

Subjective Causality

Joseph Y. Halpern¹, Evan Piermont²

¹Computer Science Department, Cornell University, Ithaca, USA.

²Department of Economics, Royal Holloway, University of London, UK.

halpern@cs.cornell.edu, evan.piermont@rhul.ac.uk

Abstract

We show that it is possible to understand and identify a decision maker's subjective causal judgements by observing her preferences over interventions. Following Pearl [2000], we represent causality using *causal models* (also called *structural equations models*), where the world is described by a collection of variables, related by equations. We show that if a preference relation over interventions satisfies certain axioms (related to standard axioms regarding counterfactuals), then we can define (i) a causal model, (ii) a probability capturing the decision-maker's uncertainty regarding the external factors in the world and (iii) a utility on outcomes such that each intervention is associated with an expected utility and such that intervention A is preferred to B iff the expected utility of A is greater than that of B . In addition, we characterize when the causal model is unique. Thus, our results allow a modeler to test the hypothesis that a decision maker's preferences are consistent with some causal model and to identify causal judgements from observed behavior.

1 Introduction

Causal judgments play an important role in decision making. When deciding between actions that intervene directly on some aspect of the world, one major source of uncertainty is the indirect effect of such actions via causal interaction. For example, when deciding the interest rate, the Federal Reserve might consider the possibility that a change in the interest rate will cause a change in unemployment, and further that this causal relationship itself might be contingent on other macroeconomic variables.

Uncovering and describing the causal relationship between variables is a task that has led to enormous effort across many different disciplines (see, e.g., [Angrist and Pischke, 2009; Cunningham, 2021; Hernán and Robins, 2020; Morgan and Winship, 2007; Parascandola and Weed, 2001; Plowright *et al.*, 2008; Pearl, 2009; Pearl, 2000; Spirtes *et al.*, 1993]; this list barely scratches the surface.) Different decision makers, on account of their private information and personal experience, might hold different

beliefs about causal relationships. In this paper, we show that it is possible to understand and identify a decision maker's subjective causal judgements by observing her preferences over interventions.

A first step to doing this is to decide how to represent causality. Most recent work has focused on using counterfactuals. In the philosophy community, following Stalnaker [1968] and Lewis [1973], counterfactuals are given semantics using possible worlds equipped with a “closer than” relation; a counterfactual such as “If ϕ were true then ψ would be true” is true at a world ω if, ψ is true at the closest world(s) to ω where ϕ is true. Pearl [2000] has championed the use of *causal models* (also called *structural equations models*, graphical models where the world is described by a collection of variables, related by equations. (These are related to models of causality in economics that go back to the work of Haavelmo [1943] and Simon [1953].) The equations model the effect of counterfactual interventions.

We use the latter approach here, although as we shall show, there are close relationships with the former approach as well. Following Pearl, we assume that the world is described by a set of variables. It is useful to split them into two sets: the *exogenous* variables, whose values are determined by factors outside the model, and the *endogenous* variables, whose values are determined by the exogenous variables and other endogenous variables. Which variables should be taken as exogenous and which should be taken as endogenous depends on the situation. For example, if the Federal Reserve is deciding whether to change the interest rates, the interest rates should clearly be viewed as endogenous. But for a company trying to decide whether to go ahead with a project that involves borrowing money at the current interest rates, the interest rates are perhaps best viewed as exogenous.

A primitive action is an intervention that sets the value of a particular variable; because the values of the exogenous variables are taken as given, we assume that only on endogenous variables can be intervened on. Following Pearl [2000], we use the **do** notation to denote such actions. For example $\text{do}[Y \leftarrow y]$ is the primitive action that sets variable Y to value y . Following Blume, Easley, and Halpern [2021] (BEH from now on), we allow more general actions to be formed from these primitive actions using **if ... then ... else**. Thus, non-primitive actions are conditional interventions of

the form

$$\text{if } \phi \text{ then } A \text{ else } B$$

, where A and B are themselves (possibly primitive) actions and ϕ is a test—a statement regarding the values of variables that is either true or false.

As is standard in decision-theoretic analyses, we assume that the decision maker has a preference relation \succsim over actions. We show that if the preference relation satisfies certain axioms (that can be understood as corresponding to standard axioms regarding counterfactuals), then we can represent the decision-maker's preference as the maximization of the expected utility of an actions relative to a causal model equipped with a probability and utility. Specifically, given a preference relation, we can define a causal model, add a probability on *contexts* (settings of the exogenous variables, which can be viewed as capturing the decision-maker's uncertainty regarding the external factors in the world), and a utility on outcomes (which are just settings of the variables in the model). In such a causal model, we can determine the expected utility of an action that involves (conditional) interventions. We show that the decision maker prefers action A to B iff the expected utility of A is greater than that of B , given the causal model, probability, and utility.

Our results allow a modeler to test the hypothesis that a decision maker's preferences are consistent with some causal model. In doing so, we provide a definition of *causally sophisticated* behavior as a benchmark of rationality for decision making in the presence of interventions.¹ When a decision maker is causally sophisticated—when her actions can be rationalized by some causal model—our results further determine when the causal model can be uniquely identified. This result provides a modeler the means to explain observed economic behavior in terms of causal judgments.

To the best of our knowledge, we are the first to examine causal decision making in the context of the structural equations, which is perhaps now the most common approach to representing causality in the social sciences and computer science. Bjorndahl and Halpern [2021] (BH from now on) addressed similar questions in the context of the closest-world approach to counterfactuals of Lewis [1973] and Stalnaker [1968]. Interestingly, our technical results use their results (and earlier results of BEH) as the basis for our representation theorem, based on a connection between the two representations of counterfactuals due to Halpern [2013]. Our results allow us to relate the two approaches at a decision-theoretic level.

Other approaches to modeling causality have also been considered in the literature. Schenone [2020] and Ellis and Thysenb [2021] take a statistical approach, taking a lack of conditional independence as a definition of causality; Alexander and Gilboa [2023] understand causality through a reduction in the Kolmogorov complexity. As we said, our approach is closer to the approach that is currently the focus of work in the social sciences and computer science. Importantly, by identifying the structural equations, we provide a more detailed insight into the causal mechanisms being considered

¹This is somewhat analogous to *probabilistic sophistication* in decision making [Machina and Schmeidler, 1992].

by the decision maker.

In our representation (like that of BH), both utility and probability are defined on valuations of variables. So, in contrast to the decision-theoretic tradition following Savage (and, of particular relevance, [Schenone, 2020]), there is no separation between states (on which probability is defined) and outcomes (on which utilities are defined): our decision maker derives utility directly from the outcome of the intervention. This permits the application of our model to the many economically relevant contexts where the effect of the intervention has direct utilitarian consequences for the decision maker (e.g., a government setting policy, or a firm deciding an investment strategy).

2 Causal Models

Let \mathcal{U} and \mathcal{V} denote the set of exogenous and endogenous variables, respectively. Let $\mathcal{R} : \mathcal{U} \cup \mathcal{V} \rightarrow 2^{\mathbb{R}}$ associate to each $Y \in \mathcal{U} \cup \mathcal{V}$ a set $\mathcal{R}(Y)$ of possible values called its range. We assume that the range for each variable is finite. Let $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathcal{R})$ denote a *signature*.

A *causal model* is a pair $M = (\mathcal{S}, \mathcal{F})$, where \mathcal{S} is a signature and \mathcal{F} is a collection of structural equations that determine the values of endogenous variables based on the values of other variables. Formally, $\mathcal{F} = \{F_X\}_{X \in \mathcal{V}}$, where

$$F_X : \prod_{Y \in \mathcal{U} \cup (\mathcal{V} - \{X\})} \mathcal{R}(Y) \rightarrow \mathcal{R}(X).$$

A model M is *recursive* if there exists a partial order on \mathcal{V} such that the structural equations for each variable is independent of the variables lower in the order. When we say that X is independent of Y , we mean that F_X does not depend on the value of Y . Given a signature \mathcal{S} , let $\mathcal{M}(\mathcal{S})$ denote the set of all models over the signature \mathcal{S} and $\mathcal{M}^{rec}(\mathcal{S})$ the set of recursive models.

