(y|\check{x})$ can be identified by randomising a surrogate variable Z , which is easier to control than X . For example, if we are interested in assessing the effect of cholesterol levels X on heart disease, Y , a reasonable experiment to conduct would be to control subjects’ diet, Z , rather than exercising direct control over cholesterol levels in subjects’ blood.

Formally, this problem amounts to transforming $\text{pr}(y|\check{x})$ into expressions in which only members of Z carry the check symbol. Using Theorem 3 it can be shown that the following conditions are sufficient for admitting a surrogate variable Z : (i) X intercepts all directed paths from Z to Y , and (ii) $\text{pr}(y|\check{x})$ is identifiable in $G_{\underline{Z}}$. Indeed, if condition (i) holds, we can write $\text{pr}(y|\check{x}) = \text{pr}(y|\check{x}, \check{z})$, because $(Y \perp\!\!\!\perp Z|X)_{G_{\overline{YZ}}}$. But $\text{pr}(y|\check{x}, \check{z})$ stands for the causal effect of X on Y in a model governed by $G_{\underline{Z}}$ which, by condition (ii), is identifiable. Figures 7(e) and 7(h) below illustrate models in which both conditions hold. Translated to our cholesterol example, these conditions require that there be no direct effect of diet on heart conditions and no confounding effect between cholesterol levels and heart disease, unless we can measure an intermediate variable between the two.

5. GRAPHICAL TESTS OF IDENTIFIABILITY

5.1. General

Figure 5 shows simple diagrams in which $\text{pr}(y|\tilde{x})$ cannot be identified due to the presence of a bow pattern, i.e. a confounding arc, shown dashed, embracing a causal link between X and Y . A confounding arc represents the existence in the diagram of a back-door path that contains only unobserved variables and has no converging arrows. For example, the path X, Z_0, B, Z_3 in Fig. 1 can be represented as a confounding arc between X and Z_3 . A bow-pattern represents an equation $Y = f_Y(X, U, \varepsilon_Y)$, where U is unobserved and dependent on X . Such an equation does not permit the identification of causal effects since any portion of the observed dependence between X and Y may always be attributed to spurious dependencies mediated by U .

The presence of a bow-pattern prevents the identification of $\text{pr}(y|\tilde{x})$ even when it is found in the context of a larger graph, as in Fig. 5(b). This is in contrast to linear models, where the addition of an arc to a bow-pattern can render $\text{pr}(y|\tilde{x})$ identifiable. For example, if Y is related to X via a linear relation $Y = bX + U$, where U is an unobserved disturbance possibly correlated with X , then $b = \partial E(Y|\tilde{x})/\partial x$ is not identifiable. However, adding an arc $Z \rightarrow X$ to the structure, that is, finding a variable Z that is correlated with X but not with U , would facilitate the computation of $E(Y|\tilde{x})$ via the instrumental-variable formula (Bowden & Turkington, 1984, p. 12; Angrist, Imbens & Rubin, 1995):

$$b := \frac{\partial}{\partial x} E(Y|\tilde{x}) = \frac{E(Y|z)}{E(X|z)} = \frac{R_{yz}}{R_{xz}}. \quad (22)$$

In nonparametric models, adding an instrumental variable Z to a bow-pattern, see Fig. 5(b), does not permit the identification of $\text{pr}(y|\tilde{x})$. This is a familiar problem in the analysis of clinical trials in which treatment assignment, Z , is randomised, hence no link enters Z , but compliance is imperfect. The confounding arc between X and Y in Fig. 5(b) represents unmeasurable factors which influence both subjects' choice of treatment, X , and response to treatment, Y . In such trials, it is not possible to obtain an unbiased estimate of the treatment effect $\text{pr}(y|\tilde{x})$ without making additional assumptions on the nature of the interactions between compliance and response (Imbens & Angrist, 1994), as is done, for example, in the approach to instrumental variables developed by Angrist et al. (1995). While the added arc $Z \rightarrow X$ permits us to calculate bounds on $\text{pr}(y|\tilde{x})$ (Robins, 1989, § 1g; Manski, 1990), and while the upper and lower bounds may even coincide for

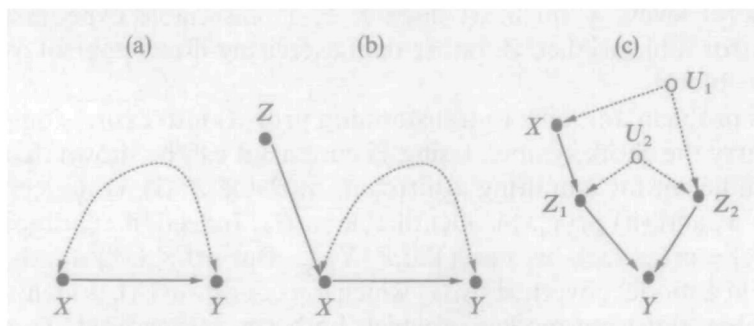


Fig. 5. (a) A bow-pattern: a confounding arc embracing a causal link $X \rightarrow Y$, thus preventing the identification of $\text{pr}(y|\tilde{x})$ even in the presence of an instrumental variable Z , as in (b). (c) A bow-less graph still prohibiting the identification of $\text{pr}(y|\tilde{x})$.

certain types of distributions $\text{pr}(x, y, z)$ (Balke & Pearl, 1994), there is no way of computing $\text{pr}(y|\tilde{x})$ for every positive distribution $\text{pr}(x, y, z)$, as required by Definition 4.

