

# When does one variable have a probabilistic causal effect on another?

David Johnston

November 12, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Variables and Probability Models</b>	<b>3</b>
2.1	Probability distributions, Markov kernels and string diagrams . .	3
2.1.1	Examples . . . . .	5
2.1.2	Example: comb insertion . . . . .	6
2.2	Semantics of observed and unobserved variables . . . . .	7
2.3	Events . . . . .	10
2.4	Probabilistic models for causal inference . . . . .	10
2.5	Order 0 gaps: probability distributions . . . . .	12
2.6	Order 1 gaps: conditional probabilities . . . . .	13
2.7	Order 2 gaps: probability combs . . . . .	16
2.8	Revisiting truncated factorisation . . . . .	19
2.9	Useful results . . . . .	20
2.9.1	Repeated variables . . . . .	20
2.9.2	Disintegrations . . . . .	22
<b>3</b>	<b>Decision theoretic causal inference</b>	<b>23</b>
3.1	Combs . . . . .	24
3.2	See-do models and classical statistics . . . . .	25
<b>4</b>	<b>Causal Bayesian Networks</b>	<b>27</b>
4.1	Proxy control . . . . .	31
<b>5</b>	<b>Potential outcomes</b>	<b>32</b>
<b>6</b>	<b>Appendix: see-do model representation</b>	<b>36</b>
<b>7</b>	<b>Appendix: Counterfactual representation</b>	<b>38</b>
7.1	Parallel potential outcomes representation theorem . . . . .	39
<b>8</b>	<b>Appendix: Connection is associative</b>	<b>42</b>

## 1 Introduction

Two widely used approaches to causal modelling are *graphical causal models* and *potential outcomes models*. Graphical causal models, which include Causal Bayesian Networks and Structural Causal Models, provide a set of *intervention* operations that take probability distributions and a graph and return a modified probability distribution (Pearl, 2009). Potential outcomes models feature *potential outcome variables* that represent the “potential” value that a quantity of interest would take under the right circumstances, a potential that may be realised if the circumstances actually arise, but will otherwise remain only a potential or *counterfactual* value (Rubin, 2005).

One challenge for both of these approaches is understanding how their causal primitives – interventions and potential outcome variables respectively – relate to the causal questions we are interested in. This challenge is related to the distinction, first drawn by (Korzybski, 1933), between “the map” and “the territory”. Causal models, like other models, are “maps” that purport to represent a “territory” that we are interested in understanding. Causal primitives are elements of the maps, and the things to which they refer are parts of the territory. The maps contain all the things that we can talk about unambiguously, so it is challenging to speak clearly about how parts of the maps relate to parts of the territory that fall outside of the maps.

For example, Hernán and Taubman (2008), who observed that many epidemiological papers have been published estimating the “causal effect” of body mass index and argued that, because *actions* affecting body mass index<sup>1</sup> are vaguely defined, potential outcome variables and causal effects themselves become ill-defined. We note that “actions targeting body mass index” are not elements of a potential outcomes model but “things to which potential outcomes should correspond”. The authors claim is that vagueness in the “territory” leads to ambiguity about elements of the “map” – and, as we have suggested, anything we can try to say about the territory is unavoidably vague. This seems like a serious problem.

In a response, Pearl (2018) argues that *interventions* (by which we mean the operation defined by a causal graphical model) are well defined, but may not always be a good model of an action. Pearl further suggests that interventions in graphical models correspond to “virtual interventions” or “ideal, atomic interventions”, and that perhaps carefully chosen interventions can be good models of actions. Shahar (2009), also in response, argued that interventions targeting body mass index applied to correctly specified graphical causal models will necessarily yield no effect on anything else which, together with Pearl’s suggestion,

---

<sup>1</sup>the authors use the term “intervention”, but they do not use it mean a formal operation on a graphical causal model, and we reserve the term for such operations to reduce ambiguity.

implies perhaps that an “ideal, atomic intervention” on body mass index cannot have any effect on anything else. If this is so, it seems that we are dealing with quite a serious case of vagueness – there is a whole body of literature devoted to estimating a “causal effect” that, it is claimed, is necessarily equal to zero! Authors of the original literature on the effects of BMI might counter that they were estimating something different that wasn’t necessarily zero, but as far as we are concerned such a response would only underscore the problem of ambiguity.