We use the standard vector notation to denote sets of variables or their values. If $\vec{X} = (X_1, \dots, X_n)$ is a vector, we write, in a standard extension of the usual containment relation, $\vec{X} \subseteq \mathcal{U} \cup \mathcal{V}$ to indicate that each $X_i \in \mathcal{U} \cup \mathcal{V}$. Similarly, we write $\vec{x} \in \mathcal{R}(\vec{X})$ if $\vec{x} \in \prod_{i \leq n} \mathcal{R}(X_i)$. For a vector \vec{Y} of variables, $\vec{y} \in \mathcal{R}(\vec{Y})$, and a vector $\vec{X} \subseteq \vec{Y}$, let $\vec{y}|_{\vec{X}}$ denote the restriction of \vec{y} to \vec{X} . In particular, for a single variable X , $\vec{y}|_X$ denotes the X^{th} component.

A *context* is a vector \vec{u} of values for all the exogenous variables \mathcal{U} . Let $\mathcal{C}(\mathcal{S}) = \prod_{Y \in \mathcal{U}} \mathcal{R}(Y)$ consist of all contexts for the signature \mathcal{S} . Call a pair $(M, \vec{u}) \in \mathcal{M}^{rec}(\mathcal{S}) \times \mathcal{C}(\mathcal{S})$ a *situation*. Each situation has a unique solution (i.e., a unique value for each variable that simultaneously satisfies all the structural equations and agrees with the context).

Given a model $M = (\mathcal{S}, \mathcal{F})$, we can construct a new model (over the same signature) that represents the counterfactual situation where some variables are set to specific values. If $\vec{Y} \subseteq \mathcal{V}$ and $\vec{y} \in \mathcal{R}(\vec{Y})$, then $M_{do[\vec{Y} \leftarrow \vec{y}]} = (\mathcal{S}, \mathcal{F}_{do[\vec{Y} \leftarrow \vec{y}]})$ denotes the model that is identical to M except that the equation for each variable $X \in \vec{Y}$ in $\mathcal{F}_{do[\vec{Y} \leftarrow \vec{y}]}$ is replaced by $X = \vec{y}|_X$.² Note that $M_{do[\vec{Y} \leftarrow \vec{y}]}$ is recursive if M is.

²It is more standard in the literature to write $M_{\vec{Y} \leftarrow \vec{y}}$ rather than

2.1 Syntax and Semantics

Fix a signature \mathcal{S} . For each $X \in \mathcal{U} \cup \mathcal{V}$ and $x \in \mathcal{R}(Y)$, let $X = x$ denote the atomic proposition that says that the variable X takes value x . Let $\mathcal{L}(\mathcal{S})$ denote the language constructed by starting with these propositions and closing off under negation, conjunction, and disjunction.³

A formula $\phi \in \mathcal{L}(\mathcal{S})$ is either true or false in a situation $(M, \vec{u}) \in \mathcal{M}^{rec}(\mathcal{S}) \times \mathcal{C}(\mathcal{S})$. We write $(M, \vec{u}) \models \psi$ if ψ is true in the situation (M, \vec{u}) . The \models relation is defined inductively.

- $(M, \vec{u}) \models X = x$ iff $X = x$ in the unique solution to the system of equations \mathcal{F} , starting with context \vec{u} .
- $(M, \vec{u}) \models \neg\phi$ iff not $(M, \vec{u}) \models \phi$.
- $(M, \vec{u}) \models \phi \wedge \psi$ iff $(M, \vec{u}) \models \phi$ and $(M, \vec{u}) \models \psi$.
- $(M, \vec{u}) \models \phi \vee \psi$ iff $(M, \vec{u}) \models \phi$ or $(M, \vec{u}) \models \psi$.

An *atom* is a complete description of the values of variables; its truth will completely determine truth of all of formulae in $\mathcal{L}(\mathcal{S})$. Formally, an *atom over \mathcal{S}* is a conjunction of the form $\wedge_{Y \in \mathcal{U} \cup \mathcal{V}} Y = y$, where y is some value in $\mathcal{R}(Y)$. We use \vec{a} to denote a generic atom. Let $\mathcal{A}(\mathcal{S})$ denote the set of atoms over \mathcal{S} . Notice that an atom $\vec{a} \in \mathcal{A}(\mathcal{S})$ determines the truth of all formulas in $\mathcal{L}(\mathcal{S})$; we write $\vec{a} \Rightarrow \phi$ if ϕ is true in situations satisfying \vec{a} . Let $\vec{a}_{M, \vec{u}}$ denote unique atom such that $(M, \vec{u}) \models \vec{a}$.

In a slight abuse of notation, identify each context $\vec{u} \in \mathcal{C}$ with the formula characterizing \vec{u} , that is, the formula $\wedge_{U \in \mathcal{U}} U = \vec{u}|_U$.

For our later discussion, it is useful to consider an extension of $\mathcal{L}(\mathcal{S})$ that includes formulas of the form $[\vec{Y} \leftarrow \vec{y}] \phi$, where $\phi \in \mathcal{L}(\mathcal{S})$, $\vec{Y} \subseteq \mathcal{V}$, and $\vec{y} \in \mathcal{R}(\vec{Y})$. We can view this formula as saying “after intervening to set the variables in \vec{Y} to \vec{y} , the formula ϕ holds”. Call this extended language $\mathcal{L}^+(\mathcal{S})$. We can extend the semantics that we gave to formulas in $\mathcal{L}(\mathcal{S})$ as follows:

- $(M, \vec{u}) \models [\vec{Y} \leftarrow \vec{y}] \phi$ iff $(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models \phi$.

In the sequel, we will also make use of another approach to giving semantics to counterfactuals, due to David Lewis and Robert Stalnaker [Lewis, 1973; Stalnaker, 1968]. This approach is based on the idea of “closest worlds”; roughly speaking, with this approach, $[\vec{Y} \leftarrow \vec{y}] \phi$ is true at a world ω if ϕ is true at all the worlds closest to ω where $\vec{Y} = \vec{y}$. Lewis formalizes the idea of “closest world” using a ternary relation R where, for each $\omega \in \Omega$, $R(\omega, \cdot, \cdot)$ is a partial order on Ω . We can interpret $R(\omega_1, \omega_2, \omega_3)$ as saying that ω_2 is closer to ω_1 than ω_3 is. (In this paper, we focus on the case that, for each world ω , $R(\omega, \cdot, \cdot)$ is a strict linear order.)

In more detail, a *Lewis-style model* is a tuple $M = (\Omega, R, I)$, where R is a ternary relation on Ω as above, and I is an *interpretation* that determines whether each atomic

³ $M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$. We use the latter notation because it makes it easier to express some later notions.

³It is more standard to have the atomic propositions involve only endogenous variables, but for our purposes, it is important to include exogenous variables as well. We explain why when we discuss actions below.

proposition is true or false at each state. Formally, if AP is the set of atomic propositions (as implicitly determined by a set $\mathcal{U} \cup \mathcal{V}$ of exogenous and endogenous variables), then $I : \Omega \times AP \rightarrow \{\text{true, false}\}$. We can again define a relation \models by induction, by taking

- $(M, \omega) \models X = x$ iff $I(\omega, X = x) = \text{true}$,

defining the semantics of negation, conjunction, and disjunction as above, and for interventions, taking

- $(M, \omega) \models [\vec{Y} \leftarrow \vec{y}] \phi$ iff $(M, \omega') \models \phi$, for all ω' that are minimal worlds according to the order $R(\omega, \cdot, \cdot)$ such that $(M, \omega') \models \vec{Y} = \vec{y}$; that is, the worlds closest to ω for which $\vec{Y} = \vec{y}$ holds.

As shown by Halpern [2013], recursive causal models correspond in a precise sense to a subclass of Lewis-style models; we return to this point in Section 5.

3 Decision Environment

A *primitive action* over \mathcal{S} has the form $\text{do}[Y_1 \leftarrow y_1, \dots, Y_n \leftarrow y_n]$, abbreviated as $\text{do}[\vec{Y} \leftarrow \vec{y}]$, where Y_1, \dots, Y_n is a (possibly) empty list of distinct variables in \mathcal{V} and $y_i \in \mathcal{R}(Y_i)$ for $i = 1, \dots, n$. This action represents an intervention that sets each Y_i to the value y_i . As we said earlier, we allow interventions only on endogenous variables.

The set $\mathbb{A}(\mathcal{S})$ of *actions* over \mathcal{S} is defined recursively, starting with the primitive actions, and closing under **if ... then ... else**, so that if $A, B \in \mathbb{A}$ and $\phi \in \mathcal{L}(\mathcal{S})$, then

$$\text{if } \phi \text{ then } A \text{ else } B \in \mathbb{A}.$$

We take **if ϕ then A** to be an abbreviation of **if ϕ then A else do[]** (**do[]** is the trivial action that sets no values). Note that although we restrict interventions to endogenous variables, the tests can involve exogenous variables. For example, even if interest rates are fixed (and hence taken to be exogenous), we may want to say “if the interest rate is 5% then borrow \$1,000,000” (where we take the amount borrowed to be represented by an endogenous variable).

