

# Graphical Models for Probabilistic and Causal Reasoning

Judea Pearl

Cognitive Systems Laboratory

Computer Science Department

University of California, Los Angeles, CA 90024

(310) 825-3243

(310) 825-2273 Fax

*judea@cs.ucla.edu*

## 1 INTRODUCTION

This chapter surveys the development of graphical models known as Bayesian networks, summarizes their semantical basis and assesses their properties and applications to reasoning and planning.

Bayesian networks are directed acyclic graphs (DAGs) in which the nodes represent variables of interest (e.g., the temperature of a device, the gender of a patient, a feature of an object, the occurrence of an event) and the links represent causal influences among the variables. The strength of an influence is represented by conditional probabilities that are attached to each cluster of parents-child nodes in the network.

Figure 1 illustrates a simple yet typical Bayesian network. It describes the causal relationships among the season of the year ( $X_1$ ), whether rain falls ( $X_2$ ) during the season, whether the sprinkler is on ( $X_3$ ) during that season, whether the pavement would get wet ( $X_4$ ), and whether the pavement would be slippery ( $X_5$ ). All variables in this figure are binary, taking a value of either true or false, except the root variable  $X_1$  which can take one of four values: Spring, Summer, Fall, or Winter. Here, the absence of a direct link between  $X_1$  and  $X_5$ , for example, captures our understanding that the influence of seasonal variations on the slipperiness of the pavement is mediated by other conditions (e.g., the wetness of the pavement).

As this example illustrates, a Bayesian network constitutes a model of the environment rather than, as in many other knowledge representation schemes (e.g., logic, rule-based systems and neural networks), a model of the reasoning process. It simulates, in fact, the causal mechanisms that operate in the environment, and thus allows the investigator to answer a variety of queries, including: associational queries, such as “Having observed  $A$ , what can we expect of  $B$ ?”; abductive queries, such as “What is the most plausible explanation for a given set of observations?”; and control queries; such as “What will happen if we intervene

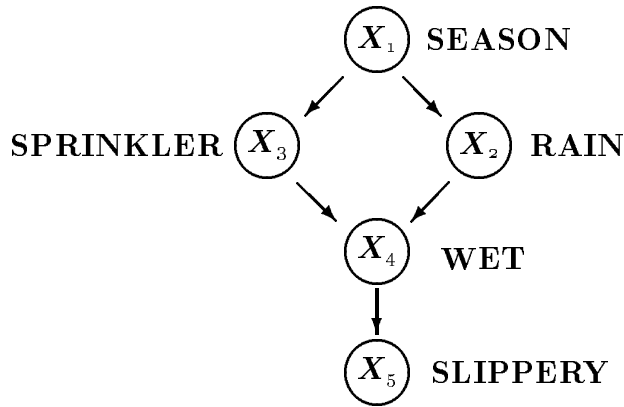


Figure 1: A Bayesian network representing causal influences among five variables.

and act on the environment?”. Answers to the first type of query depend only on probabilistic knowledge of the domain, while answers to the second and third types rely on the causal knowledge embedded in the network. Both types of knowledge, associative and causal, can effectively be represented and processed in Bayesian networks.

The associative facility of Bayesian networks may be used to model cognitive tasks such as object recognition, reading comprehension, and temporal projections. For such tasks, the probabilistic basis of Bayesian networks offers a coherent semantics for coordinating top-down and bottom-up inferences, thus bridging information from high-level concepts and low-level percepts. This capability is important for achieving selective attention, that is, selecting the most informative next observation before actually making the observation. In certain structures, the coordination of these two modes of inference can be accomplished by parallel and distributed processes that communicate through the links in the network.

However, the most distinctive feature of Bayesian networks, stemming largely from their causal organization, is their ability to represent and respond to changing configurations. Any local reconfiguration of the mechanisms in the environment can be translated, with only minor modification, into an isomorphic reconfiguration of the network topology. For example, to represent a disabled sprinkler, we simply delete from the network all links incident to the node “Sprinkler”. To represent a pavement covered by a tent, we simply delete the link between “Rain” and “Wet”. This flexibility is often cited as the ingredient that marks the division between deliberative and reactive agents, and that enables the former to manage novel situations instantaneously, without requiring retraining or adaptation.

## 2 HISTORICAL BACKGROUND

Networks employing directed acyclic graphs (DAGs) have a long and rich tradition, starting with the geneticist Sewall Wright (1921). He developed a method called *Path Analysis* [Wright, 1934], which later became an established representation of causal models in economics [Wold, 1964], sociology [Blalock, 1971, Kenny, 1979], and psychology [Duncan, 1975]. Good (1961) used DAGs to represent causal hierarchies of binary variables with disjunctive causes. *Influence diagrams* represent another application of DAG representation [Howard and Matheson, 1981]. Developed for decision analysis, they contain both event nodes and decision nodes. *Recursive models* is the name given to such networks by statisti-

cians seeking meaningful and effective decompositions of contingency tables [Lauritzen, 1982, Wermuth and Lauritzen, 1983, Kiiveri *et al.*, 1984].

The role of the network in the applications above was primarily to provide an efficient description for probability functions; once the network was configured, all subsequent computations were pursued by symbolic manipulation of probability expressions. The potential for the network to work as a computational architecture, and hence as a model of cognitive activities, was noted in [Pearl, 1982], where a distributed scheme was demonstrated for probabilistic updating on tree-structured networks. The motivation behind this particular development was the modeling of distributed processing in reading comprehension [Rumelhart, 1976], where both top-down and bottom-up inferences are combined to form a coherent interpretation. This dual mode of reasoning is at the heart of Bayesian updating, and in fact motivated Reverend Bayes's original 1763 calculations of posterior probabilities (representing explanations), given prior probabilities (representing causes), and likelihood functions (representing evidence).

Bayesian networks have not attracted much attention in the logic and cognitive modeling circles, but they did in expert systems. The ability to coordinate bi-directional inferences filled a void in expert systems technology of the late 1970s, and it is in this area that Bayesian networks truly flourished. Over the past ten years, Bayesian networks have become a tool of great versatility and power, and they are now the most common representation scheme for probabilistic knowledge [Shafer and Pearl, 1990, Shachter, 1990, Oliver and Smith, 1990, Neapolitan, 1990]. They have been used to aid in the diagnosis of medical patients [Heckerman, 1991, Andersen *et al.*, 1989, Heckerman *et al.*, 1990, Peng and Reggia, 1990] and malfunctioning systems [Agogino *et al.*, 1988], to understand stories [Charniak and Goldman, 1991], to filter documents [Turtle and Croft, 1991], to interpret pictures [Levitt *et al.*, 1990], to perform filtering, smoothing, and prediction [Abramson, 1991], to facilitate planning in uncertain environments [Dean and Wellman, 1991], and to study causation, nonmonotonicity, action, change, and attention. Some of these applications are described in a tutorial article by [Charniak, 1991]; others can be found in [Pearl, 1988], [Shafer and Pearl, 1990] and [Goldszmidt and Pearl, 1996].

## 3 BAYESIAN NETWORKS AS CARRIERS OF PROBABILISTIC INFORMATION

### 3.1 Formal Semantics

Given a DAG  $G$  and a joint distribution  $P$  over a set  $X = \{X_1, \dots, X_n\}$  of discrete variables, we say that  $G$  *represents*  $P$  if there is a one-to-one correspondence between the variables in  $X$  and the nodes of  $G$ , such that  $P$  admits the recursive product decomposition

$$P(x_1, \dots, x_n) = \prod_i P(x_i \mid \mathbf{pa}_i) \quad (1)$$

where  $\mathbf{pa}_i$  are the direct predecessors (called *parents*) of  $X_i$  in  $G$ . For example, the DAG in Figure 1 induces the decomposition

$$P(x_1, x_2, x_3, x_4, x_5) = P(x_1) P(x_2|x_1) P(x_3|x_1) P(x_4|x_2, x_3) P(x_5|x_4) \quad (2)$$

The recursive decomposition in Eq. (1) implies that, given its parent set  $\mathbf{pa}_i$ , each variable  $X_i$  is conditionally independent of all its other predecessors  $\{X_1, X_2, \dots, X_{i-1}\} \setminus \mathbf{pa}_i$ . Using Dawid's notation [Dawid, 1979], we can state this set of independencies as

$$X_i \perp\!\!\!\perp \{X_1, X_2, \dots, X_{i-1}\} \setminus \mathbf{pa}_i \mid \mathbf{pa}_i \quad i = 2, \dots, n \quad (3)$$

