

**Computer Science Department Technical Report
University of California
Los Angeles, CA 90024-1596**

KNOWLEDGE DISCOVERY VS. DATA COMPRESSION

**Judea Pearl
Rina Dechter
Thomas Verma**

**July 1991
CSD-910035**

KNOWLEDGE DISCOVERY VS. DATA COMPRESSION *

by

Judea Pearl, Rina Dechter⁽¹⁾ and Thomas Verma

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

Introduction

The purpose of this paper is to summarize recent results in the theory of causal modeling that may have bearing on research in knowledge discovery.

While everyone agrees that data is not knowledge, there is significant diversity of opinions as to what features of the data should be captured before they qualify as genuine knowledge. Naturally, if the data contains regularities, we expect our knowledge to exploit these regularities to yield a more parsimonious representation of the data, i.e., a *summary*. Another facility we expect our knowledge to provide is generalization power -- the power to predict the behavior of some attributes from measurement on others.

Since every summarization of data carries the potential for generalization, most machine learning programs have pursued the objective of *compressing* the data to fit a given format (e.g., a decision tree or a CNF expression) with the hope that the summarization will also provide the necessary generalization. We argue that while this strategy may be adequate for concept formation, classification and prediction, it is not adequate for causal modeling. Causal knowledge embodies a much stronger form of generalization than that achievable by data compression techniques and, hence, new techniques are needed for extracting such knowledge from data. The theories described in Appendix I and II provide a framework for developing these techniques.

Background

Given that statistical analysis cannot distinguish causation from covariation, and assuming that the bulk of our knowledge obtains from passive observations, the questions arise how causal knowledge is ever acquired by humans and how it should be acquired by learning robots. The theories described in the accompanied appendices provide a characterization of the asymmetries that prompt people to perceive causal structures in empirical data and leads to algorithms that emulate this perception.

* This work was partially supported by the National Science Foundation, Grant #IRI-8821444 and by the Air Force Office of Scientific Research, Grant #AFOSR-90-0136.

(1) Current Affiliation: Information and Computer Science, UCI, Irvine, CA, 92717

Both theories are based on a minimal-model semantics; the first (Appendix I) is grounded in probabilistic framework while the second (Appendix II) in categorical, constraint-based framework. This semantics permits the determination of causal directionality without resorting to chronological information and, additionally, provides a distinction between genuine and spurious causes even in the presence of unmeasurable factors.

Using the language of directed acyclic graphs (DAGs), the probabilistic definition of causation reads as follows:

"A variable C has a direct causal influence on a variable E if there is a directed path from C to E in all minimal causal models consistent with the data."

A causal model (i.e., a DAG) is consistent with the data if it can be annotated with parameters so as to define a probability distribution that fits the data (to a given level of fitness). A causal model M is minimal if the set of data consistent with M is not a superset of that consistent with some other model.

Appendix I shows that this formal definition is in line with our common intuition about causation and, in particular, with the perception of causation as a stipulation for future control. For example, the minimal-model semantics sanctions the following rule (Definition 14):

" X has a causal influence on Y if there exists a third variable Z , preceding X , such that Z and Y are dependent and Z and Y are independent given X ."

In this case, Z acts as a virtual control that influences Y via X . The difference is only that Z need not be manipulated under the direct control of the analyst, but can be identified within the data itself. This criterion provides the basis for discovering causal relationships in databases comprised of bare observations, without resorting to controlled experiments. The theory provides similar rules for detecting hidden causes, with and without temporal information.

Appendix II approaches causation from a different viewpoint, appealing to the feature of modularity as its defining characteristic. In simple terms, modularity accounts for the fact that we can ignore the future, though not the past, when it comes to analyzing the present. A *causal ordering* is any ordering of the variables along which the feasible domain of each variable is determined solely by its explicit relationships with its predecessors; relationships with its successors can be ignored.

Remarkably, the minimal-model semantics above can be used to define an *intrinsic directionality* among variables, one that depends only on the tuples in the database but is invariant to the representation used in specifying the database. In this semantics we define a DAG D to be compatible with the database if the latter can be decomposed by the following rule: Fixing its parents in D , each variable must remain unaffected by all other variables, except perhaps its descendants in D .

In both the probabilistic and categorical frameworks, the minimal-model semantics yields operational methods of discovering causal relationships from passive observations, that is, observations obtained without the controlled manipulations of quantities.

Applications to Machine Learning

While our method does not guarantee that the relationships discovered would necessarily correspond to stable physical mechanisms (no method can guarantee that), it nevertheless constitutes an effective filter to ensure that *most* of the learned relationships are stable rather than spurious.

How does it fit into machine learning and knowledge discovery?

If we examine current programs for machine learning we find that what they learn are mostly SITUATION \rightarrow ACTION rules; given a situation S and some (often implicit) goal G , choose an action A that is likely to bring you closer to the goal. This kind of rules represent how an agent should react to changes in the world, but do not represent stable relationships in the world outside the agent. The weaknesses of *reactive* representations are several. First, reactive rules are unstable, as they are vulnerable to many exceptions that cannot easily be encoded. For instance, the rule above should no longer be applicable in situations where the (implicit) goal can be satisfied by easier means than the action A . Second, reactive representations do not tell the agent how to react to novel situations that have not been explicitly encountered in the past. For example, the rule: "If HUNGRY, then GO HUNTING" should be nullified by facts such as:

F_1 - "there is food in the refrigerator", or,

F_2 - "Your wife is on her way to the supermarket"

Causal rules, in contrast, are much more stable. For example, the causal chain

HUNTING \rightarrow FOOD \rightarrow NOT HUNGRY

remains intact regardless of whether facts F_1 and F_2 hold in the database. As a result, a planning program based on a causal model of the environment would be able to automatically pose the intermediate subgoal "get food", and select the appropriate action depending on the available facts.

A major obstacle to learning causal rules has been the confusion between genuine and spurious causes. In other words, a filter is needed to prevent the learning robot from inferring wrong rules such as "HAVING FOOD causes HAVING DRINK", even though past experience reveals that, invariably, drinks became available after food. Such spurious associations, although basic to reactive learning, should be filtered out when knowledge is organized so as to facilitate control over rapidly changing environment. The criteria described in Appendix I constitute such a filter as they distinguish genuine from spurious causation (see Definitions 12-14).

Whereas in traditional learning tasks we attempt to generalize from one set of instances to another, the causal modeling task is to generalize from behavior under one set of conditions to behavior under another set, often a novel one. For example when we say that X causes Y we claim that *every* means Z capable of producing a change in X will also be capable of producing changes in Y . We make this claim without ever observing Z . Instead, we might have observed the behavior of another variable Z' (a virtual control) from which we inferred that the dependence between X and Y is genuinely causal. Thus, causal claims can be thought of as containing an implicit quantification over situations and actions. This then is the reason that causal rules are more stable than their reactive counterparts which, in turn, may explain why causal explanation are much more satisfactory than evidential explanations (e.g., "the cup fell because it broke".) This also may account for the fact that causal rules are regarded as intrinsic to the external objects in the domain, and have become the building blocks of declarative knowledge, while reactive rules serve merely as auxiliary control tools to improve the efficiency of the reasoning agent.

Learning tasks that require the stability of causal theories should benefit, therefore, from shifting their attention from reactive to causal rules. This occurs in applications such as process control, medical diagnosis, financial and economic forecasting and analysis. It would be interesting to examine whether the criteria defined in Appendix II, when incorporated into existing machine learning programs would improve the stability of theories discovered by such programs.

A THEORY OF INFERRED CAUSATION

(In Allen J.A., Fikes, R., and Sandewall, E. (Eds.) *Principles of Knowledge Representation and Reasoning: Proceedings of the Second International Conference*. San Mateo, CA: Morgan Kaufmann. April, 1991, pp. 441-452).

Judea Pearl

< judea@cs.ucla.edu >
Cognitive Systems Laboratory
Computer Science Department
University of California
Los Angeles, CA 90024

T.S. Verma

< verma@cs.ucla.edu >
Cognitive Systems Laboratory
Computer Science Department
University of California
Los Angeles, CA 90024

Abstract

This paper concerns the empirical basis of causation, and addresses the following issues:

1. the clues that might prompt people to perceive causal relationships in uncontrolled observations.
2. the task of inferring causal models from these clues, and
3. whether the models inferred tell us anything useful about the causal mechanisms that underly the observations.

We propose a minimal-model semantics of causation, and show that, contrary to common folklore, genuine causal influences can be distinguished from spurious covariations following standard norms of inductive reasoning. We also establish a sound characterization of the conditions under which such a distinction is possible. We provide an effective algorithm for inferred causation and show that, for a large class of data the algorithm can uncover the direction of causal influences as defined above. Finally, we address the issue of non-temporal causation.

1 Introduction

The study of causation is central to the understanding of human reasoning. Tasks involving changing environments require causal theories which make formal distinctions between causation and logical implication [Geffner, 1989, Lifschitz, 1987, Pearl, 1988a, Shoham, 1988]. In applications such as diagnosis [Patil et al., 1982, Reiter, 1987], qualitative physics [Bobrow, 1985], and plan recognition [Kautz, 1987, Wilensky, 1983], a central task is that of finding a satisfactory *explanation* to a given set of observations, and the meaning of explanation is intimately related to the notion of causation.

Most AI works have given the term "cause" a procedural semantics, attempting to match the way people use it in reasoning tasks, but were not concerned with the experience that prompts people to believe that "a causes b", as opposed to, say, "b causes a" or "c causes both a and b." The question of choosing an appropriate causal ordering received some attention in qualitative physics, where certain interactions attain directionality despite the instantaneous and symmetrical nature of the underlying equations, as in "current causing a voltage drop across the resistor" [Forbus and Gentner, 1986]. In some systems causal ordering is defined as the ordering at which subsets of variables can be solved independently of others [Iwasaki and Simon, 1986], in other systems it follows the way a disturbance is propagated from one variable to others [de Kleer and Brown, 1986]. Yet these choices are made as a matter of convenience, to fit the structure of a given theory, and do not reflect features of the empirical environment which compelled the formation of the theory.