We assume that the decision maker has a preference relation (weak order) $\succsim_{\mathcal{S}}$ over the set of all actions in $\mathbb{A}(\mathcal{S})$. (We often omit the \mathcal{S} if it is clear from context or plays no role in the discussion.) As usual, we write $A \sim B$ if $A \succsim B$ and $B \succsim A$. Call a formula ϕ *null* if its conditional preference is trivial, that is, if $(\text{if } \phi \text{ then } A) \sim (\text{if } \phi \text{ then } B)$ for all actions $A, B \in \mathbb{A}(\mathcal{S})$. Call ϕ *non-null* if it is not null.

Each action $A \in \mathbb{A}(\mathcal{S})$ defines a mapping $h_A : \mathcal{A}(\mathcal{S}) \rightarrow \mathbb{A}^{\text{prim}}(\mathcal{S})$ (where $\mathbb{A}^{\text{prim}}(\mathcal{S})$ is the set of primitive actions). The functions h_A are defined recursively as follows:

$$h_{\text{do}[\vec{Y} \leftarrow \vec{y}]}(\vec{a}) = \text{do}[\vec{Y} \leftarrow \vec{y}]$$

and

$$h_{\text{if } \phi \text{ then } A \text{ else } B}(\vec{a}) = \begin{cases} h_A(\vec{a}) & \text{if } \vec{a} \Rightarrow \phi \\ h_B(\vec{a}) & \text{if } \vec{a} \Rightarrow \neg\phi. \end{cases}$$

3.1 Representation

A *subjective causal expected utility* representation understands preferences as maximizing the expected utility of an action, relative to uncertainty regarding the values of exogenous variables. Specifically, the representation is governed by

- M : a model that dictates the causal equations
- \mathbf{p} : a probability distribution on contexts
- \mathbf{u} : a utility function on atoms.

An action (along with the model M) can be thought of as a operation that assigns values to all variables given a context. Specifically given a model M , for each $A \in \mathbb{A}(\mathcal{S})$, define $\beta_A^M : \mathcal{C}(\mathcal{S}) \rightarrow \mathcal{A}(\mathcal{S})$ as

$$\begin{aligned}\beta_A^M(\vec{u}) &= \text{the unique atom } \vec{a} \in \mathcal{A}(\mathcal{S}) \\ \text{such that } (M_{h_A(\vec{a}_{M,\vec{u}})}, \vec{u}) &\models \vec{a}.\end{aligned}$$

Unpacking this: if the context is \vec{u} , then the initial assignment of variables is expressed by the atom $\vec{a}_{M,\vec{u}}$, so the primitive action prescribed by A is $h_A(\vec{a}_{M,\vec{u}})$. Thus, the final assignment of variables (i.e., after the intervention $h_A(\vec{a}_{M,\vec{u}})$) is determined by the situation $(M_{h_A(\vec{a}_{M,\vec{u}})}, \vec{u})$: this is what is captured by β .

Definition 1. $(M, \mathbf{p}, \mathbf{u})$ is a subjective causal utility representation of $\succsim_{\mathcal{S}}$, where $M = (\mathcal{S}, \mathcal{F})$ and $\mathcal{S} = (\mathcal{U}, \mathcal{V}, \mathbb{R})$, if \mathbf{p} is a probability on $\mathcal{R}(\mathcal{U})$, \mathbf{u} is a utility function on $\mathbb{A}(\mathcal{S})$, and $A \succsim_{\mathcal{S}} B$ iff

$$\sum_{\vec{u} \in \mathcal{C}(\mathcal{S})} \mathbf{u}(\beta_A^M(\vec{u})) \mathbf{p}(\vec{u}) \geq \sum_{\vec{u} \in \mathcal{C}(\mathcal{S})} \mathbf{u}(\beta_B^M(\vec{u})) \mathbf{p}(\vec{u}). \quad (1)$$

In a representation, we think of contexts as states and atoms as outcomes, so β_A^M can be viewed as a function from states to outcomes. Thus, by associating A with β_A^M , we can think of A as a function from states to outcomes—exactly how Savage views acts.

4 Axioms

In this section, we discuss the axioms that we need to ensure the existence of a causal expected utility representation.

The first axiom is the cancellation axiom of Blume, Easley, and Halpern [2021]. To define this we need the notion of a multiset, which can be thought of as a set that allows for multiple instances of each of its elements; two multisets are equal just in case they contain the same elements with the same multiplicities. For example, the multiset $\{\{a, a, a, b, b\}\}$ is different from the multiset $\{\{a, a, b, b, b\}\}$: both multisets have five elements, but the multiplicities of a and b differ in the two multisets (assuming $a \neq b$).

Using this notation, we can state the cancellation axiom.

Axiom 1 (Cancellation). If $A_1 \dots A_n, B_1 \dots B_n \in \mathbb{A}$ and $\{\{h_{A_1}(\vec{a}) \dots h_{A_n}(\vec{a})\}\} = \{\{h_{B_1}(\vec{a}) \dots h_{B_n}(\vec{a})\}\}$ for all $\vec{a} \in \mathcal{A}(\mathcal{S})$, then $A_i \succsim B_i$ for all $i < n$ implies $B_n \succsim A_n$.

As in the prior literature, cancellation allows us to construct a state space and an additively separable utility representation of $\succsim_{\mathcal{V}}$. The state space that is constructed via the

BEH methodology does not have any causal structure. The remaining axioms ensure that we can construct a subjective utility representation that is a causal model.

As the remaining axioms speak directly to the structure of the *do* action, the following notation will be helpful: For each atom $\vec{a} \in \mathcal{A}(\mathcal{S})$, write $\mathbf{do}[\vec{Y} \leftarrow \vec{y}] \rightsquigarrow_{\vec{a}} (Z = z)$ as shorthand for the indifference relation

$$\mathbf{if} \vec{a} \mathbf{then} \mathbf{do}[\vec{Y} \leftarrow \vec{y}, Z \leftarrow z] \sim \mathbf{if} \vec{a} \mathbf{then} \mathbf{do}[\vec{Y} \leftarrow \vec{y}]. \quad (2)$$

This says, essentially, that in the context characterized by the atom \vec{a} , setting $\vec{Y} = \vec{y}$ results in $Z = z$. To understand why, notice that intervening to set \vec{Y} to \vec{y} causes Z to equal z if and only if the further intervention setting Z to z has no effect; $\mathbf{do}[\vec{Y} \leftarrow \vec{y}] \rightsquigarrow_{\vec{a}} (Z = z)$ can be thought of as capturing this situation, as the decision maker sees no additional benefit to the additional intervention on Z . Of course, it could be that Z is in fact set to some different value z' , but the decision maker is indifferent between $Z = z$ and $Z = z'$, given the values of the other variables. This indifference is irrelevant as far as the existence of a subjective utility representation goes; we can simply treat the worlds where $Z = z$ and $Z = z'$ identically.

The next axiom, *model uniqueness*, ensures that there is only one model in our representation, rather than a distribution over models; that is, the decision maker has no uncertainty about the structural equations. The axiom itself just says that, for each context \vec{u} , there is at most one atom \vec{a} compatible with \vec{u} that is non-null. Intuitively, if M represents $\succsim_{\mathcal{S}}$, then $(M, \vec{u}) \models \vec{a}$.

Axiom 2 (Model Uniqueness). For each context \vec{u} , there exists a most one atom $\vec{a} \in \mathcal{A}(\mathcal{S})$ such that $\vec{a} \Rightarrow (\mathcal{U} = \vec{u})$ and \vec{a} is non-null.

The next axiom, *Definiteness*, ensures that there must be some value $x \in \mathcal{R}(X)$ such that $X = x$ after intervening to set \vec{Y} to \vec{y} . This is essentially the content of the axiom of the same name introduced by Galles and Pearl [1998], also used by Halpern [2000].

Axiom 3 (Definiteness). For each atom \vec{a} , vector of \vec{Y} of endogenous variables, $\vec{y} \in \mathcal{R}(\vec{Y})$, and endogenous variable $X \notin \vec{Y}$, there exists some $x \in \mathcal{R}(X)$ such that $\mathbf{do}[\vec{Y} \leftarrow \vec{y}] \rightsquigarrow_{\vec{a}} (X = x)$.