Such a set of independencies is called *Markovian*, since it reflects the Markovian condition for state transitions: each state is rendered independent of the past, given its immediately preceding state. For example, the DAG of Figure 1 implies the following Markovian independencies:

$$X_2 \perp\!\!\!\perp \{0\} \mid X_1, \quad X_3 \perp\!\!\!\perp X_2 \mid X_1, \quad X_4 \perp\!\!\!\perp X_1 \mid \{X_2, X_3\}, \quad X_5 \perp\!\!\!\perp \{X_1, X_2, X_3\} \mid X_4 \quad (4)$$

In addition to these, the decomposition of Eq. (1) implies many more independencies, the sum total of which can be identified from the DAG using the graphical criterion of *d*-separation [Pearl, 1988]:

**Definition 3.1** (*d*-separation) *Let a path in a DAG be a sequence of consecutive edges, of any directionality. A path  $p$  is said to be  $d$ -separated (or blocked) by a set of nodes  $Z$  iff:*

- (i)  *$p$  contains a chain  $i \rightarrow j \rightarrow k$  or a fork  $i \leftarrow j \rightarrow k$  such that the middle node  $j$  is in  $Z$ , or,*
- (ii)  *$p$  contains an inverted fork  $i \rightarrow j \leftarrow k$  such that neither the middle node  $j$  nor any of its descendants (in  $G$ ) are in  $Z$ .*

*If  $X, Y$ , and  $Z$  are three disjoint subsets of nodes in a DAG  $G$ , then  $Z$  is said to  $d$ -separate  $X$  from  $Y$ , denoted  $(X \perp\!\!\!\perp Y|Z)_G$ , iff  $Z$   $d$ -separates every path from a node in  $X$  to a node in  $Y$ .*

In Figure 1, for example,  $X = \{X_2\}$  and  $Y = \{X_3\}$  are  $d$ -separated by  $Z = \{X_1\}$ ; the path  $X_2 \leftarrow X_1 \rightarrow X_3$  is blocked by  $X_1 \in Z$ , while the path  $X_2 \rightarrow X_4 \leftarrow X_3$  is blocked because  $X_4$  and all its descendants are outside  $Z$ . Thus,  $(X_2 \perp\!\!\!\perp X_3|X_1)_G$  holds in Figure 1. However,  $X$  and  $Y$  are not  $d$ -separated by  $Z' = \{X_1, X_5\}$ , because the path  $X_2 \rightarrow X_4 \leftarrow X_3$  is rendered active by virtue of  $X_5$ , a descendant of  $X_4$ , being in  $Z'$ . Consequently,  $(X_2 \perp\!\!\!\perp X_3|\{X_1, X_5\})_G$  does not hold; in words, learning the value of the consequence  $X_5$  renders its causes  $X_2$  and  $X_3$  dependent, as if a pathway were opened along the arrows converging at  $X_4$ .

The  $d$ -separation criterion has been shown to be both necessary and sufficient relative to the set of distributions that are represented by a DAG  $G$  [Verma and Pearl, 1990, Geiger *et al.*, 1990]. In other words, there is a one-to-one correspondence between the set of independencies implied by the recursive decomposition of Eq. (1) and the set of triples

$(X, Z, Y)$  that satisfy the  $d$ -separation criterion in  $G$ . Furthermore, the  $d$ -separation criterion can be tested in time linear in the number of edges in  $G$ . Thus, a DAG can be viewed as an efficient scheme for representing Markovian independence assumptions and for deducing and displaying all the logical consequences of such assumptions.

An important property that follows from the  $d$ -separation characterization is a criterion for determining when two dags are observationally equivalent, that is, every probability distribution that is represented by one of the dags is also represented by the other:

**Theorem 3.2** [Verma and Pearl, 1990] *Two dags are observationally equivalent if and only if they have the same sets of edges and the same sets of  $v$ -structures, that is, head-to-head arrows with non-adjacent tails .*

The soundness of the  $d$ -separation criterion holds not only for probabilistic independencies but for any abstract notion of conditional independence that obeys the semi-graphoid axioms [Verma and Pearl, 1990, Geiger *et al.*, 1990]. Additional properties of DAGs and their applications to evidential reasoning in expert systems are discussed in [Pearl, 1988, Pearl *et al.*, 1990, Geiger, 1990, Lauritzen and Spiegelhalter, 1988, Spiegelhalter *et al.*, 1993, Pearl, 1993a].

## 3.2 Inference Algorithms

The first algorithms proposed for probability updating in Bayesian networks used message-passing architecture and were limited to trees [Pearl, 1982] and singly connected networks [Kim and Pearl, 1983]. The idea was to assign each variable a simple processor, forced to communicate only with its neighbors, and to permit asynchronous back-and-forth message-passing until equilibrium was achieved. Coherent equilibrium can indeed be achieved this way, but only in singly connected networks, where an equilibrium state occurs in time proportional to the diameter of the network.

Many techniques have been developed and refined to extend the tree-propagation method to general, multiply connected networks. Among the most popular are Shachter's (1988) method of node elimination, Lauritzen and Spiegelhalter's (1988) method of clique-tree propagation, and the method of loop-cut conditioning [Pearl, 1988, Chapter 4.3].

Clique-tree propagation, the most popular of the three methods, works as follows. Starting with a directed network representation, the network is transformed into an undirected graph that retains all of its original dependencies. This graph, sometimes called a Markov network [Pearl, 1988, Chapter 3.1], is then triangulated to form local clusters of nodes (cliques) that are tree-structured. Evidence propagates from clique to clique by ensuring that the probability of their intersection set is the same, regardless of which of the two cliques is considered in the computation. Finally, when the propagation process subsides, the posterior probability of an individual variable is computed by projecting (marginalizing) the distribution of the hosting clique onto this variable.

While the task of updating probabilities in general networks is NP-hard [Rosenthal, 1977, Cooper, 1990], the complexity for each of the three methods cited above is exponential in the size of the largest clique found in some triangulation of the network. It is fortunate that these complexities can be estimated prior to actual processing; when the estimates exceed reasonable bounds, an approximation method such as stochastic simulation

[Pearl, 1987, Henrion, 1988] can be used instead. Learning techniques have also been developed for systematic updating of the conditional probabilities  $P(x_i|\mathbf{pa}_i)$  so as to match empirical data [Spiegelhalter and Lauritzen, 1990].

### 3.3 System's properties

By providing graphical means for representing and manipulating probabilistic knowledge, Bayesian networks overcome many of the conceptual and computational difficulties of earlier knowledge-based systems [Pearl, 1988]. Their basic properties and capabilities can be summarized as follows:

1. Graphical methods make it easy to maintain consistency and completeness in probabilistic knowledge bases. They also define modular procedures of knowledge acquisition that reduce significantly the number of assessments required [Pearl, 1988, Heckerman, 1991].
2. Independencies can be dealt with explicitly. They can be articulated by an expert, encoded graphically, read off the network, and reasoned about; yet they forever remain robust to numerical imprecision [Geiger, 1990, Geiger *et al.*, 1990, Pearl *et al.*, 1990].
3. Graphical representations uncover opportunities for efficient computation. Distributed updating is feasible in knowledge structures which are rich enough to exhibit intercausal interactions (e.g., “explaining away”) [Pearl, 1982, Kim and Pearl, 1983] And, when extended by clustering or conditioning, tree-propagation algorithms are capable of updating networks of arbitrary topology [Lauritzen and Spiegelhalter, 1988, Shachter, 1986, Pearl, 1988].
4. The combination of predictive and abductive inferences resolves many problems encountered by first-generation expert systems and renders belief networks a viable model for cognitive functions requiring both top-down and bottom-up inferences [Pearl, 1988, Shafer and Pearl, 1990].
5. The causal information encoded in Bayesian networks facilitates the analysis of action sequences, their consequences, their interaction with observations, their expected utilities and, hence, the synthesis of plans and strategies under uncertainty [Dean and Wellman, 1991, Pearl, 1993b, Pearl, 1994b].
6. The isomorphism between the topology of Bayesian networks and the stable mechanisms which operate in the environment facilitates modular reconfiguration of the network in response to changing conditions, and permits deliberative reasoning about novel situations.