An empirical semantics for causation is important for several reasons. First, an intelligent system attempting to build a workable model of its environment cannot rely exclusively on preprogrammed causal knowledge, but must be able to translate direct observations to cause-and-effect relationships. Second, by tracing empirical origins we stand to obtain an independent gauge for deciding which of the many logics proposed for causal reasoning is sound and/or complete, and which provides a proper account of causal utterances such as "a explains b", "a suggests b", "a tends to cause b", and "a actually caused b", etc.

While the notion of causation is often associated with those of necessity and functional dependence, causal expressions often tolerate exceptions, primarily due to missing variables and coarse descriptions. We say, for example, "reckless driving causes accidents" or "you will fail this course because of your laziness". Suppes [Suppes, 1970] has argued convincingly that most causal utterances in ordinary conversation reflect prob-

abilistic, not categorical relations¹. Thus, probability theory should provide a natural language for capturing causation [Reichenbach, 1956, Good, 1983]. This is especially true when we attempt to infer causation from (noisy) observations – probability calculus remains an unchallenged formalism when it comes to translating statistical data into a system of revisable beliefs.

However, given that statistical analysis is driven by covariation, not causation, and assuming that most human knowledge derives from statistical observations, we must still identify the clues that prompt people to perceive causal relationships in the data, and we must find a computational model that emulates this perception.

Temporal precedence is normally assumed essential for defining causation, and it is undoubtedly one of the most important clues that people use to distinguish causal from other types of associations. Accordingly, most theories of causation invoke an explicit requirement that a cause precedes its effect in time [Good, 1983, Reichenbach, 1956, Shoham, 1988, Suppes, 1970]. Yet temporal information alone cannot distinguish genuine causation from spurious associations caused by unknown factors. In fact the statistical and philosophical literature has adamantly warned analysts that, unless one knows in advance all causally relevant factors, or unless one can carefully manipulate some variables, no genuine causal inferences are possible [Cartwright, 1989, Cliff, 1983, Eells and Sober, 1983, Fisher, 1953, Gardenfors, 1988, Holland, 1986, Skyrms, 1986]². Neither condition is realizable in normal learning environments, and the question remains how causal knowledge is ever acquired from experience.

This paper introduces a minimal-model semantics of causation which provides a plausible account for how causal models could be inferred from observations. Using this semantics we show that genuine causal influences can in many cases be distinguished from spurious covariations and, moreover, the direction of causal influences can often be determined without resorting to chronological information. (Although, when available, chronological information can significantly simplify the modeling task.) Such semantics should be applicable, therefore, to the organization of concurrent events or events whose chronological precedence cannot be determined with precision, (e.g. “old age explains disabilities”) in the spirit of Glymour [Glymour et al., 1987] and Simon [Simon, 1954].

¹See [Dechter and Pearl, 1990] for a treatment of causation in the context of categorical data.

²Some of the popular quotes are: “No causation without manipulation”, [Holland, 1986], “No causes in, no causes out”, [Cartwright, 1989] “No computer program can take account of variables that are not in the analysis”, [Cliff, 1983].

This paper is organized as follows. In Section 2 we define the notions of causal models and causal theories, and describe the task of causal modeling as an identification game scientists play against Nature. In Section 3 we introduce the minimal-model semantics of causation and exemplify its operability and plausibility on a simple example. Section 4 identifies conditions under which effective algorithms exist that uncover the structure of causal influences as defined above. One such algorithm (called IC) is introduced in Section 5, and is shown to be sound for the class of stable distributions, even when some variables are not observable³. Section 6 extracts from the IC-algorithm the essential conditions under which causal influences are identified and proposes these as independent definitions of genuine influences and spurious associations, with and without temporal information. Section 7 provides an intuitive justification for the definitions proposed in Section 6, showing that our theory conforms to the common understanding of causation as a stipulation of stable behavior under external interventions. The definitions are shown to be in line with accepted standards of controlled experimentation, save for requiring the identification of “virtual” experimental conditions within the data itself. In Section 8 we invoke the “virtual control” metaphor to elucidate how causal relationships can still be ascertained in the absence of temporal information. We then offer an explanation for the puzzling, yet universal agreement between the temporal and the statistical aspects of causation.

2 The Causal Modeling Framework

We view the task of causal modeling as an identification game which scientists play against Nature. Nature possesses stable causal mechanisms which, on a microscopic level are deterministic functional relationships between variables, some of which are unobservable. These mechanisms are organized in the form of an acyclic schema which the scientist attempts to identify.

Definition 1 *A causal model of a set of variables U is a directed acyclic graph (dag), in which each node corresponds to a distinct element of U .*

The nodes of the dag correspond to the variables under analysis, while the links denote direct causal influences among the variables. The causal model serves as a blue print for forming a “causal theory” – a precise specification of how each variable is influenced by its parents in the dag. Here we assume that Nature is at liberty to impose arbitrary functional relationships between each effect and its causes and then to perturb these relationships by introducing arbitrary (yet mutually independent) disturbances. These disturbances reflect “hidden” or unmeasurable conditions and exceptions

³Proofs can be found in [Verma, 1991].

which Nature chooses to govern by some undisclosed probability function.

Definition 2 A causal theory is a pair $T = \langle D, \Theta_D \rangle$ consisting of a causal model D and a set of parameters Θ_D compatible with D . Θ_D assigns a function $x_i = f_i[\text{pa}(x_i), \epsilon_i]$ and a probability measure g_i , to each $x_i \in U$, where $\text{pa}(x_i)$ are the parents of x_i in D and each ϵ_i is a random disturbance distributed according to g_i , independently of the other ϵ 's and of any preceding variable x_j ; $0 < j < i$

This requirement of independence renders the disturbances "local" to each parents-child family; disturbances that influence several families simultaneously will be treated explicitly as "latent" variables (see Definition 3).

Once a causal theory T is formed, it defines a joint probability distribution $P(T)$ over the variables in the system, and this distribution reflects some features of the causal model (e.g., each variable must be independent of its grandparents, given the values of its parents). Nature then permits the scientist to inspect a select subset $O \subseteq U$ of "observed" variables, and to ask questions about the probability distribution over the observables, but hides the underlying causal theory as well as the structure of the causal model. We investigate the feasibility of recovering the topology of the dag, D , from features of the probability distribution.⁴

3 Model preferences (Occam's razor)

In principle, U being unknown, there is an unbounded number of models that would fit a given distribution, each invoking a different set of "hidden" variables and each connecting the observed variables through different causal relationships. Therefore with no restriction on the type of models considered, the scientist is unable to make any meaningful assertions about the structure underlying the phenomena. Likewise, even assuming $U = O$ but lacking temporal information, he/she can never rule out the possibility that the underlying model is a complete (acyclic) graph; a structure that, with the right choice of parameters can mimic (see Definition 4) the behavior of any other

⁴This formulation invokes several idealizations of the actual task of scientific discovery. It assumes, for example, that the scientist obtains the distribution directly, rather than events sampled from the distribution. This assumption is justified when a large sample is available, sufficient to reveal all the dependencies embedded in the distribution. Additionally, we assume that the observed variables actually appear in the original causal theory and are not some aggregate thereof. Aggregation might result in feedback loops which we do not discuss in this paper. Our theory also takes variables as the primitive entities in the language, not events which permits us to include "enabling" and "preventing" relationships as part of the mechanism.

model, regardless of the variable ordering. However, following the standard method of scientific induction, it is reasonable to rule out any model for which we find a simpler, less expressive model, equally consistent with the data (see Definition 6). Models that survive this selection are called "minimal models" and with this notion, we can construct our definition of inferred causation:

"A variable X is said to have a causal influence on a variable Y if a strictly directed path from X to Y exists in every minimal model consistent with the data"

Definition 3 Given a set of observable variables $O \subseteq U$, a latent structure is a pair $L = \langle D, O \rangle$ where D is a causal model over U .

Definition 4 One latent structure $L = \langle D, O \rangle$ is preferred to another $L' = \langle D', O \rangle$ (written $L \preceq L'$) iff D' can mimic D over O , i.e. for every Θ_D there exists a $\Theta_{D'}$, s.t. $P_{[O]}(\langle D', \Theta_{D'} \rangle) = P_{[O]}(\langle D, \Theta_D \rangle)$

Two latent structures are equivalent, written $L' \equiv L$, iff $L \preceq L'$ and $L \succeq L'$.

Note that the preference for simplicity imposed by Definition 4 is gauged by the expressive power of a model, not by its syntactic description. For example, one latent structure L_1 may invoke many more parameters than L_2 and still be preferred, if L_2 is capable of accommodating a richer set of probability distributions over the observables. One reason scientists prefer simpler models is that such models are more constrained, thus more falsifiable; they provide the scientist with less opportunities to overfit the data hindsightedly and, therefore attain greater credibility [Pearl, 1978, Popper, 1959].

We also note that the set of dependencies induced by a causal model provides a measure of its expressive power, i.e., its power of mimicking other models. Indeed, L_1 cannot be preferred to L_2 if there is even one observable dependency that is induced by L_1 and not by L_2 . Thus, tests for preference and equivalence can often be reduced to tests of induced dependencies which, in turn, can be determined directly from the topology of the dags, without ever concerning ourselves with the set of parameters. (For example, see Theorem 1 below and [Frydenberg, 1989, Pearl et al., 1989, Verma and Pearl, 1990]).

Definition 5 A latent structure L is minimal with respect to a class \mathcal{L} of latent structures iff for every $L' \in \mathcal{L}$, $L \equiv L'$ whenever $L' \preceq L$.

Definition 6 $L = \langle D, O \rangle$ is consistent with a distribution \hat{P} over O if D can accommodate some theory that generates \hat{P} , i.e. there exists a Θ_D s.t. $P_{[O]}(\langle D, \Theta_D \rangle) = \hat{P}$

Clearly, a necessary (and often sufficient) condition for

L to be consistent with \hat{P} , is that the structure of L can account for all the dependencies embodied in \hat{P} .

Definition 7 (Inferred Causation) Given \hat{P} , a variable C has a causal influence on E iff there exists a directed path $C \rightarrow^* E$ in every minimal latent structure consistent with \hat{P} .

We view this definition as normative, because it is based on one of the least disputed norms of scientific investigation: Occam's razor in its semantical casting. However, as with any scientific inquiry, we make no claims that this definition is guaranteed to always identify stable physical mechanisms in nature; it identifies the only mechanisms we can plausibly infer from non-experimental data.