In general, the addition of arcs to a causal diagram can impede, but never assist, the identification of causal effects in nonparametric models. This is because such addition reduces the set of d -separation conditions carried by the diagram and, hence, if a causal effect derivation fails in the original diagram, it is bound to fail in the augmented diagram as well. Conversely, any causal effect derivation that succeeds in the augmented diagram, by a sequence of symbolic transformations, as in Corollary 1, would succeed in the original diagram.

Our ability to compute $\text{pr}(y|\tilde{x})$ for pairs (x, y) of singleton variables does not ensure our ability to compute joint distributions, such as $\text{pr}(y_1, y_2|\tilde{x})$. Figure 5(c), for example, shows a causal diagram where both $\text{pr}(z_1|\tilde{x})$ and $\text{pr}(z_2|\tilde{x})$ are computable, but $\text{pr}(z_1, z_2|\tilde{x})$ is not. Consequently, we cannot compute $\text{pr}(y|\tilde{x})$. This diagram is the smallest graph that does not contain a bow-pattern and still presents an uncomputable causal effect.

5.2. Identifying models

Figure 6 shows simple diagrams in which the causal effect of X on Y is identifiable. Such models are called identifying because their structures communicate a sufficient number of assumptions to permit the identification of the target quantity $\text{pr}(y|\tilde{x})$. Latent variables are not shown explicitly in these diagrams; rather, such variables are implicit in the confounding arcs, shown dashed. Every causal diagram with latent variables can be converted to an equivalent diagram involving measured variables interconnected by arrows and confounding arcs. This conversion corresponds to substituting out all latent variables from the structural equations of (3) and then constructing a new diagram by

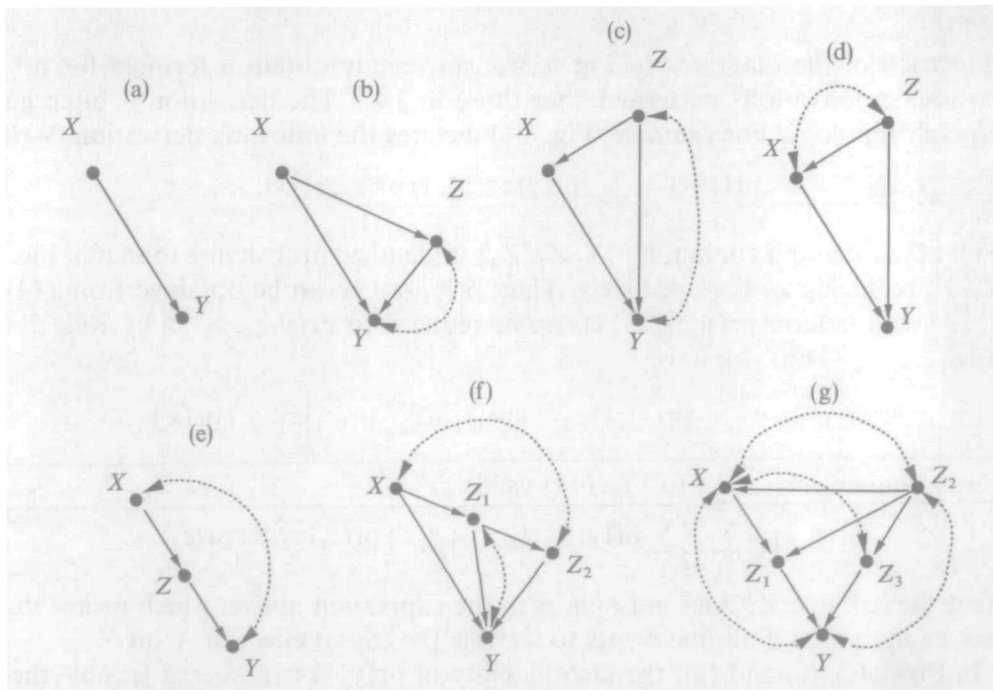


Fig. 6. Typical models in which the effect of X on Y is identifiable. Dashed arcs represent confounding paths, and Z represents observed covariates.

connecting any two variables X_i and X_j by (i) an arrow from X_j to X_i whenever X_j appears in the equation for X_i , and (ii) a confounding arc whenever the same ε term appears in both f_i and f_j . The result is a diagram in which all unmeasured variables are exogenous and mutually independent. Several features should be noted from examining the diagrams in Fig. 6.