The next axiom, *Centeredness*, dictates that intervening to set variables to their actual values does not change the values of other variables; that is, trivial interventions are indeed trivial. It is named after *centering* property considered by Lewis [1973], which ensures that the closest world to a world ω is ω itself. Let $\vec{a}|_{\vec{Y}}$ denote the restriction of the atom \vec{a} to the conjuncts in \vec{Y} .

Axiom 4 (Centeredness). For each atom \vec{a} , vector of endogenous variables \vec{Y} , and endogenous variable $X \notin \vec{Y}$, we have $\mathbf{do}[\vec{Y} \leftarrow \vec{a}|_{\vec{Y}}] \rightsquigarrow_{\vec{a}} (X = \vec{a}|_X)$.

Finally, as we are interested in recursive causal models, we require an axiom that forces the dependency order on endogenous variables to be acyclic. Towards this, for $X, Y \in \mathcal{V}$, say

that X is unaffected by Y (given \vec{a}) if

$$\text{do}[\vec{Z} \leftarrow \vec{z}] \rightsquigarrow_{\vec{a}} (X = x) \text{ iff } \text{do}[\vec{Z} \leftarrow \vec{z}, Y \leftarrow y] \rightsquigarrow_{\vec{a}} (X = x) \quad (3)$$

for all $\vec{Z} \in \mathcal{V} \setminus \{X, Y\}$, $\vec{z} \in \mathcal{R}(\vec{Z})$, $y \in \mathcal{R}(y)$, and $x \in \mathcal{R}(x)$. So X is unaffected by Y if there is no intervention on Y that changes the decision maker's perception of X (conditional on atom \vec{a}). If this relation does not hold, then X is affected by Y , written $Y \rightsquigarrow_{\vec{a}} X$. Let $\rightsquigarrow = \cup_{\vec{a}} \rightsquigarrow_{\vec{a}}$. Our final axiom states that \rightsquigarrow is acyclic; it is inspired by the corresponding axiom in [Halpern, 2000].

Axiom 5 (Recursivity). \rightsquigarrow is acyclic.

Note that Axiom A5 implies that $\rightsquigarrow_{\vec{a}}$ is acyclic for each atom \vec{a} .

As we now show, a preference order satisfies these axioms iff it has subjective causal expected utility representation $(M, \mathbf{p}, \mathbf{u})$. Moreover, if \rightsquigarrow_S satisfies Axiom A3*, then M is unique over the set of non-null contexts.

There may, in general, be many different subjective causal expected utility representations of the same preference relations. This is because, when the decision maker is indifferent between distinct atoms, the preference relation cannot distinguish between them, and hence cannot distinguish between structural equations that yield distinct but equally valued atoms.

Consider the following strengthening of A3:

Axiom* 3 (Strong Definiteness). For each non-null atom \vec{a} , vector of \vec{Y} of endogenous variables, $\vec{y} \in \mathcal{R}(\vec{Y})$, and endogenous variable $X \notin \vec{Y}$, there exists a unique $x \in \mathcal{R}(X)$ such that $\text{do}[\vec{Y} \leftarrow \vec{y}] \rightsquigarrow_{\vec{a}} (X = x)$.

Strengthening A3 to A3*, is both necessary and sufficient to avoid multiple representations.

Definition 2. A subjective causal utility representation of a preference relation of $(M', \mathbf{p}', \mathbf{u}')$, non-null contexts \vec{u} , and formulas $\phi \in \mathcal{L}^+(\mathcal{S})$,

$$(M, \vec{u}) \models \phi \quad \text{if and only if} \quad (M', \vec{u}) \models \phi.$$

Theorem 2. A subjective causal utility representation $(M, \mathbf{p}, \mathbf{u})$ of \rightsquigarrow_S is identified if and only if \rightsquigarrow_S satisfies Axiom A3*.

5 Proof of Theorems 1 and 2

Before getting into the details of the proof, we describe it at a higher level.

Given a preference order \rightsquigarrow_S , the first step is to construct a Lewis-style model M_{\rightsquigarrow_S} that represents \rightsquigarrow_S . We want M_{\rightsquigarrow_S} to be recursive in a precise sense. So for each atom $\vec{a} \in \mathcal{A}(\mathcal{S})$, we define a strict linear order $<_{\vec{a}}$ of $\mathcal{A}(\mathcal{S})$, intended to represent the closeness of atoms to \vec{a} . We then show that, for each atom \vec{a} , intervening to set $\vec{Y} \leftarrow \vec{y}$ has the same consequence as intervening to set all variables to their values in the $<_{\vec{a}}$ -closest atom to \vec{a} in which $\vec{Y} = \vec{y}$ holds (see Lemma 3). This

is the only property that BH used to prove their representation theorem, which used Lewis-style models (see their Lemma 5 and Theorem 1). It follows that we get the desired Lewis-style representation of M_{\rightsquigarrow_S} . Moreover, in M_{\rightsquigarrow_S} , $<_{\vec{a}}$ does in fact represent the closeness of atoms to \vec{a} . These observations allow us to appeal to results of Halpern [2013], and convert M_{\rightsquigarrow_S} to a causal model that also represents \rightsquigarrow_S .

So, to begin, we want to define $<_{\vec{a}}$. Axiom A5 ensures that we can extend \rightsquigarrow (and hence $\rightsquigarrow_{\vec{a}}$ for each atom \vec{a}) to a strict linear order on $\mathcal{U} \cup \mathcal{V}$, denoted \rightsquigarrow , such that for all $U \in \mathcal{U}$ and $X \in \mathcal{V}$ we have $U \rightsquigarrow X$. For atoms $\vec{b}, \vec{c} \in \mathcal{A}(\mathcal{S})$, let $Y_{\vec{b}, \vec{c}} \in \mathcal{U} \cup \mathcal{V}$ be the \rightsquigarrow -minimal variable on which \vec{b} and \vec{c} disagree. For each variable $X \in \mathcal{U} \cup \mathcal{V}$, define $\ll_{\vec{a}}^X$ to be some fixed strict linear order of $\mathcal{R}(X)$ whose minimal element is the value of X in atom \vec{a} . We define $<_{\vec{a}}$ to be, loosely speaking, a lexicographic order over atoms, ordering first over variables, using \rightsquigarrow , and then over values, using $\ll_{\vec{a}}$. We discuss the intuition behind this order after giving its definition:

Take $<_{\vec{a}}$ to be a strict linear order over $\mathcal{A}(\mathcal{S})$ such that

OR1. $\vec{b} <_{\vec{a}} \vec{c}$ if $Y_{\vec{a}, \vec{c}} \rightsquigarrow Y_{\vec{a}, \vec{b}}$, and

OR2. If $Y_{\vec{a}, \vec{b}} = Y_{\vec{a}, \vec{c}}$, then let Y denote $Y_{\vec{a}, \vec{c}}$, let \vec{Z} be the set of endogenous variables (strictly) \rightsquigarrow -less than Y , and let $\vec{z} = \vec{b}|_{\vec{Z}} (= \vec{c}|_{\vec{Z}})$. Then $\vec{b} <_{\vec{a}} \vec{c}$ if $\text{do}[\vec{Z} \leftarrow \vec{z}] \rightsquigarrow_{\vec{a}} \vec{b}|_Y$ and either

(i) not $\text{do}[\vec{Z} \leftarrow \vec{z}] \rightsquigarrow_{\vec{a}} \vec{c}|_Y$, or

(ii) $\text{do}[\vec{Z} \leftarrow \vec{z}] \rightsquigarrow_{\vec{a}} \vec{c}|_Y$ and $y_{\vec{b}} \ll_{\vec{a}}^Y y_{\vec{c}}$, where $y_{\vec{b}}$ (resp., $y_{\vec{c}}$) denotes the value of Y in \vec{b} (resp., \vec{c}).

In general, (OR1.) and (OR2.) do not completely determine the order; there may be atoms for which neither (OR1.) nor (OR2.) hold. The remainder of the order can be completed arbitrarily.

A few points of intuition regarding this order. In the Lewis-style model M_{\rightsquigarrow_S} that we construct, the states are represented by pairs of atoms. The operator $\rightsquigarrow_{\vec{a}}$ lets us probe the structural equations through the effect of interventions. When considering which of \vec{b} or \vec{c} is closer to \vec{a} according to $\rightsquigarrow_{\vec{a}}$, there are two criteria of lexicographic importance: (i) which of \vec{b} or \vec{c} coincides with \vec{a} longer (where *longer* means “for more variables, starting with the exogenous context and proceeding via \rightsquigarrow ”), and (ii) if both \vec{b} and \vec{c} deviate from \vec{a} at the same variable Y , does one coincide with the counterfactual assessment of \vec{a} given by $\rightsquigarrow_{\vec{a}}$, and if both do, is one closer to \vec{a} than the other according to $\ll_{\vec{a}}^Y$. These two criteria are captured by (OR1.) and (OR2.), respectively.