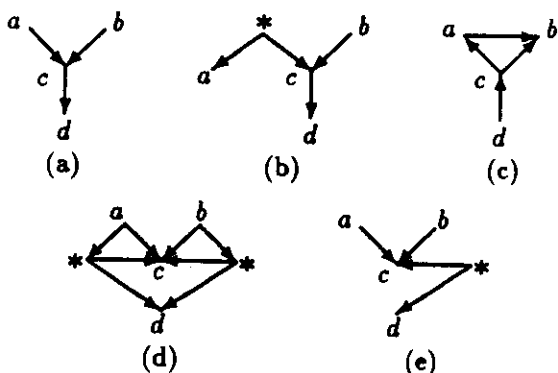


Figure 1: Causal models illustrating the soundness of $c \rightarrow d$. The node (*) represents a hidden variable.

As an example of a causal relation that is identified by the definition above, imagine that observations taken over four variables $\{a, b, c, d\}$ reveal two vanishing dependencies: “ a is independent of b ” and “ d is independent of $\{a, b\}$ given c ”. Assume further that the data reveals *no other* independence, except those that logically follow from these two. This dependence pattern would be typical for example, of the following variables: $a = \text{having cold}$, $b = \text{having hay-fever}$, $c = \text{having to sneeze}$, $d = \text{having to wipe ones nose}$. It is not hard to see that any model which explains the dependence between c and d by an arrow from d to c , or by a hidden common cause (*) between the two, cannot be minimal, because any such model would be able to out-mimic the one shown in Figure 1(a) which reflects all observed independencies. For example, the model of Figure 1(c), unlike that of Figure 1(a), accommodates distributions with arbitrary relations between a and b . Similarly, Figure 1(d) is not minimal as it fails to impose the conditional independence between d and $\{a, b\}$ given c . In contrast, Figure 1(e) is not consistent with the data since it imposes a marginal independence between $\{a, b\}$ and d , which was not observed.

4 Proof Theory and Stable Distributions

It turns out that while the minimality principle is sufficient for forming a normative and operational theory of causation, it does not guarantee that the search through the vast space of minimal models would be computationally practical. If Nature truly conspires to conceal the structure of the underlying model she could still annotate that model with a distribution that matches many minimal models, having totally disparate structures. To facilitate an effective proof theory, we rule out such eventualities, and impose a restriction on the distribution called “stability” (or “dag-isomorphism” in [Pearl, 1988b]). It conveys the assumption that all vanishing dependencies are structural, not formed by incidental equalities of numerical parameters⁵.

Definition 8 Let $I(P)$ denote the set of all conditional independence relationships embodied in P . A causal theory $T = \langle D, \Theta_D \rangle$ generates a stable distribution iff it contains no extraneous independences, i.e. $I(P(\langle D, \Theta_D \rangle)) \subseteq I(P(\langle D, \Theta'_D \rangle))$ for any set of parameters Θ'_D .

With the added assumption of stability, every distribution has a unique causal model (up to equivalence), as long as there are no hidden variables. This uniqueness follows from the fact the structural constraints that an underlying dag imposes upon the probability distribution are equivalent to a finite set of conditional independence relationships asserting that, given its parents, each variable is conditionally independent of all its non-descendants. Therefore two causal models are equivalent (i.e. they can mimic each other) if and only if they relay the same dependency information. The following theorem, which is founded upon the dependency information, states necessary and sufficient conditions for equivalence of causal models which contain no hidden variables.

Theorem 1 [Verma and Pearl, 1990] When $U = O$, two causal models are equivalent iff their dags have the same links and same set of uncoupled head-to-head nodes⁶.

The search for the minimal model then boils down to recovering the structure of the underlying dag from queries about the dependencies portrayed in that dag. This search is exponential in general, but simplifies significantly when the underlying

⁵It is possible to show that, if the parameters are chosen at random from any reasonable distribution, then any unstable distribution has measure zero [Spirtes et al., 1989]. Stability precludes deterministic constraints. Less restrictive assumptions are treated in [Geiger et al., 1990].

⁶i.e. converging arrows emanating from non-adjacent nodes, such as $a \rightarrow c \leftarrow b$ in Figure 1(a).

structure is sparse (see [Spirites and Glymour, 1991, Verma and Pearl, 1990] for such algorithms).

Unfortunately, the constraints that a latent structure impose upon the distribution cannot be completely characterized by any set of dependency statements. However, the maximal set of sound constraints can be identified [Verma and Pearl, 1990] and it is this set that permits us to recover sound fragments of latent structures.

5 Recovering Latent Structures

When Nature decides to "hide" some variables, the observed distribution \hat{P} need no longer be stable relative to the observable set O , i.e. \hat{P} may result from many equivalent minimal latent structures, each containing any number of hidden variables. Fortunately, rather than having to search through this unbounded space of latent structures, it turns out that for every latent structure L , there is a dependency-equivalent latent structure called the projection of L on O in which every unobserved node is a root node with exactly two observed children:

Definition 9 A latent structure $L_{[O]} = \langle D_{[O]}, O \rangle$ is a projection of another latent structure L iff

1. Every unobservable variable of $D_{[O]}$ is a parentless common cause of exactly two non-adjacent observable variables.
2. For every stable distribution P generated by L , there exists a stable distribution P' generated by $L_{[O]}$ such that $I(P_{[O]}) = I(P'_{[O]})$.

Theorem 2 Any latent structure has at least one projection (identifiable in linear time).

It is convenient to represent projections by bi-directional graph with only the observed variables as vertices (i.e., leaving the hidden variables implicit). Each bi-directed link in such a graph represents a common hidden cause of the variables corresponding to the link's end points.

Theorem 2 renders our definition of inferred causation (Definition 7) operational; we will show (Theorem 3) that if a certain link exists in a distinguished projection of any minimal model of \hat{P} , it must indicate the existence of a causal path in every minimal model of \hat{P} . Thus the search reduces to finding a projection of any minimal model of \hat{P} and identifying the appropriate links. Remarkably, these links can be identified by a simple procedure, the IC-algorithm, that is not more complex than that which recovers the unique minimal model in the case of fully observable structures.

IC-Algorithm (Inductive Causation)

Input: \hat{P} a sampled distribution.

Output: $\text{core}(\hat{P})$ a marked hybrid acyclic graph.

1. For each pair of variables a and b , search for a set S_{ab} such that (a, S_{ab}, b) is in $I(\hat{P})$, namely a and b are independent in \hat{P} , conditioned on S_{ab} . If there is no such S_{ab} , place an undirected link between the variables.
2. For each pair of non-adjacent variables a and b with a common neighbor c , check if $c \in S_{ab}$. If it is, then continue. If it is not, then add arrowheads pointing at c , (i.e. $a \rightarrow c \leftarrow b$).
3. Form $\text{core}(\hat{P})$ by recursively adding arrowheads according to the following two rules:⁷
If \overline{ab} and there is a strictly directed path from a to b then add an arrowhead at b .
If a and b are not adjacent but \overline{ac} and $c \rightarrow b$, then direct the link $c \rightarrow b$.
4. If \overline{ab} then mark every uni-directed link $b \rightarrow c$ in which c is not adjacent to a .

The result of this procedure is a substructure called $\text{core}(\hat{P})$ in which every marked uni-directed arrow $X \rightarrow Y$ stands for the statement: " X has a causal influence on Y (in all minimal latent structures consistent with the data)". We call these relationships "genuine" causal influences (e.g. $c \rightarrow d$ in previous Figure 1).

Definition 10 For any latent structure L , $\text{core}(L)$ is defined as the hybrid graph⁸ satisfying (1) two nodes are adjacent in $\text{core}(L)$ iff they are adjacent or they have a common unobserved cause in every projection of L , and (2) a link between a and b has an arrowhead pointing at b iff $a \rightarrow b$ or a and b have a common unobserved cause in every projection of L .

Theorem 3 For any latent structure $L = \langle D, O \rangle$ and an associated theory $T = \langle D, \Theta_D \rangle$ if $P(T)$ is stable then $\text{core}(L) = \text{core}(P_{[O]}(T))$.

Corollary 1 If every link of the directed path $C \rightarrow \dots \rightarrow E$ is marked in $\text{core}(\hat{P})$ then C has a causal influence on E according to \hat{P} .

6 Probabilistic Definitions for Causal Relations

The IC-algorithm takes a distribution \hat{P} and outputs a dag, some of its links are marked uni-directional

⁷ \overline{ab} denotes adjacency, i.e. $a - b$, $a \rightarrow b$, $a \leftarrow b$ or $a \leftrightarrow b$.
 $a\bar{b}$ denotes either $a \rightarrow b$ or $a \leftarrow b$.

⁸In a hybrid graph links may be undirected, uni-directed or bi-directed.

(denoting genuine causation), some are unmarked unidirectional (denoting potential causation), some are bidirectional (denoting spurious association) and some are undirected (denoting relationships that remain undetermined). The conditions which give rise to these labelings constitute operational definitions for the various kinds of causal relationships. In this section we present explicit definitions of potential and genuine causation, as they emerge from Theorem 3 and the IC-algorithm. Note that in all these definitions, the criterion for causation between two variables, X and Y , will require that a third variable Z exhibit a specific pattern of interactions with X and Y . This is not surprising, since the very essence of causal claims is to stipulate the behavior of X and Y under the influence of a third variable, one that corresponds to an external control of X . Therefore, our definitions are in line with the paradigm of “no causation without manipulation” [Holland, 1986]). The difference is only that the variable Z , acting as a virtual control of X , must be identified within the data itself. The IC-algorithm provides a systematic way of searching for variables Z that qualify as virtual controls.

Detailed discussions of these definitions in terms of virtual control are given in Sections 7 and 8.

Definition 11 (Potential Cause) *A variable X has a potential causal influence on another variable Y (inferable from \hat{P}), if*

1. X and Y are dependent in every context.
2. There exists a variable Z and a context S such that
 - (i) X and Z are independent given S
 - (ii) Z and Y are dependent given S

Note that this definition precludes a variable X from being a potential cause of itself or of any other variable which functionally determines X .