### 3.4 Recent Developments

#### 3.4.1 Causal discovery

One of the most exciting prospects in recent years has been the possibility of using the theory of Bayesian networks to discover causal structures in raw statistical data. Several

systems have been developed for this purpose [Pearl and Verma, 1991, Spirtes *et al.*, 1993], which systematically search and identify causal structures with hidden variables from empirical data. Technically, because these algorithms rely merely on conditional independence relationships, the structures found are valid only if one is willing to accept weaker forms of guarantees than those obtained through controlled randomized experiments: minimality and stability [Pearl and Verma, 1991]. Minimality guarantees that any other structure compatible with the data is necessarily less specific, and hence less falsifiable and less trustworthy, than the one(s) inferred. Stability ensures that any alternative structure compatible with the data must be less stable than the one(s) inferred; namely, slight fluctuations in experimental conditions will render that structure no longer compatible with the data. With these forms of guarantees, the theory provides criteria for identifying genuine and spurious causes, with or without temporal information.

Alternative methods of identifying structure in data assign prior probabilities to the parameters of the network and use Bayesian updating to score the degree to which a given network fits the data [Cooper and Herskovits, 1991, Heckerman *et al.*, 1994]. These methods have the advantage of operating well under small sample conditions, but encounter difficulties coping with hidden variables.

### 3.4.2 Plain beliefs

In mundane decision making, beliefs are revised not by adjusting numerical probabilities but by tentatively accepting some sentences as “true for all practical purposes”. Such sentences, often named *plain beliefs*, exhibit both logical and probabilistic character. As in classical logic, they are propositional and deductively closed; as in probability, they are subject to retraction and to varying degrees of entrenchment [Spohn, 1988, Goldszmidt and Pearl, 1992].

Bayesian networks can be adopted to model the dynamics of plain beliefs by replacing ordinary probabilities with non-standard probabilities, that is, probabilities that are infinitesimally close to either zero or one. This amounts to taking an “order of magnitude” approximation of empirical frequencies, and adopting new combination rules tailored to reflect this approximation. The result is an integer-addition calculus, very similar to probability calculus, with summation replacing multiplication and minimization replacing addition. A plain belief is then identified as a proposition whose negation obtains an infinitesimal probability (i.e., an integer greater than zero). The connection between infinitesimal probabilities and nonmonotonic logic is described in [Pearl, 1994a] and [Goldszmidt and Pearl, 1996].

This combination of infinitesimal probabilities with the causal information encoded by the structure of Bayesian networks facilitates linguistic communication of belief commitments, explanations, actions, goals, and preferences, and serves as the basis for current research on qualitative planning under uncertainty [Darwiche and Pearl, 1994, Goldszmidt and Pearl, 1992, Pearl, 1993b, Darwiche and Goldszmidt, 1994]. Some of these aspects will be presented in the next section.

## 4 BAYESIAN NETWORKS AS CARRIERS OF CAUSAL INFORMATION

The interpretation of DAGs as carriers of independence assumptions does not necessarily imply causation and will in fact be valid for any set of Markovian independencies along any ordering (not necessarily causal or chronological) of the variables. However, the patterns of independencies portrayed in a DAG are typical of causal organizations and some of these patterns can only be given meaningful interpretation in terms of causation. Consider, for example, two independent events,  $E_1$  and  $E_2$ , that have a common effect  $E_3$ . This triple represents an intransitive pattern of dependencies:  $E_1$  and  $E_3$  are dependent,  $E_3$  and  $E_2$  are dependent, yet  $E_1$  and  $E_2$  are independent. Such a pattern cannot be represented in undirected graphs because connectivity in undirected graphs is transitive. Likewise, it is not easily represented in neural networks, because  $E_1$  and  $E_2$  should turn dependent once  $E_3$  is known. The DAG representation provides a convenient language for intransitive dependencies via the converging pattern  $E_1 \rightarrow E_3 \leftarrow E_2$ , which implies the independence of  $E_1$  and  $E_2$  as well as the dependence of  $E_1$  and  $E_3$  and of  $E_2$  and  $E_3$ . The distinction between transitive and intransitive dependencies is the basis for the causal discovery systems of [Pearl and Verma, 1991] and [Spirtes *et al.*, 1993] (see Section 3.4.1).

However, the Markovian account still leaves open the question of how such intricate patterns of independencies relate to the more basic notions associated with causation, such as influence, manipulation, and control, which reside outside the province of probability theory. The connection is made in the mechanism-based account of causation.

The basic idea behind this account goes back to structural equations models [Wright, 1921, Haavelmo, 1943, Simon, 1953] and it was adapted in [Pearl and Verma, 1991] for defining probabilistic causal theories, as follows. Each child-parents family in a DAG  $G$  represents a deterministic function

$$X_i = f_i(\mathbf{pa}_i, \epsilon_i) \tag{5}$$

where  $\mathbf{pa}_i$  are the parents of variable  $X_i$  in  $G$ , and  $\epsilon_i$ ,  $0 < i < n$ , are mutually independent, arbitrarily distributed random disturbances. Characterizing each child-parent relationship as a deterministic function, instead of the usual conditional probability  $P(x_i | \mathbf{pa}_i)$ , imposes equivalent independence constraints on the resulting distributions and leads to the same recursive decomposition that characterizes DAG models (see Eq. (1)). However, the functional characterization  $X_i = f_i(\mathbf{pa}_i, \epsilon_i)$  also specifies how the resulting distributions would change in response to external interventions, since each function is presumed to represent a stable mechanism in the domain and therefore remains constant unless specifically altered. Thus, once we know the identity of the mechanisms altered by the intervention and the nature of the alteration, the overall effect of an intervention can be predicted by modifying the appropriate equations in the model of Eq. (5) and using the modified model to compute a new probability function of the observables.

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 *atomic* intervention amounts to replacing the old functional mechanism  $X_i = f_i(\mathbf{pa}_i, \epsilon_i)$  with a new mechanism  $X_i = x'_i$  governed by some external force that sets the value  $x'_i$ . If we imagine that each variable  $X_i$  could potentially be subject to the influence of such an external force, then we can view each Bayesian network as



an efficient code for predicting the effects of atomic interventions and of various combinations of such interventions, without representing these interventions explicitly.

## 4.1 Causal theories, actions, causal effect, and identifiability

**Definition 4.1** *A causal theory is a 4-tuple*

$$T = \langle V, U, P(\mathbf{u}), \{f_i\} \rangle$$

where

- (i)  $V = \{X_1, \dots, X_n\}$  is a set of observed variables
- (ii)  $U = \{U_1, \dots, U_m\}$  is a set of unobserved variables which represent disturbances, abnormalities or assumptions,
- (iii)  $P(\mathbf{u})$  is a distribution function over  $U_1, \dots, U_m$ , and
- (iv)  $\{f_i\}$  is a set of  $n$  deterministic functions, each of the form

$$X_i = f_i(PA_i, u) \quad i = 1, \dots, n \quad (6)$$

where  $PA_i$  is a subset of  $V$  not containing  $X_i$ .

The variables  $PA_i$  (connoting “parents”) are considered the direct causes of  $X_i$  and they define a directed graph  $G$  which may, in general, be cyclic. Unlike the probabilistic definition of “parents” in Bayesian networks (Eq. (1)),  $PA_i$  is selected from  $V$  by considering functional mechanisms in the domain, not by conditional independence considerations. We will assume that the set of equations in (6) has a unique solution for  $X_1, \dots, X_n$ , given any value of the disturbances  $U_1, \dots, U_m$ . Therefore the distribution  $P(\mathbf{u})$  induces a unique distribution on the observables, which we denote by  $P_T(\mathbf{v})$ .

We will consider concurrent actions of the form  $do(X = x)$ , where  $X \subseteq V$  is a set of variables and  $x$  is a set of values from the domain of  $X$ . In other words,  $do(X = x)$  represents a combination of actions that forces the variables in  $X$  to attain the values  $x$ .

**Definition 4.2** *(effect of actions) The effect of the action  $do(X = x)$  on a causal theory  $T$  is given by a subtheory  $T_x$  of  $T$ , where  $T_x$  obtains by deleting from  $T$  all equations corresponding to variables in  $X$  and substituting the equations  $X = x$  instead.*

The framework provided by Definitions 4.1 and 4.2 permits the coherent formalization of many subtle concepts in causal discourse, such as causal influence, causal effect, causal relevance, average causal effect, identifiability, counterfactuals, exogeneity, and so on. Examples are:

\*\*\*  **$X$  influences  $Y$**  in context  $u$  if there are two values of  $X$ ,  $x$  and  $x'$ , such that the solution for  $Y$  under  $U = u$  and  $do(X = x)$  is different from the solution under  $U = u$  and  $do(X = x')$ .