(i) Since the removal of any arc or arrow from a causal diagram can only assist the identifiability of causal effects, $\text{pr}(y|\tilde{x})$ will still be identified in any edge-subgraph of the diagrams shown in Fig. 6. Likewise, the introduction of mediating observed variables onto any edge in a causal graph can assist, but never impede, the identifiability of any causal effect. Therefore, $\text{pr}(y|\tilde{x})$ will still be identified from any graph obtained by adding mediating nodes to the diagrams shown in Fig. 6.

(ii) The diagrams in Fig. 6 are maximal, in the sense that the introduction of any additional arc or arrow onto an existing pair of nodes would render $\text{pr}(y|\tilde{x})$ no longer identifiable.

(iii) Although most of the diagrams in Fig. 6 contain bow-patterns, none of these patterns emanates from X as is the case in Fig. 7(a) and (b) below. In general, a necessary condition for the identifiability of $\text{pr}(y|\tilde{x})$ is the absence of a confounding arc between X and any child of X that is an ancestor of Y .

(iv) Figures 6(a) and (b) contain no back-door paths between X and Y , and thus represent experimental designs in which there is no confounding bias between the treatment, X , and the response, Y ; that is, X is strongly ignorable relative to Y (Rosenbaum & Rubin, 1983); hence, $\text{pr}(y|\tilde{x}) = \text{pr}(y|x)$. Likewise, Figs 6(c) and (d) represent designs in which observed covariates, Z , block every back-door path between X and Y ; that is X is conditionally ignorable given Z (Rosenbaum & Rubin, 1983); hence, $\text{pr}(y|\tilde{x})$ is obtained by standard adjustment for Z , as in (6):

$$\text{pr}(y|\tilde{x}) = \sum_z \text{pr}(y|x, z) \text{pr}(z).$$

(v) For each of the diagrams in Fig. 6, we can readily obtain a formula for $\text{pr}(y|\tilde{x})$, using symbolic derivations patterned after those in § 4.4. The derivation is often guided by the graph topology. For example, Fig. 6(f) dictates the following derivation. Writing

$$\text{pr}(y|\tilde{x}) = \sum_{z_1, z_2} \text{pr}(y|z_1, z_2, \tilde{x}) \text{pr}(z_1, z_2|\tilde{x}),$$

we see that the subgraph containing $\{X, Z_1, Z_2\}$ is identical in structure to that of Fig. 6(e), with Z_1, Z_2 replacing Z, Y , respectively. Thus, $\text{pr}(z_1, z_2|\tilde{x})$ can be obtained from (14) and (21). Likewise, the term $\text{pr}(y|z_1, z_2, \tilde{x})$ can be reduced to $\text{pr}(y|z_1, z_2, x)$ by Rule 2, since $(Y \perp\!\!\!\perp X | Z_1, Z_2)_{G_{\tilde{x}}}$. Thus, we have

$$\text{pr}(y|\tilde{x}) = \sum_{z_1, z_2} \text{pr}(y|z_1, z_2, x) \text{pr}(z_1|x) \sum_{x'} \text{pr}(z_2|z_1, x') \text{pr}(x'). \quad (23)$$

Applying a similar derivation to Fig. 6(g) yields

$$\text{pr}(y|\tilde{x}) = \sum_{z_1} \sum_{z_2} \sum_{x'} \text{pr}(y|z_1, z_2, x') \text{pr}(x') \text{pr}(z_1|z_2, x) \text{pr}(z_2). \quad (24)$$

Note that the variable Z_3 does not appear in the expression above, which means that Z_3 need not be measured if all one wants to learn is the causal effect of X on Y .

(vi) In Figs 6(e), (f) and (g), the identifiability of $\text{pr}(y|\tilde{x})$ is rendered feasible through observed covariates, Z , that are affected by the treatment X , that is descendants of X . This stands contrary to the warning, repeated in most of the literature on statistical

experimentation, to refrain from adjusting for concomitant observations that are affected by the treatment (Cox, 1958, p. 48; Rosenbaum, 1984; Pratt & Schlaifer, 1988; Wainer, 1989). It is commonly believed that, if a concomitant Z is affected by the treatment, then it must be excluded from the analysis of the total effect of the treatment (Pratt & Schlaifer, 1988). The reasons given for the exclusion is that the calculation of total effects amounts to integrating out Z , which is functionally equivalent to omitting Z to begin with. Figures 6(e), (f) and (g) show cases where one wants to learn the total effects of X and, still, the measurement of concomitants that are affected by X , for example Z or Z_1 , is necessary. However, the adjustment needed for such concomitants is nonstandard, involving two or more stages of the standard adjustment of (6): see (9), (23) and (24).

(vii) In Figs 6(b), (c) and (f), Y has a parent whose effect on Y is not identifiable, yet the effect of X on Y is identifiable. This demonstrates that local identifiability is not a necessary condition for global identifiability. In other words, to identify the effect of X on Y we need not insist on identifying each and every link along the paths from X to Y .