Definition 12 (Genuine Cause) *A variable X has a genuine causal influence on another variable Y if there exists a variable Z such that either:*

1. X is a potential cause of Y and there exists a context S satisfying:
 - (i) Z is a potential cause of X
 - (ii) Z and Y are dependent given S .
 - (iii) Z and Y are independent given $S \cup X$,
 or,
2. X is a genuine cause of Z and Z is a genuine cause of Y .

Definition 13 (Spurious Association) *Two variables X and Y are spuriously associated if they are*

dependent in some context S and there exists two other variables Z_1 and Z_2 such that:

1. Z_1 and X are dependent given S
2. Z_2 and Y are dependent given S
3. Z_1 and Y are independent given S
4. Z_2 and X are independent given S

Succinctly, using the predicates I and $\neg I$ to denote independence and dependence respectively, the conditions above can be written:

1. $\neg I(Z_1, X|S)$
2. $\neg I(Z_2, Y|S)$
3. $I(Z_1, Y|S)$
4. $I(Z_2, X|S)$

Definition 11 was formulated in [Pearl, 1990] as a relation between events (rather than variables) with the added condition $P(Y|X) > P(Y)$ in the spirit of [Good, 1983, Reichenbach, 1956, Suppes, 1970]. Condition 1 in Definition 12 may be established either by statistical methods (per Definition 11) or by other sources of information e.g., experimental studies or temporal succession (i.e. that Z precedes X in time).

When temporal information is available, as it is assumed in the most theories of causality ([Granger, 1988, Spohn, 1983, Suppes, 1970]), then Definitions 12 and 13 simplify considerably because every variable preceding and adjacent to X now qualifies as a “potential cause” of X . Moreover, adjacency (i.e. condition 1 of Definition 11) is not required as long as the context S is confined to be earlier than S . These considerations lead to simpler conditions distinguishing genuine from spurious causes as shown next.

Definition 14 (Genuine Causation with temporal information) *A variable X has a causal influence on Y if there is a third variable Z and a context S , both occurring before X such that:*

1. $\neg I(Z, Y|S)$
2. $I(Z, Y|S \cup X)$

Definition 15 (Spurious Association with temporal information) *Two variables X and Y are spuriously associated if they are dependent in some context S , X precedes Y and there exists a variable Z satisfying:*

1. $I(Z, Y|S)$
2. $\neg I(Z, X|S)$

7 Causal Intuition and Virtual Experiments

This section explains how the formulation introduced above conforms to common intuition about causation and, in particular, how symmetric probabilistic dependencies can be transformed into judgements about causal influences. We shall first uncover the intuition behind Definition 14, assuming the availability of temporal information, then (in Section 8) generalize to non-temporal data, per Definition 12.

The common intuition about causation is captured by the heuristic definition [Rubin, 1989]: “ X is a cause for Y if an external agent interfering only with X can affect Y ”.

Thus, causal claims are much bolder than those made by probability statements; not only do they summarize relationships that hold in the distribution underlying the data, but they also predict relationships that should hold when the distribution undergoes changes, such as those inferable from external intervention. The claim “ X causes Y ” asserts the existence of a *stable* dependence between X and Y , one that cannot be attributed to some prior cause common to both, and one that should be preserved when an exogenous control is applied to X .

This intuition requires the formalization of three notions:

1. That the intervening agent be “external” (or “exogenous”)
2. That the agent can “affect” Y
3. That the agent interferes “only” with X

If we label the behavior of the intervening agent by a variable Z , then these notions can be given the following probabilistic explications:

1. **Externality of Z :** Variations in Z must be independent of any factors W which precede X , i.e.,

$$I(Z, W) \quad \forall W : t_w < t_x \quad (1)$$

2. **Control:** For Z to effect changes in Y (via X) we require that Z and Y be dependent, written:

$$\neg I(Z, Y) \quad (2)$$

3. **Locality:** To ensure that Z interferes “only” with X , i.e., that its entire effect on Y is mediated by X , we use the conditional independence assertion:

$$I(Z, Y | X) \quad (3)$$

to read “ Z is independent of Y , given X ”.

Note that (2) and (3) imply (by the axioms of conditional independence [Pearl, 1988b]) that X and Y are dependent, namely, $\neg I(X, Y)$.

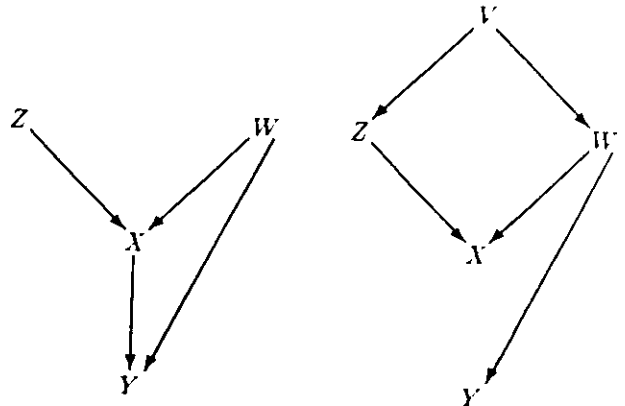


Figure 2

Figure 3

Conditions (1) through (3) constitute the traditional premises behind controlled statistical experiments, with (1) reflecting the requirement that units selected for the experiment be chosen at random from the population under study. They guarantee that any dependency observed between X and Y cannot be explained away by holding fixed some factor W preceding X , hence it must be attributed to genuine causation. The sufficiency of these premises is clearly not a theorem of probability theory, as it relies on temporal relationships among the variables. However, it can be derived from probability theory together with Reichenbach’s principle [Reichenbach, 1956], stating that every dependence $\neg I(X, Y)$ requires a causal explanation, namely either one of the variables causes the other, or there must be a variable W preceding X and Y such that $I(X, Y | W)$ (see Figure 2). Indeed, if there is no back path from Z to Y through W (Eq. (1)) and no direct path from Z to Y avoiding X (Eq. (3)) then there must be a causal path from X to Y that is responsible for the dependence in Eq. (2)⁹.

In non-experimental situations it is not practical to detach X completely from its natural surrounding and to subject it to the exclusive control of an exogenous (and randomized) variable Z . Instead, one could view some of X ’s natural causes as “virtual controls” and, provided certain conditions are met, use the latter to reveal non-spurious causal relationship between X and Y . In so doing we compromise, of course, condition (1), because we can no longer guarantee that those natural causes of X are not themselves affected by other causes which, in turn, might influence Y (see Figure 3). However, it turns out that for stable distributions, conditions (2) and (3) are sufficient to guarantee that the association between X and Y is non-spurious, thus justifying Definition 14 for genuine causation.

The intuition goes as follows (see Figure 3): If the de-

⁹Cartwright [Cartwright, 1989] offers a sufficiency proof in the context of linear models.

pendency between Z and Y (and similarly, between X and Y) is spurious, namely, X and Y are merely manifestations of some common cause W , there is no reason then for X to screen-off Y from Z , and condition (2) should be violated. In case condition (2) is accidentally satisfied by some strange combination of parameters, it is bound to be “unstable”, as it will be perturbed with any slight change of experimental conditions.

Conditions (2) and (3) are identical to those in Definition 14, save for the context S which is common to both. The inclusion of the fixed context S is legitimized by noting that if $P(X, Y, Z)$ is a marginal of a stable distribution, then so is the conditional distribution $P(X, Y, Z|S = s)$, as long as S corresponds to variables which precede X .

Definition 14 constitutes an alternative way of recovering causal structures, more flexible than the IC-algorithm; we search the data for three variables Z, X, Y (in this temporal order) that satisfy the two conditions in some context $S = s$, and when such a triple is found, X is proclaimed to have a genuine causal influence on Y . Clearly, permitting an arbitrary context S increases the number of genuine causal influences that can be identified in any given data; marginal independencies and even 1-place conditional independencies are rare phenomenon.

Note that failing to satisfy the test for genuine causation does not mean that such relationship is necessarily absent between the quantities under study. Rather, it means that the data available cannot substantiate the claim of genuine causation. To further test such claims one may need to either conduct experimental studies, or consult a richer data set where virtual control variables are found.

In testing this modeling scheme on real life data, we have examined the observations reported in Sewal Wright’s seminal paper “Corn and Hog Correlations” [Wright, 1925]. As expected, corn-price (X) can clearly be identified as a cause of hog-price (Y), not the other way around. The reason lies in the existence of the variable corn-crop (Z) that, by satisfying the conditions of Definition 14 (with $S = \emptyset$), acts as a virtual control of X (see Figure 2). To test for the possibility of reciprocal causation, one can try to find a virtual controller for Y , for example, the amount of hog-breeding (Z'). However, it turns out that Z' is not screened off from X by Y (possibly because corn prices exert direct influence over farmer’s decision to breed more hogs), hence, failing condition 3, Y disqualifies as a genuine cause of X . Such distinctions are important to policy makers in deciding, for example, which commodity, corn or hog, should be subsidized or taxed.

8 Non-Temporal Causation and Statistical Time

When temporal information is unavailable the condition that Z precede X (Definition 14) cannot be tested directly and must be replaced by an equivalent condition, based on dependence information. As it turns out, the only reason we had to require that Z precede X is to rule out the possibility that Z is a causal consequence of X ; if it were a consequence of X then the dependency between Z and Y could easily be explained away by a common cause W of X and Y (see Figure 2).

The information that permits us to conclude that one variable is not a causal consequence of another comes in the form of an “intransitive triplet”, such as the variables a, b and c in Figure 1(a) satisfying: $I(a, b)$, $\neg I(a, c)$ and $\neg I(b, c)$. The argument goes as follows: If we create conditions (fixing S_{ab}) where two variables, a and b , are each correlated with a third variable c but are independent of each other, then the third variable cannot act as a cause of a or b , (recall that in stable distributions, common causes induce dependence among their effects); it must be either their common effect, $a - c - b$, or be associated with a and b via common causes, forming a pattern such as $a - c - b$. This is indeed the eventuality that permits our algorithm to begin orienting edges in the graph (step 2), and assign arrowheads pointing at c . It is also this intransitive pattern which is used to ensure that X is not a consequence of Y (in Definition 11) and that Z is not a consequence of X (in Definition 12). In Definition 14 we have two intransitive triplets, (Z_1, X, Y) and (X, Y, Z_2) , thus ruling out direct causal influence between X and Y , implying spurious associations as the only explanation for their dependence.