\*\*\*  $X$  **can potentially influence**  $Y$  if there exist both a subtheory  $T_z$  of  $T$  and a context  $U = u$  in which  $X$  influences  $Y$ .

\*\*\* Event  $X = x$  **is the (singular) cause of event**  $Y = y$  if (i)  $X = x$  and  $Y = y$  are true, and (ii) in every context  $u$  compatible with  $X = x$  and  $Y = y$ , and for all  $x' \neq x$ , the solution of  $Y$  under  $do(X = x')$  is not equal to  $y$ .

The definitions above are deterministic. Probabilistic causality emerges when we define a probability distribution  $P(u)$  for the  $U$  variables, which, under the assumption that the equations have a unique solution, induces a unique distribution on the endogenous variables for each combination of atomic interventions.

**Definition 4.3** (*causal effect*) Given two disjoint subsets of variables,  $X \subseteq V$  and  $Y \subseteq V$ , the causal effect of  $X$  on  $Y$ , denoted  $P_T(y|\hat{x})$ , is a function from the domain of  $X$  to the space of probability distributions on  $Y$ , such that

$$P_T(y|\hat{x}) = P_{T_x}(y) \tag{7}$$

for each realization  $x$  of  $X$ . In other words, for each  $x \in \text{dom}(X)$ , the causal effect  $P_T(y|\hat{x})$  gives the distribution of  $Y$  induced by the action  $do(X = x)$ .

Note that causal effects are defined relative to a given causal theory  $T$ , though the subscript  $T$  is often suppressed for brevity.

**Definition 4.4** (*identifiability*) Let  $Q(T)$  be any computable quantity of a theory  $T$ ;  $Q$  is identifiable in a class  $M$  of theories if for any pair of theories  $T_1$  and  $T_2$  from  $M$ ,  $Q(T_1) = Q(T_2)$  whenever  $P_{T_1}(v) = P_{T_2}(v)$ .

Identifiability is essential for estimating quantities  $Q$  from  $P$  alone, without specifying the details of  $T$ , so that the general characteristics of the class  $M$  suffice. The question of interest in planning applications is the identifiability of the causal effect  $Q = P_T(y|\hat{x})$  in the class  $M_G$  of theories that share the same causal graph  $G$ . Relative to such classes we now define:

**Definition 4.5** (*causal-effect identifiability*) The causal effect of  $X$  on  $Y$  is said to be identifiable in  $M_G$  if the quantity  $P(y|\hat{x})$  can be computed uniquely from the probabilities of the observed variables, that is, if for every pair of theories  $T_1$  and  $T_2$  in  $M_G$  such that  $P_{T_1}(v) = P_{T_2}(v)$ , we have  $P_{T_1}(y|\hat{x}) = P_{T_2}(y|\hat{x})$ .

The identifiability of  $P(y|\hat{x})$  ensures that it is possible to infer the effect of action  $do(X = x)$  on  $Y$  from two sources of information:

- (i) passive observations, as summarized by the probability function  $P(v)$ ,
- (ii) the causal graph,  $G$ , which specifies, qualitatively, which variables make up the stable mechanisms in the domain or, alternatively, which variables participate in the determination of each variable in the domain.

Simple examples of identifiable causal effects will be discussed in the next subsection.

## 4.2 Acting vs. Observing

Consider the example depicted in Figure 1. The corresponding theory consists of five functions, each representing an autonomous mechanism:

$$\begin{aligned}
 X_1 &= U_1 \\
 X_2 &= f_2(X_1, U_2) \\
 X_3 &= f_3(X_1, U_3) \\
 X_4 &= f_4(X_3, X_2, U_4) \\
 X_5 &= f_5(X_4, U_5)
 \end{aligned} \tag{8}$$

To represent the action “turning the sprinkler ON”,  $do(X_3 = \text{ON})$ , we delete the equation  $X_3 = f_3(x_1, u_3)$  from the theory of Eq. (8), and replace it with  $X_3 = \text{ON}$ . The resulting subtheory,  $T_{X_3=\text{ON}}$ , contains all the information needed for computing the effect of the actions on other variables. It is easy to see from this subtheory that the only variables affected by the action are  $X_4$  and  $X_5$ , that is, the descendant, of the manipulated variable  $X_3$ .

The probabilistic analysis of causal theories becomes particularly simple when two conditions are satisfied:

1. The theory is recursive, i.e., there exists an ordering of the variables  $V = \{X_1, \dots, X_n\}$  such that each  $X_i$  is a function of a subset  $PA_i$  of its predecessors

$$X_i = f_i(PA_i, U_i), \quad PA_i \subseteq \{X_1, \dots, X_{i-1}\} \tag{9}$$

2. The disturbances  $U_1, \dots, U_n$  are mutually independent, that is,

$$P(u) = \prod_i P(u_i) \tag{10}$$

These two conditions, also called Markovian, are the basis of the independencies embodied in Bayesian networks (Eq. 1) and they enable us to compute causal effects directly from the conditional probabilities  $P(x_i|\mathbf{pa}_i)$ , without specifying the functional form of the functions  $f_i$ , or the distributions  $P(u_i)$  of the disturbances. This is seen immediately from the following observations: The distribution induced by any Markovian theory  $T$  is given by the product in Eq. (1)

$$P_T(x_1, \dots, x_n) = \prod_i P(x_i|\mathbf{pa}_i) \tag{11}$$

where  $\mathbf{pa}_i$  are (values of) the parents of  $X_i$  in the diagram representing  $T$ . At the same time, the subtheory  $T_{x'_j}$ , representing the action  $do(X_j = x'_j)$  is also Markovian, hence it also induces a product-like distribution

$$P_{T_{x'_j}}(x_1, \dots, x_n) = \begin{cases} \prod_{i \neq j} P(x_i|\mathbf{pa}_i) = \frac{P(x_1, \dots, x_n)}{P(x_j|\mathbf{pa}_j)} & \text{if } x_j = x'_j \\ 0 & \text{if } x_j \neq x'_j \end{cases} \tag{12}$$

where the partial product reflects the surgical removal of the

$$X_j = f_j(\mathbf{pa}_j, U_j)$$

from the theory of equation (9) (see [Pearl, 1993a]).

In the example of Figure 1, the pre-action distribution is given by the product

$$P_T(x_1, x_2, x_3, x_4, x_5) = P(x_1)P(x_2|x_1)P(x_3|x_1)P(x_4|x_2, x_3)P(x_5|x_4) \quad (13)$$

while the surgery corresponding to the action  $do(X_3 = \text{ON})$  amounts to deleting the link  $X_1 \rightarrow X_3$  from the graph and fixing the value of  $X_3$  to ON, yielding the post-action distribution:

$$P_T(x_1, x_2, x_4, x_5 | do(X_3 = \text{ON})) = P(x_1) P(x_2|x_1) P(x_4|x_2, X_3 = \text{ON}) P(x_5|x_4) \quad (14)$$

Note the difference between the action  $do(X_3 = \text{ON})$  and the observation  $X_3 = \text{ON}$ . The latter is encoded by ordinary Bayesian conditioning, while the former by conditioning a mutilated graph, with the link  $X_1 \rightarrow X_3$  removed. This mirrors indeed the difference between seeing and doing: after observing that the sprinkler is ON, we wish to infer that the season is dry, that it probably did not rain, and so on; no such inferences should be drawn in evaluating the effects of the deliberate action “turning the sprinkler ON”. The amputation of  $X_3 = f_3(X_1, U_3)$  from (8) ensures the suppression of any abductive inferences from  $X_3$ , the action’s recipient.

Note also that Equations (11) through (14) are independent of  $T$ , in other words, the pre-actions and post-action distributions depend only on observed conditional probabilities but are independent of the particular functional form of  $\{f_i\}$  or the distribution  $P(\mathbf{u})$  which generate those probabilities. This is the essence of identifiability as given in Definition 4.5, which stems from the Markovian assumptions (9) and (10). Section 4.3 will demonstrate that certain causal effects, though not all, are identifiable even when the Markovian property is destroyed by introducing dependencies among the disturbance terms.