5.3. Nonidentifying models

Figure 7 presents typical diagrams in which the total effect of X on Y , $\text{pr}(y|\dot{x})$, is not identifiable. Noteworthy features of these diagrams are as follows.

(i) All graphs in Fig. 7 contain unblockable back-door paths between X and Y , that is, paths ending with arrows pointing to X which cannot be blocked by observed nondescendants of X . The presence of such a path in a graph is, indeed, a necessary test for nonidentifiability. It is not a sufficient test, though, as is demonstrated by Fig. 6(e), in which the back-door path (dashed) is unblockable, yet $\text{pr}(y|\dot{x})$ is identifiable.

(ii) A sufficient condition for the nonidentifiability of $\text{pr}(y|\dot{x})$ is the existence of a confounding path between X and any of its children on a path from X to Y , as shown in Figs 7(b) and (c). A stronger sufficient condition is that the graph contain any of the patterns shown in Fig. 7 as an edge-subgraph.

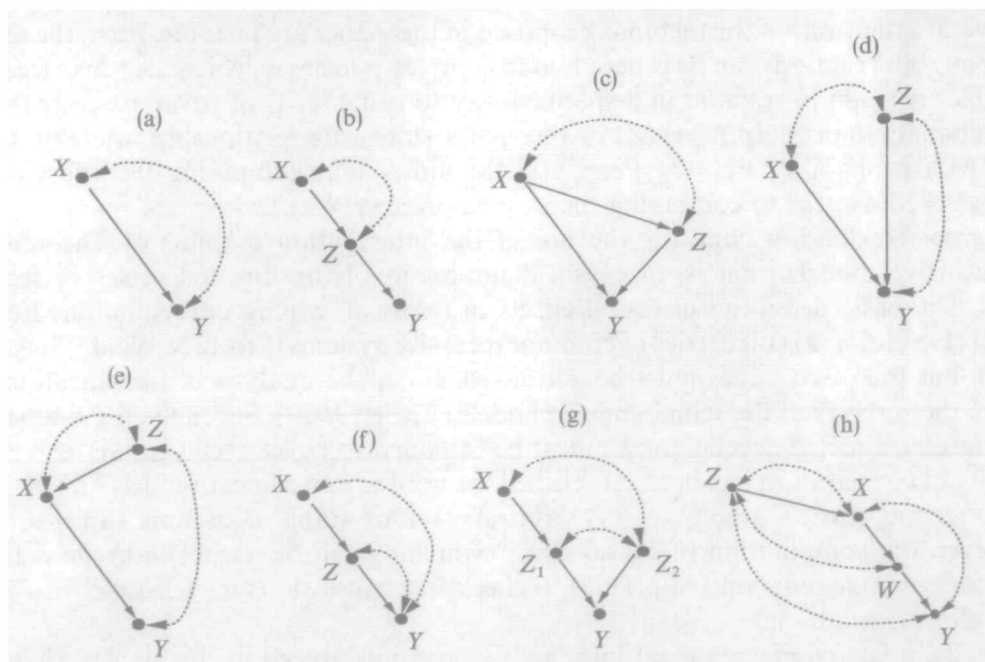


Fig. 7. Typical models in which $\text{pr}(y|\dot{x})$ is not identifiable.

(iii) Figure 7(g) demonstrates that local identifiability is not sufficient for global identifiability. For example, we can identify $\text{pr}(z_1|\tilde{x})$, $\text{pr}(z_2|\tilde{x})$, $\text{pr}(y|z_1)$ and $\text{pr}(y|z_2)$, but not $\text{pr}(y|\tilde{x})$. This is one of the main differences between nonparametric and linear models; in the latter, all causal effects can be determined from the structural coefficients, each coefficient representing the causal effect of one variable on its immediate successor.

6. DISCUSSION

The basic limitation of the methods proposed in this paper is that the results must rest on the causal assumptions shown in the graph, and that these cannot usually be tested in observational studies. In related papers (Pearl, 1994a, 1995) we show that some of the assumptions, most notably those associated with instrumental variables, see Fig. 5(b), are subject to falsification tests. Additionally, considering that any causal inferences from observational studies must ultimately rely on some kind of causal assumptions, the methods described in this paper offer an effective language for making those assumptions precise and explicit, so they can be isolated for deliberation or experimentation and, once validated, integrated with statistical data.

A second limitation concerns an assumption inherent in identification analysis, namely, that the sample size is so large that sampling variability may be ignored. The mathematical derivation of causal-effect estimands should therefore be considered a first step toward supplementing estimates of these with confidence intervals and significance levels, as in traditional analysis of controlled experiments. Having nonparametric estimates for causal effects does not imply that one should refrain from using parametric forms in the estimation phase of the study. For example, if the assumptions of Gaussian, zero-mean disturbances and linearity are deemed reasonable, then the estimand in (9) can be replaced by $E(Y|\tilde{x}) = R_{xz}\beta_{zy \cdot x}x$, where $\beta_{zy \cdot x}$ is the standardised regression coefficient, and the estimation problem reduces to that of estimating coefficients. More sophisticated estimation techniques are given by Rubin (1978), Robins (1989, § 17), and Robins et al. (1992, pp. 331–3).