This interpretation of the intransitive triple is in line with the “virtual control” view of causation. For example, one of the reasons people insist that the rain causes the grass to become wet, and not the other way around, is that they can find other means of getting the grass wet, totally independent of the rain. Transferred to our chain $a - c - b$, we can preclude c from being a cause of a if we find another means (b) of potentially controlling c without affecting a [Pearl, 1988a, p. 396].

Determining the direction of causal influences from nontemporal data raises some interesting philosophical questions about the nature of time and causal explanations. For example, can the orientation assigned to the arrow $X \rightarrow Y$ in Definition 14 ever clash with temporal information (say by a subsequent discovery that Y precedes X)? Alternatively, since the rationale behind Definition 14 is based on strong intuitions about how causal influences should behave (statistically), it is apparent that such clashes, if they occur, are rather rare. The question arises then, why? Why should orientations determined solely by statistical dependencies

have anything to do with the flow of time?

In human discourse, causal explanations indeed carry two connotations, temporal and statistical. The temporal aspect is represented by the convention that a cause should precede its effect. The statistical aspect expects causal explanations (once accounted for) to screen off their effects, i.e., render their effects conditionally independent¹⁰. More generally, causal explanations are expected to obey many of the rules that govern paths in a directed acyclic graphs (e.g., the intransitive triplet criterion for potential causation, Section 7). This leads to the observation that, if agreement is to hold between the temporal and statistical aspects of causation, natural statistical phenomena must exhibit some basic temporal bias. Indeed, we often encounter phenomenon where knowledge of a present state renders the variables of the future state conditionally independent (e.g., multi-variables economic time series as in Eq. (4) below). We rarely find the converse phenomenon, where knowledge of the present state would render the components of the past state conditionally independent. The question arises whether there is any compelling reason for this temporal bias.

A convenient way to articulate this bias is through the notion of “Statistical Time”.

Definition 16 (Statistical Time) *Given an empirical distribution P , a statistical time of P is any ordering of the variables that agrees with at least one minimal causal model consistent with P .*

We see, for example, that a scalar Markov-chain process has many statistical times; one coinciding with the physical time, one opposite to it and the others correspond to any time ordering of the variables away from some chosen variable. On the other hand a process governed by two coupled Markov chains,

$$\begin{aligned} X_t &= \alpha X_{t-1} + \beta Y_{t-1} + \xi_t \\ Y_t &= \gamma X_{t-1} + \delta Y_{t-1} + \xi'_t, \end{aligned} \quad (4)$$

has only one statistical time – the one coinciding with

¹⁰This principle, known as Reichenbach’s “conjunctive fork” or “common-cause” criterion [Reichenbach, 1956, Suppes and Zanotti, 1981] has been criticized by Salmon [Salmon, 1984], who showed that some events would qualify as causal explanations though they fail to meet Reichenbach’s criterion. Salmon admits, however, that when a conjunctive forks does occur, the screening off variable is expected to be the cause of the other two, not the effect [Salmon, 1984, p. 167]. He notes that it is difficult to find physically meaningful examples where a response variable renders its two causes conditionally independent (although this would not violate any axiom of probability theory). This asymmetry is further evidence that humans tend to reject causal theories that yield unstable distributions.

the physical time¹¹. Indeed, running the IC-algorithm on samples taken from such a process, while suppressing all temporal information, quickly identifies the components of X_{t-1} and Y_{t-1} as genuine causes of X_t and Y_t . This can be seen from Definition 11, where X_{t-2} qualifies as a potential cause of X_{t-1} using $Z = Y_{t-2}$ and $S = \{X_{t-3}, Y_{t-3}\}$, and Definition 12, where X_{t-1} qualifies as a genuine cause of X_t using $Z = X_{t-2}$ and $S = \{Y_{t-1}\}$ of X_t .

The temporal bias postulated earlier can be expressed as follows:

Conjecture 1 (Temporal Bias) *In most natural phenomenon, the physical time coincides with at least one statistical time.*

Reichenbach [Reichenbach, 1956] attributed the asymmetry associated with his conjunctive fork to the second law of thermodynamics. We are not sure at this point whether the second law can provide a full account of the temporal bias as defined above, since the influence of the external noise ξ_t and ξ'_t renders the process in (4) nonconservative¹². What is clear, however, is that the temporal bias is *language dependent*. For example, expressing Eq.(4) in a different coordinate system (say, using a unitary transformation $(X', Y') = U(X, Y)$), it is possible to make the statistical time (in the (X', Y') representation) run contrary to the physical time. This suggests that the apparent agreement between the physical and statistical times is a byproduct of human choice of linguistic primitives and, moreover, that the choice is compelled by a survival pressure to facilitate predictions at the expense of diagnosis and planning.

9 Conclusions

The theory presented in this paper should dispel the belief that statistical analysis can never distinguish genuine causation from spurious covariation. This belief, shaped and nurtured by generations of statisticians [Fisher, 1953, Keynes, 1939, Ling, 1983, Niles, 1922] has been a major hindrance in the way of developing a satisfactory, non-circular account of causation. In the words of Gardenfors [Gardenfors, 1988, page 193]:

In order to distinguish genuine from spurious causes, we must already know the causally relevant background factors. ... Further, the extra amount of information is substantial: In order to determine whether C is a cause of E, *all* causally relevant background factors must be available. It seems clear that we

¹¹ ξ_t and ξ'_t are assumed to be two independent, white noise time series. Also $\alpha \neq \delta$ and $\gamma \neq \beta$.

¹²We are grateful to Seth Lloyd for this observation.

often have determinate beliefs about causal relations between events, even if we do not know exactly which factors are causally relevant to the events in question¹³.

This paper shows that such extra information is often unnecessary: Under the assumptions of model-minimality (and/or stability), there are patterns of dependencies that should be sufficient to uncover genuine causal relationships. These relationships cannot be attributed to hidden causes lest we violate one of the basic maxims of scientific methodology: the semantical version of Occam's razor. Adherence to this maxim explains why humans reach consensus regarding the directionality and nonspuriousness of causal relationships, in the face of opposing alternatives, perfectly consistent with experience. Echoing Cartwright [Cartwright, 1989] we summarize our claim with the slogan "No Causes In, Some Causes Out".

From a methodological viewpoint, our theory should settle some of the on going disputes regarding the validity of path-analytic approaches to causal modeling in the social sciences [Freedman, 1987, Ling, 1983]. It shows that the basic philosophy governing path-analytic methods is legitimate, faithfully adhering to the traditional norms of scientific investigation. At the same time our results also explicate the assumptions upon which these methods are based, and the conditions that must be fulfilled before claims made by these methods can be accepted. Specifically, our analysis makes it clear that causal modeling must begin with *vanishing (conditional) dependencies* (i.e. missing links in their graphical representations). Models that embody no vanishing dependencies contain no virtual control variables, hence, the causal component of their claims cannot be substantiated by observational studies. With such models, the data can be used only for estimating the parameters of the causal links once we are absolutely sure of the causal structure, but the structure itself, and especially the directionality of the links, cannot be inferred from the data. Unfortunately, such models are often employed in the social and behavioral sciences e.g. [Kenny, 1979].

On the practical side, we have shown that the assumption of model minimality, together with that of "stability" (no accidental independencies) lead to an effective algorithm of recovering causal structures, transparent as well as latent. Simulation studies conducted at our laboratory show that networks containing tens of variables require less than 5000 samples to have their structure recovered by the algorithm. For example, 1000 samples taken from the process shown in Eq. (5), each containing ten successive X,Y pairs, were sufficient for recovering its double-chain structure (and the correct direction of time). The greater the noise,

the quicker the recovery.

Another result of practical importance is the following: Given a proposed causal theory of some phenomenon, our algorithm can identify in linear time those causal relationships that could potentially be substantiated by observational studies, and those whose directionality non-spuriousness can only be determined by controlled, manipulative experiments.

It should also be interesting to explore how the new criteria for causation could benefit current research in machine learning. In some sense, our method resembles a search through elements of a version space [Mitchell, 1982], where each hypothesis stands for a causal theory. Unfortunately, this is where the resemblance ends. The prevailing paradigm in the machine learning literature has been to define each hypothesis (or theory, or concept) as a subset of observable instances; once we observe the entire extension of this subset, the hypothesis is defined unambiguously. This is not the case in causal modeling. Even if the training sample exhausts the hypothesis subset (in our case, this corresponds to observing P precisely), we are still left with a vast number of equivalent causal theories, each stipulating a drastically different set of causal claims. Fitness to data, therefore, is an insufficient criterion for validating causal theories. Whereas in traditional learning tasks we attempt to generalize from one set of instances to another, the causal modeling task is to generalize from behavior under one set of conditions to behavior under another set. Causal models should therefore be chosen by a criterion that challenges their stability against changing conditions, and these show up in the data in the form of virtual control variables. Thus, the dependence patterns identified by definition 11 through 14 constitute islands of stability as well as virtual validation tests for causal models. It would be interesting to examine whether these criteria, when incorporated into existing machine learning programs would improve the stability of theories discovered by such programs.

Acknowledgement

We are grateful to Clark Glymour for posing the problem of equivalence in latent structures. Some of the problems treated in this paper were independently explored by Glymour, Spirtes and Schienes [Spirtes et al., 1989, Spirtes and Glymour, 1991], and we thank them for sharing this information with us. Discussions and correspondence with P. Bentler, D. Geiger, C. Granger, M. Hanssens, J. de Leeuw, S. Lloyd, R. Otte, A. Paz, B. Skyrms and P. Suppes are greatly appreciated.

References

[Bobrow, 1985] Bobrow, D. (1985). *Qualitative Rea-*

¹³See also Cartwright [Cartwright, 1989] for a similar position, and for a survey of the literature.