Generalization to multiple actions and conditional actions are reported in [Pearl and Robins, 1995]. Multiple actions  $do(X = x)$ , where  $X$  is a compound variable result in a distribution similar to (12), except that all factors corresponding to the variables in  $X$  are removed from the product in (11). Stochastic conditional strategies [Pearl, 1994b] of the form

$$do(X_j = x_j) \text{ with probability } P^*(x_j | \mathbf{pa}_j^*) \quad (15)$$

where  $\mathbf{pa}_j^*$  is the support of the decision strategy, also result in a product decomposition similar to (11), except that each factor  $P(x_j | \mathbf{pa}_j)$  is *replaced* with  $P^*(x_j | \mathbf{pa}_j^*)$ .

The surgical procedure described above is not limited to probabilistic analysis. The causal knowledge represented in Figure 1 can be captured by logical theories as well, for example,

$$\begin{aligned} x_2 &\iff [(X_1 = \text{Winter}) \vee (X_1 = \text{Fall}) \vee ab_2] \wedge \neg ab'_2 \\ x_3 &\iff [(X_1 = \text{Summer}) \vee (X_1 = \text{Spring}) \vee ab_3] \wedge \neg ab'_3 \\ x_4 &\iff (x_2 \vee x_3 \vee ab_4) \wedge \neg ab'_4 \\ x_5 &\iff (x_4 \vee ab_5) \wedge \neg ab'_5 \end{aligned} \quad (16)$$

where  $x_i$  stands for  $X_i = \text{true}$ , and  $ab_i$  and  $ab'_i$  stand, respectively, for triggering and inhibiting abnormalities. The double arrows represent the assumption that the events on the r.h.s.

of each equation are the *only* direct causes for the l.h.s, thus identifying the surgery implied by any action.

It should be emphasized though that the models of a causal theory are not made up merely of truth value assignments which satisfy the equations in the theory. Since each equation represents an autonomous process, the content of each individual equation must be specified in any model of the theory, and this can be encoded using either the graph (as in Figure 1) or the generic description of the theory, as in (8). Alternatively, we can view a model of a causal theory to consist of a mutually consistent set of submodels, with each submodel being a standard model of a single equation in the theory.

### 4.3 Action Calculus

The identifiability of causal effects demonstrated in Section 4.1 relies critically on the Markovian assumptions (9) and (10). If a variable that has two descendants in the graph is unobserved, the disturbances in the two equations are no longer independent, the Markovian property (9) is violated and identifiability may be destroyed. This can be seen easily from Eq. (12); if any parent of the manipulated variable  $X_j$  is unobserved, one cannot estimate the conditional probability  $P(x_j|\mathbf{pa}_j)$ , and the effect of the action  $do(X_j = x_j)$  may not be predictable from the observed distribution  $P(x_1, \dots, x_n)$ . Fortunately, certain causal effects are identifiable even in situations where members of  $\mathbf{pa}_j$  are unobservable [Pearl, 1993a] and, moreover, polynomial tests are now available for deciding when  $P(x_i|\hat{x}_j)$  is identifiable, and for deriving closed-form expressions for  $P(x_i|\hat{x}_j)$  in terms of observed quantities [Galles and Pearl, 1995].

These tests and derivations are based on a symbolic calculus [Pearl, 1994b, 1995] to be described in the sequel, in which interventions, side by side with observations, are given explicit notation, and are permitted to transform probability expressions. The transformation rules of this calculus reflect the understanding that interventions perform “local surgeries” as described in Definition 4.2, i.e., they overrule equations that tie the manipulated variables to their pre-intervention causes.

Let  $X, Y$ , and  $Z$  be arbitrary disjoint sets of nodes in a DAG  $G$ . We denote by  $G_{\overline{X}}$  the graph obtained by deleting from  $G$  all arrows pointing to nodes in  $X$ . Likewise, we denote by  $G_{\underline{X}}$  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 use the notation  $G_{\overline{X}\underline{Z}}$ . Finally, the expression  $P(y|\hat{x}, z) \triangleq P(y, z|\hat{x})/P(z|\hat{x})$  stands for the probability of  $Y = y$  given that  $Z = z$  is observed and  $X$  is held constant at  $x$ .

**Theorem 4.6** Let  $G$  be the directed acyclic graph associated with a Markovian causal theory, and let  $P(\cdot)$  stand for the probability distribution induced by that theory. For any disjoint subsets of variables  $X, Y, Z$ , and  $W$  we have:

**Rule 1** (Insertion/deletion of observations):

$$P(y|\hat{x}, z, w) = P(y|\hat{x}, w) \quad \text{if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{X}}} \quad (17)$$

**Rule 2** (Action/observation exchange):

$$P(y|\hat{x}, \hat{z}, w) = P(y|\hat{x}, z, w) \text{ if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{XZ}}} \quad (18)$$

**Rule 3** (Insertion/deletion of actions):

$$P(y|\hat{x}, \hat{z}, w) = P(y|\hat{x}, w) \text{ if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{X}, \overline{Z(W)}}} \quad (19)$$

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

Each of the inference rules above follows from the basic interpretation of the “ $\hat{x}$ ” operator as a replacement of the causal mechanism that connects  $X$  to its pre-action parents by a new mechanism  $X = x$  introduced by the intervening force. The result is a submodel characterized by the subgraph  $G_{\overline{X}}$  (named “manipulated graph” in Spirtes et al. (1993)) which supports all three rules.

**Corollary 4.7** A causal effect  $q: P(y_1, \dots, y_k|\hat{x}_1, \dots, \hat{x}_m)$  is identifiable in a model characterized by a graph  $G$  if there exists a finite sequence of transformations, each conforming to one of the inference rules in Theorem 4.6, which reduces  $q$  into a standard (i.e., hat-free) probability expression involving observed quantities.  $\square$

Although Theorem 4.6 and Corollary 4.7 require the Markovian property, they can also be applied to non Markovian, recursive theories because such theories become Markovian if we consider the unobserved variables as part of the analysis, and represent them as nodes in the graph. To illustrate, assume that variable  $X_1$  in Figure 1 is unobserved, rendering the disturbances  $U_3$  and  $U_2$  dependent since these terms now include the common influence of  $X_1$ . Theorem 4.6 tells us that the causal effect  $P(x_4|\hat{x}_3)$  is identifiable, because:

$$P(x_4|\hat{x}_3) = \sum_{x_2} P(x_4|\hat{x}_3, x_2)P(x_2|\hat{x}_3)$$

Rule 3 permits the deletion

$$P(x_2|\hat{x}_3) = P(x_2), \text{ because } (X_2 \perp\!\!\!\perp X_3)_{G_{\overline{X_3}}},$$

while Rule 2 permits the exchange

$$P(x_4|\hat{x}_3, x_2) = P(x_4|x_3, x_2), \text{ because } (X_4 \perp\!\!\!\perp X_3|X_2)_{G_{\overline{X_3}}}.$$

This gives

$$P(x_4|\hat{x}_3) = \sum_{x_2} P(x_4|x_3, x_2)P(x_2)$$

which is a “hat-free” expression, involving only observed quantities.

In general, it can be shown (Pearl, 1995) that:

1. The effect of interventions can often be identified (from nonexperimental data) without resorting to parametric models,

2. The conditions under which such nonparametric identification is possible can be determined by simple graphical tests<sup>1</sup>, and,
3. When the effect of interventions is not identifiable, the causal graph may suggest non-trivial experiments which, if performed, would render the effect identifiable.

The ability to assess the effect of interventions from nonexperimental data has immediate applications in the medical and social sciences, since subjects who undergo certain treatments often are not representative of the population as a whole. Such assessments are also important in AI applications where an agent needs to predict the effect of the next action on the basis of past performance records, and where that action has never been enacted out of free will, but in response to environmental needs or to other agent’s requests.

#### 4.4 Historical Remarks

An explicit translation of interventions to “wiping out” equations from linear econometric models was first proposed by Strotz & Wold (1960) and later used in Fisher (1970) and Sobel (1990). Extensions to action representation in nonmonotonic systems were reported in [Goldszmidt and Pearl, 1992, Pearl, 1993a]. Graphical ramifications of this translation were explicated first in Spirtes et al. (1993) and later in Pearl (1993b). A related formulation of causal effects, based on event trees and counterfactual analysis was developed by Robins (1986, pp. 1422-25). Calculi for actions and counterfactuals based on this interpretation are developed in [Pearl, 1994b] and [Balke and Pearl, 1994], respectively.