An example may help clarify: Let \vec{a} be the atom $U = 0 \wedge X = 0 \wedge Y = 0 \wedge Z = 0$, and assume that the order \rightsquigarrow on variables is $U \rightsquigarrow X \rightsquigarrow Y \rightsquigarrow Z$. Now consider the following atoms:

$$\vec{b} = U = 0 \wedge X = 0 \wedge Y = 0 \wedge Z = 1$$

$$\vec{c} = U = 0 \wedge X = 0 \wedge Y = 1 \wedge Z = 1$$

$$\vec{c}' = U = 0 \wedge X = 0 \wedge Y = 1 \wedge Z = 0.$$

(OR1.) states that if an atom coincides with \vec{a} longer, it is closer to \vec{a} , so \vec{b} is closer to \vec{a} than either \vec{c} or \vec{c}' . Since \vec{c} and \vec{c}'

agree up to Z , their closeness to \vec{a} depends on the value(s) of z such that $\text{do}[X \leftarrow 0, Y \leftarrow 1] \sim_{\vec{a}} (Z = z)$. In particular, if we can have $z = 1$ (which is the value of Z according to \vec{c}) and not $z = 0$, then \vec{c} is closer to \vec{a} than \vec{c}' . Intuitively, this is because, according to $\sim_{\vec{a}}$, if we set X to 0 and Y to 1, Z would be 1, so \vec{c} is consistent with the equations we plan to use in the context encoded by \vec{a} , while \vec{c}' is not. Similarly, if we can have $z = 0$ and not $z = 1$, then \vec{c}' is closer. Because of the possibility of indifference, both $z = 0$ and $z = 1$ could be consistent with $\sim_{\vec{a}}$. In this case, the order $\ll_{\vec{a}}^Y$ serves as a consistent method of breaking ties.⁴

Let $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$ denote the $\sim_{\vec{a}}$ -minimal atom satisfying $\vec{Y} = \vec{y}$.

Lemma 1. Suppose that \vec{Y} is a set of endogenous variables, $X \in \mathcal{V} \setminus \vec{Y}$, and $\vec{Z}_X \subseteq \mathcal{V}$ is the (possibly empty) set of endogenous variables strictly \sim -less than X . Then we have

$$\text{M1. } \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\mathcal{U}} = \vec{a}|_{\mathcal{U}}, \text{ and}$$

$$\text{M2. The value of } X \text{ in } \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]} \text{ is the } \ll_{\vec{a}}^X \text{-minimal element of }$$

$$m(\vec{a}, X) =$$

$$\{x \in \mathcal{R}(X) : \text{do}[\vec{Z}_X \leftarrow \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\vec{Z}_X}] \sim_{\vec{a}} (X = x)\}.$$

Proof. Since \mathcal{U} forms the initial segment of \sim , it follows from (OR1.) that (M1.) must hold. In more detail, if $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\mathcal{U}} \neq \vec{a}|_{\mathcal{U}}$, then let \vec{b} be such that $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\mathcal{V}} = \vec{b}|_{\mathcal{V}}$ and $\vec{b}|_{\mathcal{U}} = \vec{a}|_{\mathcal{U}}$. Since \mathcal{U} forms the initial segment of \sim , it follows from (OR1.) that $\vec{b} <_{\vec{a}} \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\mathcal{U}}$, yet \vec{b} satisfies $\vec{Y} = \vec{y}$. This is a contradiction.

To see that (M2.) holds, suppose by way of contradiction that it does not hold. Then there exists a variable X such that the value of X in $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$ is not the $\ll_{\vec{a}}^X$ -minimal element of $m(\vec{a}, X)$ (note that $m(\vec{a}, X)$ is non-empty by Axiom A3). Let \vec{b} be the atom that coincides with $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$ for all variables except X , and the value of X in \vec{b} is the $\ll_{\vec{a}}^X$ -minimal element of $m(\vec{a}, X)$. There are two cases:

- $\vec{a}|_{\vec{Z}_X} = \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\vec{Z}_X}$: From Axiom A4, it follows that $\text{do}[\vec{Z}_X \leftarrow \vec{z}_X] \sim_{\vec{a}} \vec{a}|_X$. By the definition of $\ll_{\vec{a}}^X$, the value of X in \vec{a} is the $\ll_{\vec{a}}^X$ -minimal element of $m(\vec{a}, X)$. By assumption, the value of X in $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_X$ is not the $\ll_{\vec{a}}^X$ -minimal element of $m(\vec{a}, X)$, while the value of X in \vec{b} is the $\ll_{\vec{a}}^X$ -minimal element of $m(\vec{a}, X)$. It follows that $\vec{b} <_{\vec{a}} \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$, a contradiction.
- $\vec{a}|_{\vec{Z}_X} \neq \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\vec{Z}_X}$: Since $\vec{a}|_{\vec{Z}_X} \neq \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\vec{Z}_X} = \vec{b}|_{\vec{Z}_X}$, we have $Y_{\vec{a}, \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}} = Y_{\vec{a}, \vec{a}, \vec{b}}$, so by OR2., $\vec{b} <_{\vec{a}} \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$, and again we have a contradiction.

⁴ $\ll_{\vec{a}}^Y$ is arbitrary except that its initial element is the value of Y in \vec{a} , reflecting the fact that beyond coinciding with \vec{a} if possible, this tie-breaking is arbitrary. Nonetheless, an order still needs to be fixed to ensure tie-breaking is consistent across different interventions.

□

Lemma 2. For all atoms $\vec{a} \in \mathcal{A}(\mathcal{S})$, (disjoint) vectors of endogenous variables \vec{Y} and \vec{Z} , $\vec{y} \in \mathcal{R}(\vec{Y})$, $\vec{z} \in \mathcal{R}(\vec{Z})$, endogenous variables $X \notin \vec{Z} \cup \vec{Y}$ such that for all $Y \in \vec{Y}$, $X \sim_{\vec{a}} Y$, and $x \in \mathcal{R}(X)$, we have

$$\begin{aligned} \text{do}[\vec{Z} \leftarrow \vec{z}] &\sim_{\vec{a}} (X = x) \text{ iff} \\ \text{do}[\vec{Z} \leftarrow \vec{z}, \vec{Y} \leftarrow \vec{y}] &\sim_{\vec{a}} (X = x). \end{aligned}$$

Proof. By assumption, for all $Y \in \vec{Y}$, we have $Y \not\sim_{\vec{a}} X$. The lemma then follows from a simple induction argument on the variables of \vec{Y} , appealing to (3) in each step. □

The following lemma shows that, starting with \vec{a} , the effect of setting \vec{Y} to \vec{y} is the same as setting all the endogenous variables to the value that results from setting \vec{Y} to \vec{y} . Intuitively, this is because for a variable $X \notin \vec{Y}$, we are setting it to the value it already has after setting \vec{Y} to \vec{y} , so nothing further changes.

Lemma 3. For all atoms $\vec{a} \in \mathcal{A}(\mathcal{S})$, vectors of endogenous variables \vec{Y} , and $\vec{y} \in \mathcal{R}(\vec{Y})$, we have that

$$\text{if } \vec{a} \text{ then do}[\vec{Y} \leftarrow \vec{y}] \sim \text{if } \vec{a} \text{ then do}[\mathcal{V} \leftarrow \vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}|_{\mathcal{V}}].$$

Proof. We prove this by induction on the variables in \mathcal{V} with respect to the order \sim . Fix some $X \in \mathcal{V} \setminus \vec{Y}$, let $\vec{Z} \subseteq \mathcal{V}$ denote the (possibly empty) vector of endogenous variables strictly \sim -less than X , and let \vec{z} denote the values of the variables \vec{Z} in $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$. To simplify notation, let x denote the value of X in $\vec{a}_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$, let $\vec{Y}' \subseteq \vec{Y}$ consist of those variables \sim -greater than X , and let \vec{y}' be the restriction of \vec{y} to \vec{Y}' . We now show that if

$$\text{if } \vec{a} \text{ then do}[\vec{Y} \leftarrow \vec{y}] \sim \text{if } \vec{a} \text{ then do}[\vec{Z} \leftarrow \vec{z}, \vec{Y}' \leftarrow \vec{y}'] \quad (4)$$

then

$$\begin{aligned} \text{if } \vec{a} \text{ then do}[\vec{Y} \leftarrow \vec{y}] \\ \sim \\ \text{if } \vec{a} \text{ then do}[\vec{Z} \leftarrow \vec{z}, \vec{Y}' \leftarrow \vec{y}', X \leftarrow x]. \end{aligned} \quad (5)$$