- soning about Physical Systems. MIT Press, Cambridge, MA.
- [Cartwright, 1989] Cartwright, N. (1989). *Nature Capacities and Their Measurements*. Clarendon Press, Oxford.
- [Cliff, 1983] Cliff, N. (1983). Some cautions concerning the application of causal modeling methods. *Multivariate behavioral research*, 18:115 - 126.
- [de Kleer and Brown, 1986] de Kleer, J. and Brown, J. S. (1986). Theories of causal ordering. *Artificial Intelligence*, 29(1):33 - 62.
- [Dechter and Pearl, 1990] Dechter, R. and Pearl, J. (1990). Directional constraint networks: A relational framework for causal modeling. Technical Report R-153, UCLA Cognitive Systems Laboratory.
- [Eells and Sober, 1983] Eells, E. and Sober, E. (1983). Probabilistic causality. *Philosophy of Science*, 50:35 - 57.
- [Fisher, 1953] Fisher, R. A. (1953). *Design of Experiments*. Oliver and Boyd, London.
- [Forbus and Gentner, 1986] Forbus, K. D. and Gentner, D. (1986). Causal reasoning about quantities. *Proceedings Cognitive Science Society*, pages 196 - 207.
- [Freedman, 1987] Freedman, D. (1987). As others see us: A case study in path analysis (with discussion). *Journal of Educational Statistics*, 12:101 - 223.
- [Frydenberg, 1989] Frydenberg, M. (1989). The chain graph markov property. Technical Report 186, Department of Theoretical Statistics, University of Aarhus, Denmark.
- [Gardenfors, 1988] Gardenfors, P. (1988). Causation and the dynamics of belief. In Harper, W. and Skyrms, B., editors, *Causation in Decision, Belief Change and Statistics II*, pages 85 - 104. Kluwer Academic Publishers.
- [Geffner, 1989] Geffner, H. (1989). *Default Reasoning: Causal and Conditional Theories*. PhD thesis, UCLA Computer Science Department, Los Angeles, CA.
- [Geiger et al., 1990] Geiger, D., Paz, A., and Pearl, J. (1990). Learning causal trees from dependence information. In *Proceedings, AAAI-90*, pages 770 - 776, Boston, MA.
- [Glymour et al., 1987] Glymour, C., Scheines, R., Spirtes, P., and Kelly, K. (1987). *Discovering Causal Structure*. Academic Press, New York.
- [Good, 1983] Good, I. J. (1983). A causal calculus. *British Journal for Philosophy of Science*, 11 and 12 and 13:305 - 328 and 43 - 51 and 88. reprinted as Ch. 21 in *Good Thinking* University of Minnesota Press, Minneapolis, MN.
- [Granger, 1988] Granger, C. W. J. (1988). Causality testing in a decision science. In Harper, W. and Skyrms, B., editors, *Causation in Decision, Belief Change and Statistics I*, pages 1 - 20. Kluwer Academic Publishers.
- [Holland, 1986] Holland, P. (1986). Statistics and causal inference. *Journal of the American Statistical Association*, 81:945 - 960.
- [Iwasaki and Simon, 1986] Iwasaki, Y. and Simon, H. A. (1986). Causality in device behavior. *Artificial Intelligence*, 29(1):3 - 32.
- [Kautz, 1987] Kautz, H. (1987). *A formal Theory of Plan Recognition*. PhD thesis. University of Rochester, Rochester, N.Y.
- [Kenny, 1979] Kenny, D. A. (1979). *Correlation and Causality*. Wiley, New York.
- [Keynes, 1939] Keynes, J. M. (1939). Professor tinbergen's method. *Economic Journal*, 49:560.
- [Lifschitz, 1987] Lifschitz, V. (1987). Formal theories of action. In *Workshop of the Frame Problem in AI*, pages 35 - 57, Kansas.
- [Ling, 1983] Ling, R. (1983). Review of "Correlation and Causation" by D. Kenny. *Journal of the American Statistical Association*, pages 489 - 491.
- [Mitchell, 1982] Mitchell, T. M. (1982). Generalization as search. *Artificial Intelligence*, 18:203 - 226.
- [Niles, 1922] Niles, H. E. (1922). Correlation, causation, and Wright theory of "path coefficients". *Genetics*, 7:258 - 273.
- [Patil et al., 1982] Patil, R. S., Szolovitz, P., and Schwartz, W. B. (1982). Causal understanding of patient illness in patient diagnosis. In *Proceedings of AAAI-82*, pages 345 - 348.
- [Pearl, 1978] Pearl, J. (1978). On the connection between the complexity and credibility of inferred models. *International Journal of General Systems*, 4:255 - 264.
- [Pearl, 1988a] Pearl, J. (1988a). Embracing causality in formal reasoning. *Artificial Intelligence*, 35(2):259 - 71.
- [Pearl, 1988b] Pearl, J. (1988b). *Probabilistic Reasoning in Intelligent Systems*. Morgan-Kaufman, San Mateo, CA.
- [Pearl, 1990] Pearl, J. (1990). Probabilistic and qualitative abduction. In *Proceedings of AAAI Spring Symposium on Abduction*, pages 155 - 158. Stanford.
- [Pearl et al., 1989] Pearl, J., Geiger, D., and Verma, T. S. (1989). The logic of influence diagrams. In Oliver, R. M. and Smith, J. Q., editors, *Influence Diagrams, Belief Networks and Decision Analysis*, pages 67 - 87. John Wiley and Sons, Ltd., Sussex, England.
- [Popper, 1959] Popper, K. R. (1959). *The Logic of Scientific Discovery*. Basic Books, New York.

- [Reichenbach, 1956] Reichenbach, H. (1956). *The Direction of Time*. University of California Press, Berkeley.
- [Reiter, 1987] Reiter, R. (1987). A theory of diagnosis from first principles. *Artificial Intelligence*, 32(1):57 - 95.
- [Rubin, 1989] Rubin, H. (1989). Discussion of "The Logic of Influence Diagrams" by Pearl et al. In Oliver, R. M. and Smith, J. Q., editors, *Influence Diagrams, Belief Networks and Decision Analysis*, pages 83 - 85. John Wiley and Sons, Ltd., Sussex, England.
- [Salmon, 1984] Salmon, W. (1984). *Scientific explanation and the causal structure of the world*. Princeton University Press., Princeton.
- [Shoham, 1988] Shoham, Y. (1988). *Reasoning About Change*. MIT Press, Boston, MA.
- [Simon, 1954] Simon, H. (1954). Spurious correlations: A causal interpretation. *Journal American Statistical Association*, 49:469 - 492.
- [Skyrms, 1986] Skyrms, B. (1986). *Causal Necessity*. Yale University Press, New Haven, CT.
- [Spirtes and Glymour, 1991] Spirtes, P. and Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9.
- [Spirtes et al., 1989] Spirtes, P., Glymour, C., and Scheines, R. (1989). Causality from probability. Technical Report CMU-LCL-89-4, Department of Philosophy Carnegie-Mellon University.
- [Spohn, 1983] Spohn, W. (1983). Deterministic and probabilistic reasons and causes. *Erkenntnis*, 19:371 - 396.
- [Suppes, 1970] Suppes, P. (1970). *A Probabilistic Theory of Causation*. North Holland, Amsterdam.
- [Suppes and Zaniotti, 1981] Suppes, P. and Zaniotti, M. (1981). When are probabilistic explanations possible? *Synthese*, 48:191 - 199.
- [Verma, 1991] Verma, T. S. (1991). Invariant properties of causal models. Technical report, UCLA Cognitive Systems Laboratory.
- [Verma and Pearl, 1990] Verma, T. S. and Pearl, J. (1990). Equivalence and synthesis of causal models. In *Proceedings 6th Conference on Uncertainty in AI*, pages 220 - 227.
- [Wilensky, 1983] Wilensky, R. (1983). *Planning and understanding*. Addison Wesley.
- [Wright, 1925] Wright, S. (1925). Corn and hog correlations. Technical Report 1300, U.S. Department of Agriculture.

**DIRECTED CONSTRAINT NETWORKS:
A RELATIONAL FRAMEWORK FOR CAUSAL MODELING***

(To appear in Proceedings, IJCAI'91, Sydney, Australia, August 24-30, 1991.)

Rina Dechter
Information & Computer Science
University of California, Irvine
Irvine, CA. 92717
U.S.A.

Judea Pearl
Computer Science Department
University of California, Los Angeles
Los Angeles, CA. 90024
U.S.A.

Abstract

Normally, constraint networks are undirected, since constraints merely tell us which sets of values are compatible, and compatibility is a symmetrical relationship. In contrast, causal models use directed links, conveying cause-effect asymmetries. In this paper we give a relational semantics to this directionality, thus explaining why prediction is easy while diagnosis and planning are hard. We use this semantics to show that certain relations possess intrinsic directionalities, similar to those characterizing causal influences. We also use this semantics to decide when and how an unstructured set of symmetrical constraints can be configured so as to form a directed causal theory.

1. Introduction

Finding a solution to an arbitrary set of constraints is known to be an NP-hard problem. Yet certain types of constraint systems, usually those describing causal mechanisms, manage to escape this limitation and permit us to construct a solution in an extremely efficient way. Consider, for example, the task of computing the output of an acyclic circuit consisting of a large number of logical gates. In theory, each gate is merely a constraint that forbids certain input-output combinations from occurring, and the task of computing the output of the overall circuit (for a given combination of the circuit inputs) is equivalent to that of finding a solution to a set of constraints. Yet contrary to the general constraint problem, this task is remarkably simple; one need only trace the flow of causation and propagate the values of the intermediate variables from the circuit inputs down to the circuit output(s). This forward computation encounters none of the difficulties of the general constraint-satisfaction problems, thus exemplifying the simplicity inherent to causal predictions.

The aim of this paper is to identify and characterize the features that render this class of problems computationally efficient, thus explaining some of the reasons that causal models are so popular in the organization of human knowledge. Note that this efficiency is asymmetric; it only characterizes the forward computation, but fails to hold in the backward direction. For instance, the problem of finding an input combination that yields a given output (a task we normally associate with planning or diagnosis) is as hard as any constraint satisfaction problem. Thus, the second aim of our analysis is to explain how a system of constraints, each defined in terms of the totally symmetric relationship of compatibility, can give rise to such profound asymmetries as those attributed to cause-effect or input-output relationships. At first glance, we might be tempted to attribute the asymmetry to the functional nature of the constraints involved. However, functional dependency in itself cannot explain the directional asymmetry found in the analysis of causal mechanisms such as the logic circuit above. Imagine a circuit containing some faulty components, the output of which may attain one of several values. The constraints are no longer functional, yet the asymmetry persists; finding an output compatible with a given input is easy while finding an input compatible with a given output is hard. This asymmetry between prediction and planning seems to be a universal feature of all systems involving causal mechanisms [Shoham, 1988], a feature we must emulate in defining causal theories.