## 5 Counterfactuals

A counterfactual sentence has the form

*If A were true, then C would have been true?*

where  $A$ , the counterfactual antecedent, specifies an event that is contrary to one’s real-world observations, and  $C$ , the counterfactual consequent, specifies a result that is expected to hold in the alternative world where the antecedent is true. A typical example is “If Oswald were not to have shot Kennedy, then Kennedy would still be alive” which presumes the factual knowledge of Oswald’s shooting Kennedy, contrary to the antecedent of the sentence.

The majority of the philosophers who have examined the semantics of counterfactual sentences have resorted to some version of Lewis’ “closest world” approach; “ $C$  if it were  $A$ ” is true, if  $C$  is true in worlds that are “closest” to the real world yet consistent with the counterfactuals antecedent  $A$  [Lewis, 1973]. Ginsberg (1986), followed a similar strategy. While the “closest world” approach leaves the precise specification of the closeness measure

---

<sup>1</sup>These graphical tests offer, in fact, a complete formal solution to the “covariate-selection” problem in statistics: finding an appropriate set of variables that need be adjusted for in any study which aims to determine the effect of one factor upon another. This problem has been lingering in the statistical literature since Karl Pearson, the founder of modern statistics, discovered (1899) what in modern terms is called the “Simpson’s paradox”; any statistical association between two variables may be reversed or negated by including additional factors in the analysis [Aldrich, 1995].

almost unconstrained, causal knowledge imposes very specific preferences as to which worlds should be considered closest to any given world. For example, considering an array of domino tiles standing close to each other. The manifestly closest world consistent with the antecedent “tile  $i$  is tipped to the right” would be a world in which just tile  $i$  is tipped, while all the others remain erect. Yet, we all accept the counterfactual sentence “Had tile  $i$  been tipped over to the right, tile  $i + 1$  would be tipped as well” as plausible and valid. Thus, distances among worlds are not determined merely by surface similarities but require a distinction between disturbed mechanisms and naturally occurring transitions. The local surgery paradigm expounded in Section 4.1 offers a concrete explication of the closest world approach which respects causal considerations. A world  $w_1$  is “closer” to  $w$  than a world  $w_2$  is, if the the set of atomic surgeries needed for transforming  $w$  into  $w_1$  is a subset of those needed for transforming  $w$  into  $w_2$ . In the domino example, finding tile  $i$  tipped and  $i + 1$  erect requires the breakdown of two mechanism (e.g., by two external actions) compared with one mechanism for the world in which all  $j$ -tiles,  $j > i$  are tipped. This paradigm conforms to our perception of causal influences and lends itself to economical machine representation.

## 5.1 Formal underpinning

The structural equation framework offers an ideal setting for counterfactual analysis.

**Definition 5.1** (*context-based potential response*): *Given a causal theory  $T$  and two disjoint sets of variables,  $X$  and  $Y$ , the potential response of  $Y$  to  $X$  in a context  $u$ , denoted  $Y(x, u)$  or  $Y_x(u)$ , is the solution for  $Y$  under  $U = u$  in the subtheory  $T_x$ .  $Y(x, u)$  can be taken as the formal definition of the counterfactual English phrase: “the value that  $Y$  would take in context  $u$ , had  $X$  been  $x$ ,”<sup>2</sup>*

Note that this definition allows for the context  $U = u$  and the proposition  $X = x$  to be incompatible in  $T$ . For example, if  $T$  describes a logic circuit with input  $U$  it may well be reasonable to assert the counterfactual: “Given  $U = u$ ,  $Y$  would be high if  $X$  were low”, even though the input  $U = u$  may preclude  $X$  from being low. It is for this reason that one must invoke some notion of intervention (alternatively, a theory change or a “miracle” [Lewis, 1973]) in the definition of counterfactuals.

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  $Y(x)$  or  $Y_x$ . Moreover, the distribution of this random variable is easily seen to coincide with the causal effect  $P(y|\hat{x})$ , as defined in Eq. (7), i.e.,

$$P((Y(x) = y) = P(y|\hat{x})$$

The probability of a counterfactual conditional  $x \rightarrow y \mid o$  may then be evaluated by the following procedure:

---

<sup>2</sup>The term *unit* instead of *context* is often used in the statistical literature [Rubin, 1974], where it normally stands for the identity of a specific individual in a population, namely, the set of attributes  $u$  that characterize that individual. In general,  $u$  may include the time of day, the experimental conditions under study, and so on. Practitioners of the counterfactual notation do not explicitly mention the notions of “solution” or “intervention” in the definition of  $Y(x, u)$ . Instead, the phrase “the value that  $Y$  would take in unit  $u$ , had  $X$  been  $x$ ,” viewed as basic, is posited as the definition of  $Y(x, u)$ .



- Use the observations  $o$  to update  $P(u)$  thus forming a causal theory  $T^o = \langle V, U, \{f_i\}, P(u|o) \rangle$
- Form the mutilated theory  $T_x^o$  (by deleting the equation corresponding to variables in  $X$ ) and compute the probability  $P_{T^o}(y|\hat{x})$  which  $T_x^o$  induces on  $Y$ .

Unlike causal effect queries, counterfactual queries are not identifiable even in Markovian theories, but require that the functional-form of  $\{f_i\}$  be specified. In [Balke and Pearl, 1994] a method is devised for computing sharp bounds on counterfactual probabilities which, under certain circumstances may collapse to point estimates. This method has been applied to the evaluation of causal effects in studies involving noncompliance, and to the determination of legal liability.

## 5.2 Applications to Policy Analysis

Counterfactual reasoning is at the heart of every planning activity, especially real-time planning. When a planner discovers that the current state of affairs deviates from the one expected, a “plan repair” activity need be invoked to determine what went wrong and how it could be rectified. This activity amounts to an exercise of counterfactual thinking, as it calls for rolling back the natural course of events and determining, based on the factual observations at hand, whether the culprit lies in previous decisions or in some unexpected, external eventualities. Moreover, in reasoning forward to determine if things would have been different a new model of the world must be consulted, one that embodies hypothetical changes in decisions or eventualities, hence, a breakdown of the old model or theory.

The logic-based planning tools used in AI, such as STRIPS and its variants or those based on the situation calculus, do not readily lend themselves to counterfactual analysis; as they are not geared for coherent integration of abduction with prediction, and they do not readily handle theory changes. Remarkably, the formal system developed in economics and social sciences under the rubric “structural equations models” does offer such capabilities but, as will be discussed below, these capabilities are not well recognized by current practitioners of structural models. The analysis presented in this chapter could serve both to illustrate to AI researchers the basic formal features needed for counterfactual and policy analysis, and to call the attention of economists and social scientists to capabilities that are dormant within structural equations models.

Counterfactual thinking dominates reasoning in political science and economics. We say, for example, “If Germany were not punished so severely at the end of World War I, Hitler would not have come to power,” or “If Reagan did not lower taxes, our deficit would be lower today.” Such thought experiments emphasize an understanding of generic laws in the domain and are aimed toward shaping future policy making, for example, “defeated countries should not be humiliated,” or “lowering taxes (contrary to Reaganomics) tends to increase national debt.”

Strangely, there is very little formal work on counterfactual reasoning or policy analysis in the behavioral science literature. An examination of a number of econometric journals and textbooks, for example, reveals a glaring imbalance: while an enormous mathematical machinery is brought to bear on problems of estimation and prediction, policy analysis (which is the ultimate goal of economic theories) receives almost no formal treatment. Currently, the most popular methods driving economic policy making are based on so-called *reduced-form*

analysis: to find the impact of a policy involving decision variables  $X$  on outcome variables  $Y$ , one examines past data and estimates the conditional expectation  $E(Y|X=x)$ , where  $x$  is the particular instantiation of  $X$  under the policy studied.

The assumption underlying this method is that the data were generated under circumstances in which the decision variables  $X$  act as exogenous variables, that is, variables whose values are determined outside the system under analysis. However, while new decisions should indeed be considered exogenous for the purpose of evaluation, past decisions are rarely enacted in an exogenous manner. Almost every realistic policy (e.g., taxation) imposes control over some endogenous variables, that is, variables whose values are determined by other variables in the analysis. Let us take taxation policies as an example. Economic data are generated in a world in which the government is reacting to various indicators and various pressures; hence, taxation is endogenous in the data-analysis phase of the study. Taxation becomes exogenous when we wish to predict the impact of a specific decision to raise or lower taxes. The reduced-form method is valid only when past decisions are non-responsive to other variables in the system, and this, unfortunately, eliminates most of the interesting control variables (e.g., tax rates, interest rates, quotas) from the analysis.