One of the key problems in this whole discussion is how the things we have called *interventions* – which are elements of causal models – relate to the things we have called *actions*, which live outside of causal models. One way to address this difficulty is to construct a bigger causal model that can contain both “interventions” and “actions”, and we can then speak unambiguously about how one relates to another. This is precisely what we do here.

- We need to talk about variables
- We use compatibility + string diagrams
- We consider causation in terms of “proxy control”

## 2 Variables and Probability Models

### 2.1 Probability distributions, Markov kernels and string diagrams

We make extensive use of probability theory, and the following is a brief introduction to the string diagram notation we use for probabilistic reasoning. This notation comes from the study of Markov categories. Markov categories are abstract categories that represent models of the flow of information. We can form Markov categories from collections of sets – for example, discrete sets or standard measurable sets – along with the Markov kernel product as the composition operation. Markov categories come equipped with a graphical language of *string diagrams*, and a coherence theorem which states that valid proofs using string diagrams correspond to valid theorems in *any* Markov category (Selinger, 2010). More comprehensive introductions to Markov categories can be found in Fritz (2020); Cho and Jacobs (2019). Thus, while we limit ourselves to discrete sets in this paper, any derivation that uses only string diagrams is more broadly applicable.

We say, given a variable  $X : \Omega \rightarrow X$ , a probability distribution  $\mathbb{P}^X$  is a probability measure on  $(X, \mathcal{X})$ . Recall that a probability measure is a  $\sigma$ -additive function  $\mathbb{P}^X : \mathcal{X} \rightarrow [0, 1]$  such that  $\mathbb{P}^X(\emptyset) = 0$  and  $\mathbb{P}^X(X) = 1$ . Given a second variable  $Y : \Omega \rightarrow Y$ , a conditional probability  $\mathbb{Q}^{X|Y}$  is a Markov kernel  $\mathbb{Q}^{X|Y} : X \rightarrow Y$  which is a map  $Y \times \mathcal{X} \rightarrow [0, 1]$  such that

1.  $y \mapsto \mathbb{Q}^{X|Y}(A|y)$  is  $\mathcal{B}$ -measurable for all  $A \in \mathcal{X}$
2.  $A \mapsto \mathbb{Q}^{X|Y}(A|y)$  is a probability measure on  $(X, \mathcal{X})$  for all  $y \in Y$

In the context of discrete sets, a probability distribution can be defined as a vector, and a Markov kernel a matrix.

**Definition 2.1** (Probability distribution (discrete sets)). A probability distribution  $\mathbb{P}$  on a discrete set  $X$  is a vector  $(\mathbb{P}(x))_{x \in X} \in [0, 1]^{|X|}$  such that  $\sum_{x \in X} \mathbb{P}(x) = 1$ . For  $A \subset X$ , define  $\mathbb{P}(A) = \sum_{x \in A} \mathbb{P}(x)$ .

**Definition 2.2** (Markov kernel (discrete sets)). A Markov kernel  $\mathbb{K} : X \rightarrow Y$  is a matrix  $(\mathbb{K}(y|x))_{x \in X, y \in Y} \in [0, 1]^{|X| \times |Y|}$  such that  $\sum_{y \in Y} \mathbb{K}(y|x) = 1$  for all  $x \in X$ . For  $B \subset Y$  define  $\mathbb{K}(B|x) = \sum_{y \in B} \mathbb{K}(y|x)$ .

In the graphical language, Markov kernels are drawn as boxes with input and output wires, and probability measures (which are kernels with the domain  $\{*\}$ ) are represented by triangles:

$$\mathbb{K} := \boxed{\mathbb{K}} \quad (1)$$

$$\mathbb{P} := \triangleleft \mathbb{P} \quad (2)$$

Two Markov kernels  $\mathbb{L} : X \rightarrow Y$  and  $\mathbb{M} : Y \rightarrow Z$  have a product  $\mathbb{LM} : X \rightarrow Z$ , given in the discrete case by the matrix product  $\mathbb{LM}(z|x) = \sum_{y \in Y} \mathbb{M}(z|y)\mathbb{L}(y|x)$ . Graphically, we represent products between compatible Markov kernels by joining wires together:

$$\mathbb{LM} := X \boxed{\mathbb{K}} \boxed{\mathbb{M}} Z \quad (3)$$

The Cartesian product  $X \times Y := \{(x, y) | x \in X, y \in Y\}$ . Given kernels  $\mathbb{K} : W \rightarrow Y$  and  $\mathbb{L} : X \rightarrow Z$ , the tensor product  $\mathbb{K} \otimes \mathbb{L} : W \times X \rightarrow Y \times Z$  given by  $(\mathbb{K} \otimes \mathbb{L})(y, z|w, x) := \mathbb{K}(y|w)\mathbb{L}(z|x)$ . The tensor product is graphically represented by drawing kernels in parallel:

$$\mathbb{K} \otimes \mathbb{L} := \begin{array}{c} W \boxed{\mathbb{K}} Y \\ X \boxed{\mathbb{L}} Z \end{array} \quad (4)$$

We read diagrams from left to right (this is somewhat different to Fritz (2020); Cho and Jacobs (2019); Fong (2013) but in line with Selinger (2010)), and any diagram describes a set of nested products and tensor products of Markov kernels. There are a collection of special Markov kernels for which we can replace the generic “box” of a Markov kernel with a diagrammatic elements that are visually suggestive of what these kernels accomplish.

The identity map  $\text{id}_X : X \rightarrow X$  defined by  $(\text{id}_X)(x'|x) = \llbracket x = x' \rrbracket$ , where the Iverson bracket  $\llbracket \cdot \rrbracket$  evaluates to 1 if  $\cdot$  is true and 0 otherwise, is a bare line:

$$\text{id}_X := X - X \quad (5)$$

We choose a particular 1-element set  $\{*\}$  that acts as the identity in the sense that  $\{*\} \times A \cong A \times \{*\} \cong A$  for any set  $A$ . The erase map  $\text{del}_X : X \rightarrow \{*\}$  defined by  $(\text{del}_X)(*|x) = 1$  is a Markov kernel that “discards the input”. It is drawn as a fuse:

$$\text{del}_X := \text{---} * X \quad (6)$$

The copy map  $\text{copy}_X : X \rightarrow X \times X$  defined by  $(\text{copy}_X)(x', x''|x) = \llbracket x = x' \rrbracket \llbracket x = x'' \rrbracket$  is a Markov kernel that makes two identical copies of the input. It is drawn as a fork:

$$\text{copy}_X := X \text{---} \begin{array}{c} X \\ \diagup \diagdown \\ X \end{array} \quad (7)$$

The swap map  $\text{swap}_{X,Y} : X \times Y \rightarrow Y \times X$  defined by  $(\text{swap}_{X,Y})(y', x'|x, y) = \llbracket x = x' \rrbracket \llbracket y = y' \rrbracket$  swaps two inputs, and is represented by crossing wires:

$$\text{swap}_X := \begin{array}{c} \diagup \diagdown \\ \diagdown \diagup \end{array} \quad (8)$$

Because we anticipate that the graphical notation will be unfamiliar, we will include some examples in the next section.

### 2.1.1 Examples

When translating string diagram notation to integral notation, a number of identities can speed up the process.

For arbitrary  $\mathbb{K} : X \times Y \rightarrow Z$ ,  $\mathbb{L} : W \rightarrow Y$

$$[(\text{id}_X \otimes \mathbb{L})\mathbb{K}](A|x, w) = \int_Y \int_X \mathbb{K}(z|x', y') \mathbb{L}(dy'|w) \text{id}_X(dx'|x) \quad (9)$$

$$= \int_Y \mathbb{K}(z|x, y') \mathbb{L}(dy'|w) \quad (10)$$

That is, an identity map passes its input to the next kernel in the product.

For arbitrary  $\mathbb{K} : X \times Y \times Y \rightarrow Z$  (where we apply the above shorthand in the first line):

$$[(\text{id}_X \otimes \text{copy}_Y)\mathbb{K}](A|x, y) = \int_Y \int_Y \mathbb{K}(A|x, y', y'') \text{copy}_Y(dy' \times dy''|y) \quad (11)$$

$$= \mathbb{K}(A|x, y, y) \quad (12)$$

That is, the copy map passes along two copies of its input to the next kernel in the product.