From (M2.), it follows that $\text{do}[\vec{Z} \leftarrow \vec{z}] \sim_{\vec{a}} (X = x)$. Applying Lemma 2, we obtain

$$\text{do}[\vec{Z} \leftarrow \vec{z}, \vec{Y}' \leftarrow \vec{y}'] \sim_{\vec{a}} (X = x). \quad (6)$$

Thus,

$$\begin{aligned} \text{if } \vec{a} \text{ then do}[\vec{Z} \leftarrow \vec{z}, \vec{Y}' \leftarrow \vec{y}', X \leftarrow x] \\ \sim \text{if } \vec{a} \text{ then do}[\vec{Z} \leftarrow \vec{z}, \vec{Y}' \leftarrow \vec{y}'] \\ \sim \text{if } \vec{a} \text{ then do}[\vec{Y} \leftarrow \vec{y}], \end{aligned}$$

where the first indifference comes from from (6) and the definition of $\sim_{\vec{a}}$, and the second from (4). □

BH show in their Theorem 1 that (in our terminology), given a preference $\succsim_{\mathcal{S}}$ satisfying their Lemma 5 (which is a direct translation of our Lemma 3) and the Cancellation axiom, and a family $<_a$ of linear orders, one for each atom, that they can construct a Lewis-style model that represents $\succsim_{\mathcal{S}}$. We outline the construction below, as well as defining formally what it means for a Lewis-style model to represent $\succsim_{\mathcal{S}}$.

The states in the Lewis-style model M are pairs (\vec{a}, \vec{a}') of atoms. Roughly speaking, since we are trying to capture the effect of interventions, we can think of \vec{a} as the current world (the result of performing the intervention) and \vec{a}' as the world before the intervention was performed. The truth of formulas in $\mathcal{L}(\mathcal{S})$ is completely determined by the first component. That is, $(M, (\vec{a}, \vec{a}')) \models Y = y$ if $Y = y$ is a conjunct of \vec{a} ; we extend to negation, conjunction, and negation in the standard way. To extend this semantics to $\mathcal{L}^+(\mathcal{S})$, we need a closeness relation for each pair (\vec{a}, \vec{a}') . Let $\sqsubset_{\vec{a}, \vec{a}'}$ be an arbitrary family of strict linear orders on $\mathcal{A}(\mathcal{S}) \times \mathcal{A}(\mathcal{S})$ such that

$$\langle \vec{b}, \vec{b}' \rangle \sqsubset_{\langle \vec{a}, \vec{a}' \rangle} \langle \vec{c}, \vec{c}' \rangle \text{ whenever } \vec{b} <_{\vec{a}} \vec{c}$$

and

$$\langle \vec{b}, \vec{a} \rangle \sqsubset_{\langle \vec{a}, \vec{a}' \rangle} \langle \vec{b}, \vec{b}' \rangle \text{ for all } \vec{b}' \neq \vec{a}$$

(so that $\langle \vec{a}, \vec{a} \rangle$ is the minimal element of $\sqsubset_{\langle \vec{a}, \vec{a}' \rangle}$). This completes the description of M .

Let $\text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}]) \in \mathcal{A}(\mathcal{S}) \times \mathcal{A}(\mathcal{S})$ denote the $\sqsubset_{\langle \vec{a}, \vec{a}' \rangle}$ -minimal state such that $(M, \text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models \vec{Y} \leftarrow \vec{Y}$. Thus, $(M, \langle \vec{a}, \vec{a}' \rangle) \models \text{do}[\vec{Y} \leftarrow \vec{y}] \phi$ iff $(M, \text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models \phi$. It is easy to check that

$$\text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}]) = \langle \vec{a}_{\vec{Y} \leftarrow \vec{y}}, \vec{a} \rangle. \quad (7)$$

Thus, the first component of the minimal state that results from performing the intervention $\vec{Y} \leftarrow \vec{y}$ in a state of the form $\langle \vec{a}, \vec{a}' \rangle$ encodes the atom that results when the intervention is performed starting at \vec{a} , while the second component keeps track of the atom we started with.

BH define an analogue of our function h_A , which we denote h'_A , that associates with each of their actions an action of the form $\text{do}[\vec{Y} \leftarrow \vec{y}]$.⁵ For a state ω , the analogue of $\beta_A^M(\omega)$ is $\text{MIN}_{\omega}(h'_A(\omega))$.⁶ BH show that there exists a probability measure \mathbf{p} and utility function \mathbf{u} , both defined on the states of M , such that $A \succsim_{\mathcal{S}} B$ iff

$$\begin{aligned} & \sum_{\substack{(\vec{a}, \vec{a}') \in \\ \mathcal{A}(\mathcal{S}) \times \mathcal{A}(\mathcal{S})}} \mathbf{p}(\langle \vec{a}, \vec{a}' \rangle) \cdot \mathbf{u}(\text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(h'_A(\langle \vec{a}, \vec{a}' \rangle))) \\ & \geq \sum_{\substack{(\vec{a}, \vec{a}') \in \\ \mathcal{A}(\mathcal{S}) \times \mathcal{A}(\mathcal{S})}} \mathbf{p}(\langle \vec{a}, \vec{a}' \rangle) \cdot \mathbf{u}(\text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(h'_B(\langle \vec{a}, \vec{a}' \rangle))). \end{aligned} \quad (8)$$

This is the sense in which the Lewis-style model M represents $\succsim_{\mathcal{S}}$.

⁵BH allow for (but do not require) a richer set of interventions, they allow $\text{do}[\phi]$ for any consistent formula ϕ in the language.

⁶The states for us are contexts, whereas for BH they are pairs of atoms, so using the same symbol ω for states is somewhat of an abuse of notation; we hope that our intention is clear here.

Since BH's Lemma 5 is a direct translation of our Lemma 3, it follows from our axioms as well that this Lewis-style model M represents $\succsim_{\mathcal{S}}$. But we are not done; we need to get a causal model that represents $\succsim_{\mathcal{S}}$. So we want to convert the Lewis-style model constructed by BH to a causal model $M^C = (\mathcal{S}^C, \mathcal{F}^C)$, and construct an appropriate probability \mathbf{p}^C on the contexts of M , and utility \mathbf{u}^C on $\mathbb{A}(\mathcal{S})$.

The first step is easy. We start with a signature \mathcal{S} that determines the language; we take $\mathcal{S}^C = \mathcal{S}$. By Axiom A2, it follows that for each context $\vec{u} \in \mathcal{C}(\mathcal{S})$, there at most one atom such that such that $\vec{a} \Rightarrow (\mathcal{U} = \vec{u})$ and \vec{a} is non-null. Let \mathcal{C}^\dagger denote the set of contexts for which such a non-null atom exists, and for each context $\vec{u} \in \mathcal{C}^\dagger$, let $\vec{a}_{\vec{u}}$ denote this non-null atom. We define \mathbf{p}^C using the probability \mathbf{p} in the model M provided by BH so that it has support \mathcal{C}^\dagger :

$$\mathbf{p}^C(\vec{u}) = \frac{\sum_{\vec{a}' \in \mathcal{A}(\mathcal{S})} \mathbf{p}(\langle \vec{a}_{\vec{u}}, \vec{a}' \rangle)}{\sum_{\vec{u}' \in \mathcal{C}^\dagger} \sum_{\vec{a}' \in \mathcal{A}(\mathcal{S})} \mathbf{p}(\langle \vec{a}_{\vec{u}'}, \vec{a}' \rangle)}.$$

Observe that, since the truth of formulas in M at state $\langle \vec{a}, \vec{a}' \rangle$ is fully determined by \vec{a} , $h'_A(\langle \vec{a}, \vec{a}' \rangle)$ does not depend on \vec{a}' . In particular, $h'_A(\langle \vec{a}, \vec{a}' \rangle) = h_A(\vec{a})$. Similarly, by (7), $\text{MIN}_{\langle \vec{a}, \vec{a}' \rangle}(\cdot)$ does not depend on \vec{a}' . We can therefore rewrite the BH representation (8) as

$$\begin{aligned} & \sum_{\vec{u} \in \mathcal{C}^\dagger} \mathbf{p}^C(\vec{u}) \cdot \mathbf{u}(\text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(h_A(\vec{a}_{\vec{u}}))) \\ & \geq \sum_{\vec{u} \in \mathcal{C}^\dagger} \mathbf{p}^C(\vec{u}) \cdot \mathbf{u}(\text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(h_B(\vec{a}_{\vec{u}}))). \end{aligned} \quad (9)$$

We next want to show that there exists a set of equations \mathcal{F}^C such that $M^C = (\mathcal{S}^C, \mathcal{F}^C)$ satisfies the same formulas as M . Specifically:

Lemma 4. *There exists a set of equations \mathcal{F}^C such that for all $\phi \in \mathcal{L}(\mathcal{S})$ and all interventions $\vec{Y} \leftarrow \vec{y}$, we have that*

$$\begin{aligned} & (M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}^C, \vec{u}) \models \phi \quad \text{iff} \\ & (M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models \phi. \end{aligned} \quad (10)$$

Note that taking $\vec{Y} = \emptyset$ gives a special case of (10): $(M^C, \vec{u}) \models \phi$ iff $(M, \langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle) \models \phi$.