This difficulty is not unique to economic or social policy making; it appears whenever one wishes to evaluate the merit of a plan on the basis of the past performance of other agents. Even when the signals triggering the past actions of those agents are known with certainty, a systematic method must be devised for selectively ignoring the influence of those signals from the evaluation process. In fact, the very essence of *evaluation* is having the freedom to imagine and compare trajectories in various counterfactual worlds, where each world or trajectory is created by a hypothetical implementation of a policy that is free of the very pressures that compelled the implementation of such policies in the past.

Balke and Pearl (1995) demonstrate how linear, nonrecursive structural models with Gaussian noise can be used to compute counterfactual queries of the type: “Given an observation set  $O$ , find the probability that  $Y$  would have attained a value greater than  $y$ , had  $X$  been set to  $x$ ”. The task of inferring “causes of effects”, that is, of finding the probability that  $X = x$  is the cause for effect  $E$ , amounts to answering the counterfactual query: “Given effect  $E$  and observations  $O$ , find the probability that  $E$  would not have been realized, had  $X$  not been  $x$ ”. The technique developed in Balke and Pearl (1995) is based on probability propagation in dual networks, one representing the actual world, the other the counterfactual world. The method is not limited to linear functions but applies whenever we are willing to assume the functional form of the structural equations. The noisy OR-gate model [Pearl, 1988] is a canonical example where such functional form is normally specified. Likewise, causal theories based on Boolean functions (with exceptions), such as the one described in Eq. (16) lend themselves to counterfactual analysis in the framework of Definition 5.1.

## Acknowledgments

The research was partially supported by Air Force grant #F49620-94-1-0173, NSF grant #IRI-9420306, and Northrop/Rockwell Micro grant #94-100.

## References

- [Abramson, 1991] B. Abramson. ARCO1: An application of belief networks to the oil market. In *Proceedings of the Seventh Conference on Uncertainty in Artificial Intelligence*, San Mateo, CA, 1991. Morgan Kaufmann.
- [Agogino *et al.*, 1988] A.M. Agogino, S. Srinivas, and K. Schneider. Multiple sensor expert system for diagnostic reasoning, monitoring and control of mechanical systems. *Mechanical Systems and Signal Processing*, 2:165–185, 1988.
- [Aldrich, 1995] J. Aldrich. Correlations genuine and spurious in Pearson and Yule. Forthcoming *Statistical Science*, 1995.
- [Andersen *et al.*, 1989] S.K. Andersen, K.G. Olesen, F.V. Jensen, and F. Jensen. Hugin - a shell for building Bayesian belief universes for expert systems. In *Eleventh International Joint Conference on Artificial Intelligence*, pages 1080–1085, 1989.
- [Balke and Pearl, 1994] 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, 1995] A. Balke and J. Pearl. Counterfactuals and policy analysis in structural models. In P. Besnard and S. Hanks, editors, *Uncertainty in Artificial Intelligence 11*, pages 11–18. Morgan Kaufmann, San Francisco, CA, 1995.
- [Blalock, 1971] H.M. Blalock. *Causal Models in the Social Sciences*. Macmillan, London, 1971.
- [Charniak and Goldman, 1991] E. Charniak and R. Goldman. A probabilistic model of plan recognition. In *Proceedings, AAAI-91*. AAAI Press/The MIT Press, Anaheim, CA, 1991.
- [Charniak, 1991] E. Charniak. Bayesian networks without tears. *AI Magazine*, 12(4):50–63, 1991.
- [Cooper and Herskovits, 1991] G.F. Cooper and E. Herskovits. A Bayesian method for constructing Bayesian belief networks from databases. In B.D. D’Ambrosio, P. Smets, and P.P. Bonissone, editors, *Proceedings of Uncertainty in Artificial Intelligence Conference, 1991*, pages 86–94. Morgan Kaufmann, San Mateo, 1991.
- [Cooper, 1990] G.F. Cooper. Computational complexity of probabilistic inference using Bayesian belief networks. *Artificial Intelligence*, 42(2):393–405, 1990. research note.
- [Darwiche and Goldszmidt, 1994] A. Darwiche and M. Goldszmidt. On the relation between kappa calculus and probabilistic reasoning. In R. Lopez de Mantaras and D. Poole, editors, *Uncertainty in Artificial Intelligence*, volume 10, pages 145–153. Morgan Kaufmann, San Francisco, CA, 1994.

- [Darwiche and Pearl, 1994] A. Darwiche and J. Pearl. Symbolic causal networks for planning under uncertainty. In *Symposium Notes of the 1994 AAAI Spring Symposium on Decision-Theoretic Planning*, pages 41–47. Stanford, CA, 1994.
- [Dawid, 1979] A.P. Dawid. Conditional independence in statistical theory. *Journal of the Royal Statistical Society, Series A*, 41:1–31, 1979.
- [Dean and Wellman, 1991] T.L. Dean and M.P. Wellman. *Planning and Control*. Morgan Kaufmann, San Mateo, CA, 1991.
- [Duncan, 1975] O.D. Duncan. *Introduction to Structural Equation Models*. Academic Press, New York, 1975.
- [Fisher, 1970] F.M. Fisher. A correspondence principle for simultaneous equations models. *Econometrica*, 38:73–92, 1970.
- [Galles and Pearl, 1995] D. Galles and J. Pearl. Testing identifiability of causal effects. In P. Besnard and S. Hanks, editors, *Uncertainty in Artificial Intelligence 11*, pages 185–195. Morgan Kaufmann, San Francisco, CA, 1995.
- [Geiger *et al.*, 1990] D. Geiger, T.S. Verma, and J. Pearl. Identifying independence in Bayesian networks. In *Networks*, volume 20, pages 507–534. John Wiley and Sons, Sussex, England, 1990.
- [Geiger, 1990] D. Geiger. Graphoids: A qualitative framework for probabilistic inference. PhD thesis, University of California, Los Angeles, Department of Computer Science, 1990.
- [Ginsberg, 1986] M.L. Ginsberg. Counterfactuals. *Artificial Intelligence*, 30(35–79), 1986.
- [Goldszmidt and Pearl, 1992] M. Goldszmidt and J. Pearl. Rank-based systems: A simple approach to belief revision, belief update, and reasoning about evidence and actions. In B. Nebel, C. Rich, and W. Swartout, editors, *Proceedings of the Third International Conference on Knowledge Representation and Reasoning*, pages 661–672. Morgan Kaufmann, 1992.
- [Goldszmidt and Pearl, 1996] M. Goldszmidt and J. Pearl. Qualitative probabilities for default reasoning, belief revision, and causal modeling. *Artificial Intelligence*, 84(1–2):57–112, July 1996.
- [Good, 1961] I.J. Good. A causal calculus, I-II. *British Journal for the Philosophy of Science*, 11:305–318, 12:43–51, 1961.
- [Haavelmo, 1943] T. Haavelmo. The statistical implications of a system of simultaneous equations. *Econometrica*, 11:1–12, 1943.
- [Heckerman *et al.*, 1990] D.E. Heckerman, E.J. Horvitz, and B.N. Nathwany. Toward normative expert systems: The pathfinder project. Technical Report KSL-90-08, Medical Computer Science Group, Section on Medical Informatics, Stanford University, Stanford, CA, 1990.