Several extensions of the methods proposed in this paper are possible. First, the analysis of atomic interventions can be generalised to complex policies in which a set X of treatment variables is made to respond in a specified way to some set Z of covariates, say through a functional relationship $X = g(Z)$ or through a stochastic relationship whereby X is set to x with probability $P^*(x|z)$. Pearl (1994b) shows that computing the effect of such policies is equivalent to computing the expression $\text{pr}(y|\tilde{x}, z)$.

A second extension concerns the use of the intervention calculus of Theorem 3 in nonrecursive models, that is, in causal diagrams involving directed cycles or feedback loops. The basic definition of causal effects in terms of ‘wiping out’ equations from the model (Definition 2) still carries over to nonrecursive systems (Strotz & Wold, 1960; Sobel, 1990), but then two issues must be addressed. First, the analysis of identification must ensure the stability of the remaining submodels (Fisher, 1970). Secondly, the d -separation criterion for directed acyclic graphs must be extended to cover cyclic graphs as well. The validity of d -separation has been established for nonrecursive linear models and extended, using an augmented graph, to any arbitrary set of stable equations (Spirtes, 1995). However, the computation of causal effect estimands will be harder in cyclic networks, because symbolic reduction of $\text{pr}(y|\tilde{x})$ to check-free expressions may require the solution of nonlinear equations.

Finally, a few comments regarding the notation introduced in this paper. There have been three approaches to expressing causal assumptions in mathematical form. The most

common approach in the statistical literature invokes Rubin's model (Rubin, 1974), in which probability functions are defined over an augmented space of observable and counterfactual variables. In this model, causal assumptions are expressed as independence constraints over the augmented probability function, as exemplified by Rosenbaum & Rubin's (1983) definitions of ignorability conditions. An alternative but related approach, still using the standard language of probability, is to define augmented probability functions over variables representing hypothetical interventions (Pearl, 1993b).

The language of structural models, which includes path diagrams (Wright, 1921) and structural equations (Goldberger, 1972) represents a drastic departure from these two approaches, because it invokes new primitives, such as arrows, disturbance terms, or plain causal statements, which have no parallels in the language of probability. This language has been very popular in the social sciences and econometrics, because it closely echoes statements made in ordinary scientific discourse and thus provides a natural way for scientists to communicate knowledge and experience, especially in situations involving many variables.

Statisticians, however, have generally found structural models suspect, because the empirical content of basic notions in these models appears to escape conventional methods of explication. For example, analysts have found it hard to conceive of experiments, however hypothetical, whose outcomes would be constrained by a given structural equation. Standard probability calculus cannot express the empirical content of the coefficient b in the structural equation $Y = bX + \varepsilon_Y$ even if one is prepared to assume that ε_Y , an unobserved quantity, is uncorrelated with X . Nor can any probabilistic meaning be attached to the analyst's excluding from this equation certain variables that are highly correlated with X or Y . As a consequence, the whole enterprise of structural equation modelling has become the object of serious controversy and misunderstanding among researchers (Freedman, 1987; Wermuth, 1992; Whittaker, 1990, p. 302; Cox & Wermuth, 1993).

To a large extent, this history of controversy stems not from faults in the structural modelling approach but rather from a basic limitation of standard probability theory: when viewed as a mathematical language, it is too weak to describe the precise experimental conditions that prevail in a given study. For example, standard probabilistic notation cannot distinguish between an experiment in which variable X is observed to take on value x and one in which variable X is set to value x by some external control. The need for this distinction was recognised by several researchers, most notably Pratt & Schlaifer (1988) and Cox (1992), but has not led to a more refined and manageable mathematical notation capable of reflecting this distinction.

The 'check' notation developed in this paper permits one to specify precisely what is being held constant and what is merely measured in a given study and, using this specification, the basic notions of structural models can be given clear empirical interpretation. For example, the meaning of b in the equation $Y = bX + \varepsilon_Y$ is simply $\partial E(Y|\check{x})/\partial x$, namely, the rate of change, in x , of the expectation of Y in an experiment where X is held at x by external control. This interpretation holds regardless of whether ε_Y and X are correlated, for example, via another equation: $X = aY + \varepsilon_X$. Moreover, the notion of randomisation need not be invoked. Likewise, the analyst's decision as to which variables should be included in a given equation is based on a hypothetical controlled experiment: a variable Z is excluded from the equation for Y if it has no influence on Y when all other variables, S_{YZ} , are held constant, that is, $\text{pr}(y|\check{z}, \check{s}_{YZ}) = \text{pr}(y|\check{s}_{YZ})$. In other words, variables that are excluded from the equation $Y = bX + \varepsilon_Y$ are not conditionally independent of Y given measurements of X , but rather conditionally independent of Y given settings of X . The

operational meaning of the so-called ‘disturbance term’, ε_Y , is likewise demystified: ε_Y is defined as the difference $Y - E(Y|\xi_Y)$; two disturbance terms, ε_X and ε_Y , are correlated if $\text{pr}(y|\tilde{x}, \xi_{XY}) \neq \text{pr}(y|x, \xi_{XY})$; and so on.