Our starting point is to formulate a necessary and sufficient condition for a system of constraints to exhibit a directional asymmetry similar to that characterizing causal organizations. Basically, the criterion is that of modularity: there should exist an ordering of the variables in the system such that imposing constraints on later variables would not further constrain earlier variables. Intuitively, it captures the understanding that predictions are useless for diagnosis; e.g., given a set of findings, we cannot improve the accuracy of our diagnosis by concentrating our analysis on the patient's prospects for recovery. Likewise, in the context of the logic circuit example, modularity asserts that if we wish to add a new gate, then, as long as we do not connect to its output, we can add this gate anywhere in the circuit without

*This work was partially supported by the National Science Foundation, Grant #IRI-8821444 and by the Air Force Office of Scientific Research, Grant #AFOSR-90-0136.

perturbing the circuit's behavior. Starting with modularity as a definition of causal theories (Section 2), we show¹ that it is tantamount to enabling backtrack-free search (for a feasible solution) along any natural ordering of the theory. We then explore methods of constructing causal specifications for a given relation, that is, specifications that permit objects from the relation to be retrieved backtrack-free along some ordering. Such methods are investigated along two dimensions: inductive and pragmatic. Along the inductive dimension (Section 3), we observe the tuples of some relation ρ , and we seek to represent this set of observations by a causal theory that is as *simple* as possible. We provide a formal definition of simplicity and show that together with the insistence on backtrack-free predictions, it leads to a natural definition of *intrinsic directionality*, matching our perception of causal directionality in logical circuits and other physical devices.

Along the pragmatic dimension (Section 4), we start with an unordered collection of constraint specifications, which might represent some stable physical laws, and we seek an ordering of the variables such that the overall system constitutes a causal theory. Clearly, not every system of constraints can turn causal by a clever ordering of the variables. The criterion for the existence of such an ordering depends on both the nature of the constraints and the topology of the subsets of variables upon which the constraints are specified. Some constraint systems are amiable to causal ordering by virtue of their topology alone, *regardless* of the content of the individual constraints. These are called acyclic constraint systems, originally studied in the literature of relational databases, [Beeri et al., 1983]. In contrast, Section 4 ascribes causal ordering to a more general set of topologies, but imposes special requirements on the character of the individual constraints.

Our basic requirement for a k -variable constraint to qualify as a description of a primitive causal mechanism, is that at least one set of $k-1$ variables must behave as **inputs** (or **causes**) relative to the remaining k^{th} variable (to be regarded as an **output** or an **effect**), that is, no value combination of these $k-1$ variables can be forbidden, and each such combination must be compatible with at least one value of the k^{th} variable. Additionally, in order for the system as a whole to act as a causal system, mechanisms must be ordered in a way that prevents conflicts among their predictions, hence, we require that no two constraints should designate the same variable as an output. We provide effective procedures for: (1) deciding if such an ordering exists and, (2) identifying such ordering whenever possible. The ordering found can be used to facilitate search and retrieval, and are similar to those used to describe the operation of physical devices [Kuipers, 1984; Iwasaki and Simon, 1986; de-Kleer and Brown, 1986].

2. Definitions and Preliminaries: Constraint

¹Proofs can be found in [Dechter and Pearl, 1991].

Specifications and Causal Theories

Definition 1 (Constraint Specification): A constraint specification (CS) consists of a set of n variables $X = \{X_1, \dots, X_n\}$, each associated with a finite domain, dom_1, \dots, dom_n , and a set of constraints $\{C_1, C_2, \dots, C_r\}$ on subsets of X . Each constraint C_i is a relation on a subset of variables $S_i = \{X_{i_1}, \dots, X_{i_j}\}$, namely, it defines a subset of the Cartesian product of $dom_{i_1} \times \dots \times dom_{i_j}$. The **scheme** of a CS is the set of subsets on which constraints are defined, $scheme(CS) = \{S_1, S_2, \dots, S_r\}$, $S_i \subseteq X$, and each such subset is called a **component**. A **solution** of a given CS is an assignment of values to the variables in X such that all the constraints in the CS are satisfied. A constraint specification CS is said to define an **underlying relation** $rel(CS)$, consisting of all the solutions of CS.

Definition 2 (Causal Theories): Given a constraint specification CS, its underlying relation $\rho = rel(CS)$, and an ordering $d = (X_1, X_2, \dots, X_n)$, we say that a CS is a **causal theory** (of ρ) relative to d if for all $i \geq 1$ we have

$$\Pi_{X_1, \dots, X_i}(\rho) = \bowtie_{j(i)} C_j \quad (1)$$

where

$$j(i) = \{j: S_j \subseteq \{X_1, \dots, X_i\}\}. \quad (2)$$

$\Pi_{X_1, \dots, X_i}(\rho)$ denotes the projection of ρ on $\{X_1, \dots, X_i\}$, that is, the set of all subtuples (x_1, \dots, x_i) for which an extension $(x_1, \dots, x_i, x_{i+1}, \dots, x_n)$ exists in ρ , and \bowtie is the *join* operator. Any pair $\langle d, CS \rangle$ satisfying (1) will be called a **causal theory** (of ρ).

Although condition (1) may seem hard to verify in practice, it nevertheless provides an operational definition for causal theories. To test whether a given CS is causal relative to ordering d , we need to find the set of solutions to the given CS, project back these solutions on the strings of variables X_1, X_2, \dots, X_i , $1 \leq i \leq n$, then check whether each such projection coincides exactly with the set of solutions to a smaller CS, one consisting of only those constraints that are defined on variables taken from $\{X_1, \dots, X_i\}$. In Section 4 we will show that certain types of specifications possess syntactic features that render them inherently causal, in no need of the elaborate test prescribed by (1). For example, the specifications provided by a collection of logic gates always constitutes a causal theory relative to any ordering compatible with their standard assembly in acyclic circuits (i.e., no variable can serve as an output of two different gates). Similarly, linear inequalities and propositional clauses, under certain conditions, can be assembled into causal theories by finding appropriate orderings of the variables.

From a conceptual viewpoint, Definition 2 formalizes the notion of modularity (see Introduction) and can be given the following temporal interpretation. If we view the variables X_1, \dots, X_i as past events, the variables X_{i+1}, \dots, X_n as future events, and the constraints as physical laws, then Eq. (1) asserts that the permissible set of past scenarios is not affected by laws that pertain only to future events. In other words, the set of scenarios we get by ignoring future constraints will remain valid after including such constraints in the analysis. This interpretation is indeed at the very heart of the notion of causation, and is closely related to the principle of *chronological ignorance* described in [Shoham, 1988], although Shoham's definition of causal theories insists on functional dependencies.

We shall now show that causal theories as defined by (1) yield a computationally effective scheme of encoding relations; it guarantees that the tuples of these relations can be generated systematically, without search, by simply instantiating variables along the natural ordering of the theory.

Definition 3 (Backtrack-free): We say that a CS is *backtrack-free* along ordering $d = (X_1, \dots, X_n)$ if for every i and for every assignment x_1, \dots, x_i consistent with $\{C_j: S_j \subseteq \{X_1, \dots, X_i\}\}$ there is a value x_{i+1} of X_{i+1} such that x_1, \dots, x_i, x_{i+1} satisfies all the constraints in $\{C_j: S_j \subseteq \{X_1, \dots, X_{i+1}\}\}$. In other words, a CS is backtrack-free w.r.t. d if $rel(CS)$ can be recovered with no dead-ends along the order d .

Definition 3 is an extension of the standard notion of backtrack-free originally stated for binary constraints [Freuder, 1982], and later related to directional consistency [Dechter, 1990]. Note that, given a constraint C_i on a subset S_j of variables, definition 3 does not allow testing whether some partial instantiation of S_j is compatible with C_j . It is possible to weaken this restriction by considering all the constraints projections as part of the problem's scheme. In this paper we do not consider such projections; nevertheless, our analysis is extensible to that case as well.

Theorem 1: A constraint specification CS is backtrack-free along an ordering d if and only if it is causal relative to d .

In the practice of causal modeling, it is common to depict the structure of causal theories using directed acyclic graphs (dags), not total orders. Each such dag, called a causal model, indicates the existence of direct causal influences among sets of variables, but does not specify the precise nature of the influences. We will next give a formal definition of such models, and then explore what properties of the underlying relation are portrayed by the topology of the dag.

Definition 4 (Dags and Families): Given a directed acyclic graph (dag) D , we say that an ordering $d = (X_1, \dots, X_n)$

of the nodes in the graph respects D if all edges in D are directed from lower to higher nodes of d . A dag D defines a set of n families F_1, \dots, F_n , each family F_i is a subset consisting of a son node, X_i , and all its parent nodes, P_i , which are those directed towards X_i in D .

Definition 4' (Characteristic dag): The characteristic dag, D , of the pair (d, CS) is constructed as follows: For each component S_j in $scheme(CS)$, designate the latest variable (according to d) in S_j as a sink and direct the other variables in S_j towards it.

Figure 1 shows the characteristic dag of a CS defined on the subsets AB, AC, BD, CD, CE, DEF , along the ordering $d = (A, B, C, D, E, F)$.

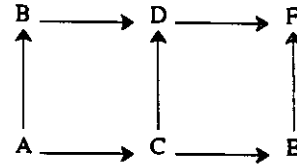


Figure 1: The characteristic dag of a CS

Lemma 1: If D is the characteristic dag of the pair (d, CS) then it is also the characteristic dag of (d', CS) whenever d' respects D and, furthermore, if $\langle d, CS \rangle$ is a causal theory, then so is $\langle d', CS \rangle$. \square

We now define causal theories and models using dags:

Definition 5: A pair $\langle D, CS \rangle$ is a causal theory if $\langle d, CS \rangle$ is a causal theory for all d respecting D .

Definition 6 (Causal model): Given a relation ρ and an arbitrary dag D , D is a **causal model** of ρ if there exists a constraint specification CS such that $\langle D, CS \rangle$ is causal theory of ρ .