- [Heckerman *et al.*, 1994] D. Heckerman, D. Geiger, and D. Chickering. Learning Bayesian networks: The combination of knowledge and statistical data. In *Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence*, Seattle, WA, pages 293–301, San Mateo, CA, July 1994. Morgan Kaufmann.
- [Heckerman, 1991] D. Heckerman. Probabilistic similarity networks. *Networks*, 20(5):607–636, 1991.
- [Henrion, 1988] M. Henrion. Propagation of uncertainty by probabilistic logic sampling in bayes’ networks. In J.F. Lemmer and L.N. Kanal, editors, *Uncertainty in Artificial Intelligence 2*, pages 149–164. Elsevier Science Publishers, North-Holland, Amsterdam, Netherlands, 1988.
- [Howard and Matheson, 1981] R.A. Howard and J.E. Matheson. Influence diagrams. *Principles and Applications of Decision Analysis*, 1981. Strategic Decisions Group, Menlo Park, CA.
- [Kenny, 1979] D.A. Kenny. *Correlation and Causality*. Wiley, New York, 1979.
- [Kiiveri *et al.*, 1984] H. Kiiveri, T.P. Speed, and J.B. Carlin. Recursive causal models. *Journal of Australian Math Society*, 36:30–52, 1984.
- [Kim and Pearl, 1983] J.H. Kim and J. Pearl. A computational model for combined causal and diagnostic reasoning in inference systems. In *Proceedings IJCAI-83*, pages 190–193, Karlsruhe, Germany, 1983.
- [Lauritzen and Spiegelhalter, 1988] S.L. Lauritzen and D.J. Spiegelhalter. Local computations with probabilities on graphical structures and their application to expert systems(with discussion). *Journal of the Royal Statistical Society, Series B*, 50(2):157–224, 1988.
- [Lauritzen, 1982] S.L. Lauritzen. *Lectures on Contingency Tables*. University of Aalborg Press, Aalborg, Denmark, 2nd ed. edition, 1982.
- [Levitt *et al.*, 1990] T.S. Levitt, J.M. Agosta, and T.O. Binford. Model-based influence diagrams for machine vision. In M. Henrion, R.D. Shachter, L.N. Kanal, and J.F. Lemmer, editors, *Uncertainty in Artificial Intelligence 5*, pages 371–388. North Holland, Amsterdam, 1990.
- [Lewis, 1973] D. Lewis. *Counterfactuals*. Basil Blackwell, Oxford, 1973.
- [Neapolitan, 1990] R.E. Neapolitan. *Probabilistic Reasoning in Expert Systems: Theory and Algorithms*. Wiley, New York, 1990.
- [Oliver and Smith, 1990] R.M. Oliver and J.Q. (Eds.) Smith. *Influence Diagrams, Belief Nets, and Decision Analysis*. John Wiley, New York, 1990.

- [Pearl and Robins, 1995] J. Pearl and J.M. Robins. Probabilistic evaluation of sequential plans from causal models with hidden variables. In P. Besnard and S. Hanks, editors, *Uncertainty in Artificial Intelligence 11*, pages 444–453. Morgan Kaufmann, San Francisco, CA, 1995.
- [Pearl and Verma, 1991] J. Pearl and T. Verma. A theory of inferred causation. In J.A. Allen, R. Fikes, and E. Sandewall, editors, *Principles of Knowledge Representation and Reasoning: Proceedings of the Second International Conference*, pages 441–452, San Mateo, CA, 1991. Morgan Kaufmann.
- [Pearl *et al.*, 1990] J. Pearl, D. Geiger, and T. Verma. The logic and influence diagrams. In R.M. Oliver and J.Q. Smith, editors, *Influence Diagrams, Belief Nets and Decision Analysis*, pages 67–87. Wiley, 1990.
- [Pearl, 1982] J. Pearl. Reverend bayes on inference engines: A distributed hierarchical approach. In *Proceedings AAAI National Conference on AI*, pages 133–136, Pittsburgh, PA, 1982.
- [Pearl, 1987] J. Pearl. Bayes decision methods. In *Encyclopedia of AI*, pages 48–56. Wiley Interscience, New York, 1987.
- [Pearl, 1988] J. Pearl. *Probabilistic Reasoning in Intelligence Systems*. Morgan Kaufmann, San Mateo, CA, 1988. (Revised 2nd printing, 1992).
- [Pearl, 1993a] J. Pearl. From Bayesian networks to causal networks. In *Proceedings of the Adaptive Computing and Information Processing Seminar*, pages 25–27, Brunel Conference Centre, London, January 1993. See also *Statistical Science*, 8(3):, 266–269, 1993.
- [Pearl, 1993b] J. Pearl. From conditional oughts to qualitative decision theory. In D. Heckerman and A. Mamdani, editors, *Proceedings of the Ninth Conference on Uncertainty in Artificial Intelligence*, pages 12–20. Morgan Kaufmann, 1993.
- [Pearl, 1994a] J. Pearl. From Adams’ conditionals to default expressions, causal conditionals, and counterfactuals. In E. Eells and B. Skyrms, editors, *Probability and Conditionals*, pages 47–74. Cambridge University Press, Cambridge, MA, 1994.
- [Pearl, 1994b] J. Pearl. A probabilistic calculus of actions. In R. Lopez de Mantaras and D. Poole, editors, *Uncertainty in Artificial Intelligence 10*, pages 454–462. Morgan Kaufmann, San Mateo, CA, 1994.
- [Pearl, 1995] J. Pearl. Causal diagrams for experimental research. *Biometrika*, 82(4):669–710, December 1995.
- [Peng and Reggia, 1990] Y. Peng and J.A. Reggia. *Abductive Inference Models for Diagnostic Problem-Solving*. Springer-Verlag, New York, 1990.
- [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.

- [Rosenthal, 1977] A. Rosenthal. A computer scientist looks at reliability computations. In Barlow et. al., editor, *Reliability and Fault Tree Analysis*, pages 133–152. SIAM, Philadelphia, 1977.
- [Rubin, 1974] D.B. Rubin. Estimating causal effects of treatments in randomized and non-randomized studies. *Journal of Educational Psychology*, 66:688–701, 1974.
- [Rumelhart, 1976] D.E. Rumelhart. Toward an interactive model of reading. Technical Report CHIP-56, University of California, La Jolla, 1976.
- [Shachter, 1986] R.D. Shachter. Evaluating influence diagrams. *Operations Research*, 34(6):871–882, 1986.
- [Shachter, 1988] R.D. Shachter. Probabilistic inference and influence diagrams. *Operations Research*, 36:589–604, 1988.
- [Shachter, 1990] R.D. Shachter. Special issue on influence diagrams. *Networks: An International Journal*, 20(5), August 1990.
- [Shafer and Pearl, 1990] G. Shafer and J. (Eds.) Pearl. *Readings in Uncertain Reasoning*. Morgan Kaufmann, San Mateo, CA, 1990.
- [Simon, 1953] H.A. Simon. Causal ordering and identifiability. In W.C. Hood and T.C. Koopmans, editors, *Studies in Econometric Method*, pages 49–74. John Wiley and Sons, New York, 1953.
- [Sobel, 1990] M.E. Sobel. Effect analysis and causation in linear structural equation models. *Psychometrika*, 55(3):495–515, 1990.
- [Spiegelhalter and Lauritzen, 1990] D.J. Spiegelhalter and S.L. Lauritzen. Sequential updating of conditional probabilities on directed graphical structures. *Networks*, 20(5):579–605, 1990.
- [Spiegelhalter et al., 1993] D.J. Spiegelhalter, S.L. Lauritzen, P.A. Dawid, and R.G. Cowell. Bayesian analysis in expert systems. *Statistical Science*, 8:219–247, 1993.
- [Spirtes et al., 1993] P. Spirtes, C. Glymour, and R. Schienes. *Causation, Prediction, and Search*. Springer-Verlag, New York, 1993.
- [Spohn, 1988] W. Spohn. A general non-probabilistic theory of inductive reasoning. In *Proceedings of the Fourth Workshop on Uncertainty in Artificial Intelligence*, pages 315–322, Minneapolis, MN, 1988.
- [Strotz and Wold, 1960] R.H. Strotz and H.O.A. Wold. Causal models in the social sciences. *Econometrica*, 28:417–427, 1960.
- [Turtle and Croft, 1991] H.R. Turtle and W.B. Croft. Evaluation of an inference network-based retrieval model. *ACM Transactions on Information Systems*, 9(3), July 1991.

- [Verma and Pearl, 1990] T. Verma and J. Pearl. Equivalence and synthesis of causal models. In *Uncertainty in Artificial Intelligence 6*, pages 220–227, Cambridge, MA, 1990. Elsevier Science Publishers.
- [Wermuth and Lauritzen, 1983] N. Wermuth and S.L. Lauritzen. Graphical and recursive models for contingency tables. *Biometrika*, 70:537–552, 1983.
- [Wold, 1964] H. Wold. *Econometric Model Building*. North-Holland, Amsterdam, 1964.
- [Wright, 1921] S. Wright. Correlation and causation. *Journal of Agricultural Research*, 20:557–585, 1921.
- [Wright, 1934] S. Wright. The method of path coefficients. *Ann. Math. Statist.*, 5:161–215, 1934.