The distinctions provided by the ‘check’ notation clarify the empirical basis of structural equations and should make causal models more acceptable to empirical researchers. Moreover, since most scientific knowledge is organised around the operation of ‘holding X fixed’, rather than ‘conditioning on X ’, the notation and calculus developed in this paper should provide an effective means for scientists to communicate subject-matter information, and to infer its logical consequences when combined with statistical data.

ACKNOWLEDGEMENT

Much of this investigation was inspired by Spirtes et al. (1993), in which a graphical account of manipulations was first proposed. Phil Dawid, David Freedman, James Robins and Donald Rubin have provided genuine encouragement and valuable advice. The investigation also benefitted from discussions with Joshua Angrist, Peter Bentler, David Cox, Arthur Dempster, David Galles, Arthur Goldberger, Sander Greenland, David Hendry, Paul Holland, Guido Imbens, Ed Leamer, Rod McDonald, John Pratt, Paul Rosenbaum, Keunkwan Ryu, Glenn Shafer, Michael Sobel, David Tritchler and Nanny Wermuth. The research was partially supported by grants from Air Force Office of Scientific Research and National Science Foundation.

APPENDIX

Proof of Theorem 3

(i) Rule 1 follows from the fact that deleting equations from the model in (8) results, again, in a recursive set of equations in which all ε terms are mutually independent. The d -separation condition is valid for any recursive model, hence it is valid for the submodel resulting from deleting the equations for X . Finally, since the graph characterising this submodel is given by $G_{\tilde{X}}$, $(Y \perp\!\!\!\perp Z|X, W)_{G_{\tilde{X}}}$ implies the conditional independence $\text{pr}(y|\tilde{x}, z, w) = \text{pr}(y|\tilde{x}, w)$ in the post-intervention distribution.

(ii) The graph G_{XZ} differs from G_X only in lacking the arrows emanating from Z , hence it retains all the back-door paths from Z to Y that can be found in G_X . The condition $(Y \perp\!\!\!\perp Z|X, W)_{G_{XZ}}$ ensures that all back-door paths from Z to Y in G_X are blocked by $\{X, W\}$. Under such conditions, setting $Z = z$ or conditioning on $Z = z$ has the same effect on Y . This can best be seen from the augmented diagram $G_{\tilde{X}}$, to which the intervention arcs $F_Z \rightarrow Z$ were added, where F_Z stands for the functions that determine Z in the structural equations (Pearl, 1993b). If all back-door paths from F_Z to Y are blocked, the remaining paths from F_Z to Y must go through the children of Z , hence these paths will be blocked by Z . The implication is that Y is independent of F_Z given Z , which means that the observation $Z = z$ cannot be distinguished from the intervention $F_Z = \text{set}(z)$.

(iii) The following argument was developed by D. Galles. Consider the augmented diagram $G_{\tilde{X}}$ to which the intervention arcs $F_Z \rightarrow Z$ are added. If $(F_Z \perp\!\!\!\perp Y|W, X)_{G_{\tilde{X}}}$, then $\text{pr}(y|\tilde{x}, z, w) = \text{pr}(y|\tilde{x}, w)$. If $(Y \perp\!\!\!\perp Z|X, W)_{G_{\tilde{X}} \setminus Z(W)}$ and $(F_Z \not\perp\!\!\!\perp Y|W, X)_{G_{\tilde{X}}}$, there must be an unblocked path from a member $F_{Z'}$ of F_Z to Y that passes either through a head-to-tail junction at Z' , or a head-to-head junction at Z' . If there is such a path, let P be the shortest such path. We will show that P will violate some premise, or there exists a shorter path, either of which leads to a contradiction.

If the junction is head-to-tail, that means that $(Y \not\perp\!\!\!\perp Z'|W, X)_{G_{\tilde{X}}}$ but $(Y \perp\!\!\!\perp Z'|W, X)_{G_{\tilde{X}} \setminus Z(W)}$. So, there must be an unblocked path from Y to Z' that passes through some member Z'' of $Z(W)$ in either a head-to-head or a tail-to-head junction. This is impossible. If the junction is head-to-head, then some descendant of Z'' must be in W for the path to be unblocked, but then Z'' would not

be in $Z(W)$. If the junction is tail-to-head, there are two options: either the path from Z' to Z'' ends in an arrow pointing to Z'' , or in an arrow pointing away from Z'' . If it ends in an arrow pointing away from Z'' , then there must be a head-to-head junction along the path from Z' to Z'' . In that case, for the path to be unblocked, W must be a descendant of Z'' , but then Z'' would not be in $Z(W)$. If it ends in an arrow pointing to Z'' , then there must be an unblocked path from Z' to Y in $G_{\bar{X}}$ that is blocked in $G_{\bar{X}\bar{Z}(W)}$. If this is true, then there is an unblocked path from $F_{Z''}$ to Y that is shorter than P , the shortest path.