It is easy to see that not every dag D could be a causal model of a given relation ρ . For example, the relation defined by the pair of logical clauses $\{X \vee Z, Y \vee Z\}$ can be modeled by either $X \rightarrow Z \leftarrow Y$ or $X \leftarrow Z \rightarrow Y$, but not by $X \rightarrow Y \leftarrow Z$. The reason is that while the former two dags form causal theories with the specification above, no such theory can be formed for the third dag, because to determine the permissible values of X and Z we must consult a later variable, Y .

To determine whether a dag D is a causal model of a given relation ρ , one need not enumerate the space of specifications for ρ . The condition is simply that D should decompose ρ by the following rule: Fixing its parents in D , each variable must remain unaffected by all other variables, except possibly by its descendants in D . This rule reflects another common feature of causation: once we learn the current status of its direct causal factors, no other information is needed for predicting the state of a given variable.

Formally, D is a causal model of ρ if for some ordering X_1, \dots, X_n respecting D and for all i , we have $\Pi_{X_1, \dots, X_{i-1}}(\rho) \bowtie \Pi_{P_i \cup \{X_i\}}(\rho) = \Pi_{X_1, \dots, X_i}(\rho)$ where P_i stands for the parents of X_i . This result follows from the theory of graphoids, as applied to database dependencies [Pearl, 1988]. The complexity of the test above is polynomial in the size of ρ , but may be exponential in the number of variables. Once D qualified as a causal model of ρ , a causal theory $\langle D, CS \rangle$ (of ρ) can be formed by simply pairing D with the projections of ρ on the families of D .

3. Synthesizing Causal Theories and Uncovering Causal Directionality

Our ultimate goal is to construct causal theories for the information we possess. In this section we analyze two tasks. First, we assume that the information we have is a database tabulating explicitly the tuples of some relation ρ , and our task is to replace the table by a more economical representation, one that enjoys the computational advantage of causal organizations. Such a task would be useful in machine learning applications, where the tuples represent a stream of observations and the causal theory forms a convenient model of the environment, facilitating modular organization and fast predictions. In our second task, the information will be given in the form of a preformulated constraint specification CS , and the problem will be to construct a causal theory without explicating the underlying relation of CS .

Task 1: (decomposition) Given a relation ρ and an ordering d , find a causal theory for ρ along d .

Barring additional requirements, a causal theory can be obtained by a trivial construction. For instance, the complete dag generated by directing an edge from each lower variable to every higher variable is clearly a causal model of ρ , and the desired causal theory can be obtained by projecting ρ onto the complete families $F_i = \{X_1, X_2, \dots, X_i\}$. We next present a scheme for constructing a causal theory on top of an edge-minimal model of ρ , that is, a dag D from which no edge can be deleted without destroying its capability to support a causal theory of ρ .

The algorithm that follows constructs an edge-minimal causal model of ρ .

build-causal-1 (ρ, d):

1. Begin
2. For $i = n$ to 2 by -1 do:
3. Find a minimal subset $P_i \subseteq \{X_1, \dots, X_{i-1}\}$ such that

$$\Pi_{X_1, \dots, X_{i-1}}(\rho) \bowtie \Pi_{P_i \cup \{X_i\}}(\rho) = \Pi_{X_1, \dots, X_i}(\rho)$$
4. Return a dag D generated by directing an arc from each node in P_i towards X_i .
5. End.

To form a causal theory, we simply pair this dag with the projections of ρ on its families.

The construction above shows that a causal theory can be found for any arbitrary ordering. However, we will next show that certain orderings possess features that render them more natural for a given relation. It is these features, we conjecture, which give rise to the perception that certain relations possess "intrinsic" directionalities.

Definition 7 (Model Preference): A causal model D_2 is said to be at least as expressive as D_1 , denoted $D_1 \leq D_2$, if for any causal theory $\langle D_1, CS_1 \rangle$ there exists a causal theory $\langle D_2, CS_2 \rangle$ such that $rel(CS_1) = rel(CS_2)$. A dag D is said to be a **minimal** causal model of ρ if it is not strictly more expressive than any other causal model of ρ . In other words, the set of relations modeled by D is not a superset of any set of relations, containing ρ , that can be modeled by some other dag.

Clearly, every minimal model must be edge-minimal, but not the converse. For example, the complete dag $Z \rightarrow X, Z \rightarrow Y, X \rightarrow Y$ is an edge-minimal causal model of the relation given by the formula $Z = X \vee Y$, but it is not a minimal model, because it is strictly more expressive than the dag $X \rightarrow Z \leftarrow Y$; the latter can model only relations where X does not constrain Y . Polynomial graphical methods for testing preference and equivalence between causal models are described in [Pearl et al., 1990]. However, finding a minimal model for a given relation may be exponentially hard.

Definition 8 (Intrinsic Directionality): Given a relation ρ , a variable X is said to be a **direct cause** of variable Y , if there exists a directed edge from X to Y in all minimal causal models of ρ .

Example 1. Consider a relation ρ specified by the table of Figure 2(a). The table is small enough to verify that the dag in 2(b) is the only minimal causal model of ρ . For example, the arrow from X to Z cannot be reversed, because ρ cannot be expressed as a set of constraints on the families of the resulting dag, $\{YZ, ZX, XYW\}$. Adding an arc $Y \rightarrow X$ to the resulting dag would permit a representation of ρ (using the scheme $\{YZ, YZX, YXW\}$), but would no longer be minimal. It is strictly more expressive than the one in 2(b), because, unlike the latter, it also models relations in which some XY pairs are forbidden. The causal theory corresponding to the dag of 2(b) is shown in 2(c), matching our intuition about the causal relationships embedded in 2(a). Note that the same minimal model ensues (though not the same theory) were we to destroy the functional dependencies by adding the tuple 1100 to the table in 2(a). However, it is no longer unique.

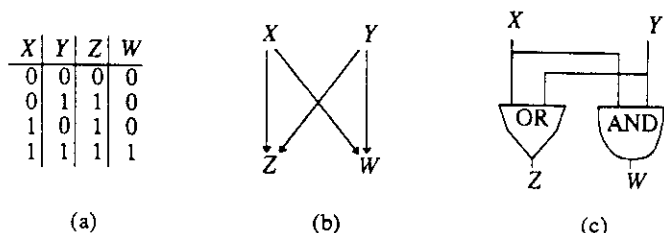


Figure 2: The directionality shown in (b) is intrinsic to the relation in (a), because (b) is a unique minimal causal model of (a).

Verma and Pearl [1990] have used minimal model semantics to construct a probabilistic definition of causal directionality. They have also developed a proof theory which, under certain conditions provides efficient algorithms for determining causal directionality without examining the vast space of minimal models [Pearl and Verma, 1991]. Whether similar conditions exist in the relational framework remains an open problem.

5. Conclusions

This paper presents a relational semantics for the directionality associated with cause-effect relationships, explaining why prediction is easy while diagnosis and planning are hard. We used this semantics to show that certain relations possess intrinsic directionalities, similar to those characterizing causal influences. We also provided an effective procedure for deciding when and how an unstructured set of constraints can be configured so as to form a directed causal theory.

These results have several applications. First, it is often more natural for a person to express causal relationships as directional, rather than symmetrical constraints. The semantics presented in this paper permits us to interpret and process directional relationships in a consistent way and to utilize the computational advantages latent in causal theories. Second, the notion of intrinsic directionality suggests automated procedures for discovering causal structures in raw observations or, at the very least, for organizing such observations into structures that enjoy the characteristics of causal theories. Finally, the set of constraint specifications that can be configured to form causal theories constitutes another "island of tractability" in constraint satisfaction problems. The procedure provided for identifying such specifications can be used to order computational sequences in qualitative physics and scheduling applications.

References

[Beeri et al., 1983] Beeri, C., Fagin, R., Maier, D., and Yannakakis, M., "On the desirability of acyclic database schemes," *Journal of ACM*, Vol. 30, No. 2, July, 1983, pp. 479-513

[de-Kleer and Brown, 1986] de-Kleer, J. and Brown, J.S., "Theories of causal ordering," *Artificial Intelligence*, Vol. 29, No. 1, 1986, pp. 33-62.

[Dechter, 1990] Dechter, R., "From local to global consistency," In *Proceedings, Eighth Canadian Conference on AI, CSCSI-90*, May 23-25, 1990, pp. 231-237.

[Dechter and Pearl, 1991] Dechter, R. and Pearl, J., "Directed constraint networks," Cognitive Systems Laboratory, University of California, Los Angeles. Technical Report R-153-L, 1991.

[Freuder, 1982] Freuder, E.C., "A sufficient condition for backtrack-free search." *JACM*, 29(1), 1982, pp. 24-32.

[Iwasaki and Simon, 1986] Iwasaki, Y. and Simon, H.A., "Causality in device behavior," *Artificial Intelligence*, Vol. 29, No. 1, 1986, pp. 3-32.

[Kuipers, 1984] Kuipers, B. "Common sense reasoning about causality: Deriving behavior from structure," *Artificial Intelligence*, Vol. 24, No. 1-3, 1984, pp. 169-203.

[Pearl, 1988] Pearl, J. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Palo Alto, CA., Morgan Kaufmann, 1988.

[Pearl et al., 1990] Pearl, J., Geiger, D., and Verma, T., "The logic of influence diagrams," in R.M. Oliver and J.Q. Smith (Eds), *Influence Diagrams, Belief Nets and Decision Analysis*, Sussex, England: John Wiley & Sons, Ltd., 1989, pp. 67-87.

[Pearl and Verma, 1991] Pearl, J. and Verma T., "A theory of inferred causation," in Allen, J.A., Fikes, R., and Sandewall, E. (Eds.), *Principles of Knowledge Representation and Reasoning: Proceedings of the Second International Conference*. San Mateo, CA: Morgan Kaufmann, (April, 1991), pp. 441-452.

[Shoham, 1988] Shoham, Y., *Reasoning About Change*, Cambridge, Massachusetts: The MIT Press, 1988.

[Verma and Pearl, 1990] Verma, T. and Pearl, J., "Equivalence and synthesis of causal models," in *Proceedings, Sixth Conference on Uncertainty in AI*, Cambridge, Massachusetts, July 27-29, 1990, pp. 220-227.