Proof. Let Y_1, \dots, Y_k be the ordering on endogenous variables according to \rightsquigarrow . We define \mathcal{F}^C for the endogenous variables by induction j so that (10) for all formulas $Y_j = y_j$.

For the base case, note that the value of the variable Y_1 does not depend on the values of any other endogenous variable. If is straightforward to define F_{Y_1} so that $F_{Y_1}(\vec{u}) = y_1$ iff $a_{\vec{u}} \Rightarrow Y_1 = y_1$. It easily follows that $(M^C, \vec{u}) \models Y_1 = y_1$ iff $(M, \langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle) \models Y_1 = y_1$. Moreover, if $Y_1 \notin \vec{Y}$, then $(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}^C, \vec{u}) \models Y_1 = y_1$ iff $(M^C, \vec{u}) \models Y_1 = y_1$. By Lemma 2 (taking $Z = \emptyset$ and $X = Y_1$) it follows that $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models Y_1 = y_1$ iff $(M, \langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle) \models Y_1 = y_1$. On the other hand, if $Y_1 \in \vec{Y}$, then $(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}^C, \vec{u}) \models Y_1 = y_1$ iff $\vec{y}|_{Y_1} = y_1$ and similarly, $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models Y_1 = y_1$ iff $\vec{y}|_{Y_1} = y_1$. Thus, we get the desired result in the base case.

For the inductive case, suppose that we have proved the result for Y_1, \dots, Y_j . We can define $F_{Y_{j+1}}$ to be function of the exogenous variables and Y_1, \dots, Y_j such that for all $k \in$

$\{1, \dots, j\}$, and $y_i \in \mathcal{R}(Y_i)$ for $i = 1, \dots, j$, we have that $F_{Y_{j+1}}(\vec{u}, y_1, \dots, y_j) = y_{k+1}$ iff $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y}_1 \leftarrow y_1, \dots, \vec{Y}_k \leftarrow y_k])) \models Y_{k+1} = y_{k+1}$. Now for arbitrary \vec{Y} , if $Y_{j+1} \notin \vec{Y}$, let $\vec{Y}' = \vec{Y} \cap \{Y_1, \dots, Y_j\}$. Let $\vec{y}' = \vec{y}|_{\vec{Y}'}$. Writing \vec{Y}'' for $\{Y_1, \dots, Y_j\}$, let \vec{y}'' be such that $(M^C_{\text{do}[\vec{Y}' \leftarrow \vec{y}']}, \vec{u}) \models \vec{Y}'' = \vec{y}''$. Then it follows from the induction hypothesis that $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y}' \leftarrow \vec{y}'])) \models \vec{Y}'' = \vec{y}''$. Moreover, it easily follows that $(M^C_{\text{do}[\vec{Y}' \leftarrow \vec{y}']}, \vec{u}) \models Y_{j+1} = y_{j+1}$ iff $(M^C_{\text{do}[\vec{Y}'' \leftarrow \vec{y}']}, \vec{u}) \models Y_{j+1} = y_{j+1}$ and $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}])) \models Y_{j+1} = y_{j+1}$ iff $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y}'' \leftarrow \vec{y}'])) \models Y_{j+1} = y_{j+1}$. Since, by the definition of $F_{Y_{j+1}}$, $(M^C_{\text{do}[\vec{Y}'' \leftarrow \vec{y}']}, \vec{u}) \models Y_{j+1} = y_{j+1}$ iff $(M, \text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y}'' \leftarrow \vec{y}'])) \models Y_{j+1} = y_{j+1}$, the result easily follows. The argument if $Y_{j+1} \in \vec{Y}$ is the same as in the base case.

The result for arbitrary formulas $\phi \in \mathcal{L}(\mathcal{S})$ is immediate (since conjunction, disjunction, and negation work the same way in both M and M^C). \square

From (7), it follows that $\text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(\text{do}[\vec{Y} \leftarrow \vec{y}]) = \langle (\vec{a}_{\vec{u}})_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{a} \rangle$; moreover, $(M, \langle (\vec{a}_{\vec{u}})_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{a} \rangle) \models (\vec{a}_{\vec{u}})_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$. Thus, by Lemma 4, $(M^C_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (\vec{a}_{\vec{u}})_{\text{do}[\vec{Y} \leftarrow \vec{y}]}$. Moreover, since $(M^c, \vec{u}) \models \vec{a}_{\vec{u}}$, that is, $\vec{a}_{M^C, \vec{u}} = \vec{a}_{\vec{u}}$, we have

$$\beta_A^{M^C}(\vec{u}) = (\vec{a}_{\vec{u}})_{h_A(\vec{a}_{\vec{u}})}. \quad (11)$$

We can now complete the proof that M^C represents the preference order. The utility function \mathbf{u} given by BH is defined on their states, which are pairs of atoms. We define \mathbf{u}^C over individual atoms as

$$\mathbf{u}^C(\vec{a}) = u(\langle \vec{a}, \vec{a}_{\vec{a}|_{\mathcal{U}}} \rangle) \quad (12)$$

if $\vec{a}|_{\mathcal{U}} \in \mathcal{C}^\dagger$ (and define it arbitrarily otherwise). Unpacking this, $\vec{a}|_{\mathcal{U}}$ is the context of the atom \vec{a} , and so, $\vec{a}_{\vec{a}|_{\mathcal{U}}}$ is the unique non-null atom with the same context as \vec{a} .

$$\begin{aligned} \mathbf{u}(\text{MIN}_{\langle \vec{a}_{\vec{u}}, \vec{a}_{\vec{u}} \rangle}(h_A(\vec{a}_{\vec{u}}))) &= \mathbf{u}(\langle (\vec{a}_{\vec{u}})_{h_A(\vec{a}_{\vec{u}})}, \vec{a}_{\vec{u}} \rangle) \quad (\text{from (7)}) \\ &= \mathbf{u}^C((\vec{a}_{\vec{u}})_{h_A(\vec{a}_{\vec{u}})}) \quad (\text{from (12)}) \\ &= \mathbf{u}^C(\beta_A^M(\vec{u})) \quad (\text{from (11)}) \end{aligned}$$

Substituting \mathbf{u}^C for \mathbf{u} in (9) thus delivers the desired representation.

To prove Theorem 2, the uniqueness claim, let $(M, \mathbf{p}, \mathbf{u})$ and $(M', \mathbf{p}', \mathbf{u}')$ be two representations of $\succ_{\mathcal{S}}$. For all $\vec{u} \in \mathcal{C}^\dagger$ and all vectors $\vec{Y}, \vec{y} \in \mathcal{R}(\vec{Y})$, endogenous variables $X \notin \vec{Y}$, and $x \in \mathcal{R}(X)$

$$(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x) \implies \text{do}[\vec{Y} \leftarrow \vec{y}] \sim_{\vec{a}_{\vec{u}}} (X = x); \quad (13)$$

similarly,

$$(M'_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x) \implies \text{do}[\vec{Y} \leftarrow \vec{y}] \sim_{\vec{a}_{\vec{u}}} (X = x). \quad (14)$$

Suppose that $\succ_{\mathcal{S}}$ satisfies A3*. Then there is a unique $x \in \mathcal{R}(X)$ that satisfies the right-hand relations of (13) and (14), so we have

$$(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x) \iff (M'_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x).$$

It easily follows that M and M' agree on all formulas in $\mathcal{L}^+(\mathcal{S})$, so M is identifiable.