If the junction through Z' is head-to-head, then either Z' is in $Z(W)$, in which case that junction would be blocked, or there is an unblocked path from Z' to Y in $G_{\bar{X}\bar{Z}(W)}$ that is blocked in $G_{\bar{X}}$. Above, we proved that this could not occur. So $(Y \perp\!\!\!\perp Z | X, W)_{G_{\bar{X}\bar{Z}(W)}}$ implies $(F_Z \perp\!\!\!\perp Y | W, X)_{G_{\bar{X}}}$, and thus $\text{pr}(y | \bar{x}, \bar{z}, w) = \text{pr}(y | \bar{x}, w)$.

REFERENCES

- ANGRIST, J. D., IMBENS, G. W. & RUBIN, D. B. (1995). Identification of causal effects using instrumental variables. *J. Am. Statist. Assoc.* To appear.
- BALKE, A. & PEARL, J. (1994). Counterfactual probabilities: Computational methods, bounds, and applications. In *Uncertainty in Artificial Intelligence*, Ed. R. Lopez de Mantaras and D. Poole, pp. 46–54. San Mateo, CA: Morgan Kaufmann.
- BOWDEN, R. J. & TURKINGTON, D. A. (1984). *Instrumental Variables*. Cambridge, MA: Cambridge University Press.
- COX, D. R. (1958). *The Planning of Experiments*. New York: John Wiley.
- COX, D. R. (1992). Causality: Some statistical aspects. *J. R. Statist. Soc. A* **155**, 291–301.
- COX, D. R. & WERMUTH, N. (1993). Linear dependencies represented by chain graphs. *Statist. Sci.* **8**, 204–18.
- DAWID, A. P. (1979). Conditional independence in statistical theory (with Discussion). *J. R. Statist. Soc. B* **41**, 1–31.
- FISHER, F. M. (1970). A correspondence principle for simultaneous equation models. *Econometrica* **38**, 73–92.
- FREEDMAN, D. (1987). As others see us: A case study in path analysis (with Discussion). *J. Educ. Statist.* **12**, 101–223.
- FRISCH, R. (1938). Statistical versus theoretical relations in economic macrodynamics. *League of Nations Memorandum*. Reproduced (1948) in *Autonomy of Economic Relations*, Universitetets Socialokonomiske Institut, Oslo.
- GALLES, D. & PEARL, J. (1995). Testing identifiability of causal effects. In *Uncertainty in Artificial Intelligence—11*, Ed. P. Besnard and S. Hanks, pp. 185–95. San Francisco, CA: Morgan Kaufmann.
- GEIGER, D., VERMA, T. S. & PEARL, J. (1990). Identifying independence in Bayesian networks. *Networks* **20**, 507–34.
- GOLDBERGER, A. S. (1972). Structural equation models in the social sciences. *Econometrica* **40**, 979–1001.
- HAAVELMO, T. (1943). The statistical implications of a system of simultaneous equations. *Econometrica* **11**, 1–12.
- HOLLAND, P. W. (1988). Causal inference, path analysis, and recursive structural equations models. In *Sociological Methodology*, Ed. C. Clogg, pp. 449–84. Washington, D.C.: American Sociological Association.
- IMBENS, G. W. & ANGRIST, J. D. (1994). Identification and estimation of local average treatment effects. *Econometrica* **62**, 467–76.
- LAURITZEN, S. L., DAWID, A. P., LARSEN, B. N. & LEIMER, H. G. (1990). Independence properties of directed Markov fields. *Networks* **20**, 491–505.
- LAURITZEN, S. L. & SPIEGELHALTER, D. J. (1988). Local computations with probabilities on graphical structures and their applications to expert systems (with Discussion). *J. R. Statist. Soc. B* **50**, 157–224.
- MANSKI, C. F. (1990). Nonparametric bounds on treatment effects. *Am. Econ. Rev., Papers Proc.* **80**, 319–23.
- PEARL, J. (1988). *Probabilistic Reasoning in Intelligent Systems*. San Mateo, CA: Morgan Kaufmann.
- PEARL, J. (1993a). Belief networks revisited. *Artif. Intel.* **59**, 49–56.
- PEARL, J. (1993b). Comment: Graphical models, causality, and intervention. *Statist. Sci.* **8**, 266–9.
- PEARL, J. (1994a). From Bayesian networks to causal networks. In *Bayesian Networks and Probabilistic Reasoning*, Ed. A. Gammerman, pp. 1–31. London: Alfred Walter.
- PEARL, J. (1994b). A probabilistic calculus of actions. In *Uncertainty in Artificial Intelligence*, Ed. R. Lopez de Mantaras and D. Poole, pp. 452–62. San Mateo, CA: Morgan Kaufmann.
- PEARL, J. (1995). Causal inference from indirect experiments. *Artif. Intel. Med. J.* To appear.
- PEARL, J. & VERMA, T. (1991). A theory of inferred causation. In *Principles of Knowledge Representation and Reasoning: Proceedings of the 2nd International Conference*, Ed. J. A. Allen, R. Fikes and E. Sandewall, pp. 441–52. San Mateo, CA: Morgan Kaufmann.