For the converse, suppose that $M \neq M'$ for some $\vec{u} \in \mathcal{C}^\dagger$. In particular, this requires that there exists some vector $\vec{Y}, \vec{y} \in \mathcal{R}(\vec{Y})$ and endogenous variable $X \notin \vec{Y}$, such that $(M_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x)$ and $(M'_{\text{do}[\vec{Y} \leftarrow \vec{y}]}, \vec{u}) \models (X = x')$ and $x \neq x'$. But then, by (13) and (14), $\succ_{\mathcal{S}}$ violates A3*.

6 Discussion

We have given a representation theorem in the spirit of Savage [1954] that helps us understand a decision maker's (subjective) causal judgements: If a decision maker's preferences among actions that involve interventions satisfy certain axioms, then we can find a causal model M , a probability over contexts in M , and a utility on states in M such that the decision maker prefers action A to B iff A has higher expected utility than B . Moreover, we have shown if we add another axiom, then M is unique. Our approach builds on earlier work by BH, who proved an analogous representation theorem using Lewis-style models. Interestingly, other than the Cancellation axiom (which is due to BEH), the axioms used by BH are completely different from ours. For example, they do not define an analogue of our $\rightsquigarrow_{\vec{a}}$ relation, which plays a critical role in our definiteness and centeredness axioms. Roughly speaking, the BH axioms build on axioms for counterfactuals used by Lewis, while ours build on axioms for causal models introduced by Galles and Pearl [1998] and Halpern [2001]. In future work, we hope to better understand the connection between the axioms (e.g., to what extent could we use the BH axioms as a basis for a representation theorem for causal models), with the hope of understanding better the connection between Lewis-style models and causal models.

Here we have considered only one-step decisions; that is, the decision-maker performs an intervention, perhaps conditional on some test. It is clearly also of interest to consider sequential decisions. In practice, plans are composed of a sequence of steps; later interventions might depend on earlier interventions. This leads us to consider a richer set of actions such as **if** ϕ_1 **then** A_1 **else** B_1 ; **if** ϕ_2 **then** A_2 **else** B_2 , where the second action (**if** ϕ_2 **then** A_2 **else** B_2) is performed after the first. Getting a representation theorem for Lewis-style models in the presence of sequential actions seems significantly more complicated than in the case of one-step actions [Bjorndahl and Halpern, 2023]. It would be of interest to see what is involved in getting such a representation theorem using causal models.

Interestingly, we expect there to be significant differences between causal models and Lewis-style once we have sequential decisions. A decision maker who uses causal models will behave in a dynamically consistent way with respect to sequential interventions involving primitive actions: if the intervention $\text{do}[\vec{Y} \leftarrow \vec{y}, \vec{Z} \leftarrow \vec{z}]$ is preferred to the intervention

$\text{do}[\vec{Y} \leftarrow \vec{y}, \vec{Z} \leftarrow \vec{z}]$, then $\text{do}[\vec{Z} \leftarrow \vec{z}]$ must be preferred to $\text{do}[\vec{Z}' \leftarrow \vec{z}']$ in the model that results from taking action $\text{do}[\vec{Y} \leftarrow \vec{y}]$. In other words, the decision maker's preferences regarding sequences of primitive actions are invariant to timing, or order, although this is not necessarily the case for arbitrary interventions. This is because the interventions $\text{do}[\vec{Y} \leftarrow \vec{y}, \vec{Z} \leftarrow \vec{z}]$ and $\text{do}[\vec{Z} \leftarrow \vec{z}, \vec{Y} \leftarrow \vec{y}]$ are identical.

By way of contrast, the order in which primitive interventions are performed can have a significant impact in Lewis-style models. In such models, the closest-world operator is local; loosely speaking, the “distance” between two worlds depends on the world at which they are being contemplated. Because of this, intervening to make ϕ true—moving to the closest ϕ -world—also changes the relative closeness of other worlds; thus, the closest $\phi \wedge \psi$ -world may bare no relation to the closest ψ -world to the closest ϕ -world. This means that we lose some dynamic consistency in Lewis-style models. In future work, we hope to explore the issue of dynamic consistency and order dependence in both causal models and Lewis-style models, both because we believe that the issue of importance in its own right, and because it will help elucidate the differences between these two approaches to modeling causality.

References

- [Alexander and Gilboa, 2023] Y. Alexander and I. Gilboa. Subjective causality. *Revue économique*, 74(4):619–633, 2023.
- [Angrist and Pischke, 2009] J. D. Angrist and J.-S. Pischke. *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton University Press, 2009.
- [Bjorndahl and Halpern, 2021] A. Bjorndahl and J. Y. Halpern. Language-based decisions. In *Theoretical Aspects of Rationality and Knowledge: Proc. Eighteenth Conference (TARK 2021)*. 2021. The proceedings are published as *Electronic Proceedings in Theoretical Computer Science*.
- [Bjorndahl and Halpern, 2023] A. Bjorndahl and J. Y. Halpern. Sequential language-based decisions. In *Theoretical Aspects of Rationality and Knowledge: Proc. Nineteenth Conference (TARK 2023)*. 2023. The proceedings are published as *Electronic Proceedings in Theoretical Computer Science*.
- [Blume *et al.*, 2021] L. E. Blume, D. Easley, and J. Y. Halpern. Constructive decision theory. *Journal of Economic Theory*, 196, 2021. An earlier version, entitled “Re-doing the Foundations of Decision Theory”, appears in *Principles of Knowledge Representation and Reasoning: Proc. Tenth International Conference (KR ’06)*.
- [Cunningham, 2021] S. Cunningham. *Causal Inference: The Mixtape*. Yale university press, 2021.
- [Ellis and Thysen, 2021] Andrew Ellis and Heidi Christina Thysen. Subjective causality in choice. Available at <http://arxiv.org/abs/2106.05957>, 2021.
- [Galles and Pearl, 1998] D. Galles and J. Pearl. An axiomatic characterization of causal counterfactuals. *Foundation of Science*, 3(1):151–182, 1998.
- [Haavelmo, 1943] T. Haavelmo. The statistical implications of a system of simultaneous equations. *Econometrica*, 11:1–12, 1943.
- [Halpern, 2000] J. Y. Halpern. Axiomatizing causal reasoning. *Journal of A.I. Research*, 12:317–337, 2000.
- [Halpern, 2001] J. Y. Halpern. Substantive rationality and backward induction. *Games and Economic Behavior*, 37:425–435, 2001.
- [Halpern, 2013] J. Y. Halpern. From causal models to possible-worlds models of counterfactuals. *Review of Symbolic Logic*, 6(2):81–101, 2013.
- [Hernán and Robins, 2020] M. A. Hernán and J. M. Robins. *Causal inference: What If*. Chapman and Hall/CRC, 2020.
- [Lewis, 1973] D. Lewis. *Counterfactuals*. Harvard University Press, Cambridge, MA, 1973.
- [Machina and Schmeidler, 1992] M. Machina and D. Schmeidler. A more robust definition of subjective probability. *Econometrica*, 60(2):745–780, 1992.
- [Morgan and Winship, 2007] S. Morgan and C. Winship. *Counterfactuals and Causal inference*. Cambridge University Press, 2007.
- [Parascandola and Weed, 2001] M. Parascandola and D. L. Weed. Causation in epidemiology. *Journal of Epidemiology & Community Health*, 55(12):905–912, 2001.
- [Pearl, 2000] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000.
- [Pearl, 2009] J. Pearl. Causal inference in statistics. *Statistics Surveys*, 3:96–146, 2009.
- [Plowright *et al.*, 2008] R. K. Plowright, S. H. Sokolow, M. E. Gorman, P. Daszak, and J. E. Foley. Causal inference in disease ecology: investigating ecological drivers of disease emergence. *Frontiers in Ecology and the Environment*, 6(8):420–429, 2008.
- [Savage, 1954] L. J. Savage. *Foundations of Statistics*. Wiley, New York, 1954.
- [Schenone, 2020] P. Schenone. Causality: a decision theoretic foundation. Available at <http://arxiv.org/abs/1812.07414>, 2020.
- [Simon, 1953] H. A. Simon. Causal ordering and identifiability. In W. C. Hood and T. C. Koopmans, editors, *Studies in Econometric Methods*, Cowles Commission for Research in Economics, Monograph No. 14, pages 49–74. Wiley, New York, 1953.
- [Spirtes *et al.*, 1993] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. Springer-Verlag, New York, 1993.
- [Stalnaker, 1968] R. C. Stalnaker. A theory of conditionals. In N. Rescher, editor, *Studies in Logical Theory*, pages 98–112. Blackwell, Oxford, U.K., 1968.