- PRAATT, J. W. & SCHLAIFER, R. (1988). On the interpretation and observation of laws. *J. Economet.* **39**, 23–52.
- ROBINS, J. M. (1986). A new approach to causal inference in mortality studies with a sustained exposure period—applications to control of the healthy workers survivor effect. *Math. Model.* **7**, 1393–512.
- ROBINS, J. M. (1989). The analysis of randomized and non-randomized AIDS treatment trials using a new approach to causal inference in longitudinal studies. In *Health Service Research Methodology: A Focus on AIDS*, Ed. L. Sechrest, H. Freeman and A. Mulley, pp. 113–59. Washington, D.C.: NCHSR, U.S. Public Health Service.
- ROBINS, J. M., BLEVINS, D., RITTER, G. & WULFSOHN, M. (1992). G-estimation of the effect of prophylaxis therapy for pneumocystis carinii pneumonia on the survival of AIDS patients. *Epidemiology* **3**, 319–36.
- ROSENBAUM, P. R. (1984). The consequences of adjustment for a concomitant variable that has been affected by the treatment. *J. R. Statist. Soc. A* **147**, 656–66.
- ROSENBAUM, P. & RUBIN, D. (1983). The central role of propensity score in observational studies for causal effects. *Biometrika* **70**, 41–55.
- RUBIN, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.* **66**, 688–701.
- RUBIN, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *Ann. Statist.* **7**, 34–58.
- RUBIN, D. B. (1990). Neyman (1923) and causal inference in experiments and observational studies. *Statist. Sci.* **5**, 472–80.
- SIMON, H. A. (1953). Causal ordering and identifiability. In *Studies in Econometric Method*, Ed. W. C. Hood and T. C. Hoopmans, Ch. 3. New York: John Wiley.
- SOBEL, M. E. (1990). Effect analysis and causation in linear structural equation models. *Psychometrika* **55**, 495–515.
- SPIEGELHALTER, D. J., LAURITZEN, S. L., DAWID, A. P. & COWELL, R. G. (1993). Bayesian analysis in expert systems (with Discussion). *Statist. Sci.* **8**, 219–47.
- SPIRITES, P. (1995). Conditional independence in directed cyclic graphical models for feedback. *Networks*. To appear.
- SPIRITES, P., GLYMOUR, C. & SCHEINES, R. (1993). *Causation, Prediction, and Search*. New York: Springer-Verlag.
- STROTZ, R. H. & WOLD, H. O. A. (1960). Recursive versus nonrecursive systems: An attempt at synthesis. *Econometrica* **28**, 417–27.
- WAINER, H. (1989). Eelworms, bullet holes, and Geraldine Ferraro: Some problems with statistical adjustment and some solutions. *J. Educ. Statist.* **14**, 121–40.
- WERMUTH, N. (1992). On block-recursive regression equations (with Discussion). *Brazilian J. Prob. Statist.* **6**, 1–56.
- WHITTAKER, J. (1990). *Graphical Models in Applied Multivariate Statistics*. Chichester: John Wiley.
- WRIGHT, S. (1921). Correlation and causation. *J. Agric. Res.* **20**, 557–85.

[Received May 1994. Revised February 1995]

Discussion of 'Causal diagrams for empirical research' by J. Pearl

BY D. R. COX

Nuffield College, Oxford, OX1 1NF, U.K.

AND NANNY WERMUTH

*Psychologisches Institut, Johannes Gutenberg-Universität Mainz, Staudingerweg 9,
D-55099 Mainz, Germany*

Judea Pearl has provided a general formulation for uncovering, under very explicit assumptions, what he calls the causal effect on y of 'setting' a variable x at a specified level, $\text{pr}(y|\bar{x})$, as assessed in a system of dependencies that can be represented by a directed acyclic graph. His Theorem 3 then provides a powerful computational scheme.

The back-door criterion requires there to be no unobserved 'common cause' for x and y that is not blocked out by observed variables, that is at least one of the intermediate variables between x and y or the common cause is to be observed. It is precisely doubt about such assumptions that makes epidemiologists, for example, wisely in our view, so cautious in distinguishing risk factors from causal effects. The front-door criterion requires, first, that there be an observed variable z such