For a collection of kernels  $\mathbb{K}^n : Y^n \rightarrow Z$ ,  $n \in [n]$ , define  $(y)^n = (y|i \in [n])$  and:

$$\text{copy}_Y^n := \begin{cases} \text{copy}_Y^{n-1}(\text{id}_{Y^{n-2}} \otimes \text{copy}_Y) & n > 2 \\ \text{copy}_Y & n = 2 \end{cases} \quad (13)$$

$$(\text{copy}_Y^2 \mathbb{K}^2)(z|y) = \mathbb{K}^2(z|y, y) \quad (14)$$

(15)

Suppose for induction

$$(\text{copy}_Y^{n-1} \mathbb{K}^{n-1})(z|y) = \mathbb{K}^{n-1}(z|(y)^{n-1}) \quad (16)$$

then

$$(\text{copy}_Y^n \mathbb{K}^n)(z|y) = (\text{copy}_Y^{n-1}(\text{id}_{Y_{n-2}} \otimes \text{copy}_Y) \mathbb{K}^n)(z|y) \quad (17)$$

$$= \sum_{y' \in Y^{n-1}} (\text{id}_{Y^{n-2}} \otimes \text{copy}_Y)(\mathbf{y}'|(y)^{n-1}) \mathbb{K}^n(z|\mathbf{y}') \quad (18)$$

$$= \mathbb{K}^n(z|(y)^n) \quad (19)$$

That is, we can define the  $n$ -fold copy map that passes along  $n$  copies of its input to the next kernel in the product.

### 2.1.2 Example: comb insertion

The following examples illustrate 2-combs and the insertion operation, both of which we will define later. As an example in translating diagrams, we show how the diagrams for a 2-comb and 2-comb with an inserted Markov kernel can be translated to integral notation.

Consider the Markov kernels  $\mathbb{K} : W \rightarrow X$ ,  $\mathbb{L} : X \times W \times Y \rightarrow Z$  and the 2-comb  $\mathbb{M} : W \times Y \rightarrow X \times Z$  defined as

$$\mathbb{M} = \begin{array}{c} W \\ Y \end{array} \xrightarrow{\quad \text{K} \quad} X \quad \text{and} \quad \begin{array}{c} W \\ Y \end{array} \xrightarrow{\quad \text{L} \quad} Z \quad (20)$$

Following the rules above, we can translate this to ordinary notation by first breaking it down into products and tensor products, and then evaluating these products

$$\mathbb{M}(A \times B|w, y) = [(\text{copy}_W \otimes \text{id}_Y)(\mathbb{K} \otimes \text{id}_{W \times Y}) \quad (21)$$

$$(\text{copy}_X \otimes \text{id}_{W \times Y})(\text{id}_X \otimes \mathbb{L})](A \times B|w, y) \quad (22)$$

$$= [(\mathbb{K} \otimes \text{id}_{W \times Y})(\text{copy}_X \otimes \text{id}_{W \times Y}) \quad (23)$$

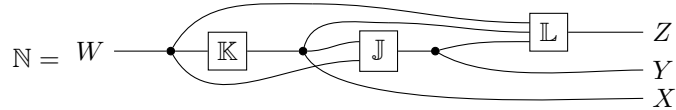
$$(\text{id}_X \otimes \mathbb{L})](A \times B|w, w, y) \quad (24)$$

$$= \int_X (\text{id}_X \otimes \mathbb{L})(A \times B|x', w, y) \mathbb{K}(dx'|w)(y, z|y', x) \quad (25)$$

$$= \int_X \text{id}_X(A|x') \mathbb{L}(B|x', w, y) \mathbb{K}(dx'|w) \quad (26)$$

$$= \int_A \mathbb{L}(B|x', w, y) \mathbb{K}(dx'|w) \quad (27)$$

If we are given additionally  $\mathbb{J} : X \times W \rightarrow Y$ , we can define a new Markov kernel  $\mathbb{N} : W \rightarrow Z$  given by “inserting”  $\mathbb{J}$  into  $\mathbb{M}$ :



$$\mathbb{N} = W \text{ --- } \text{[Diagram]} \quad (28)$$

We can translate Equation 28 to

$$\mathbb{N}(A \times B \times C|w) = [\text{copy}_W(\mathbb{K} \text{copy}_Y^3 \otimes \text{id}_W) \quad (29)$$

$$(\text{id}_Y \otimes \mathbb{J} \otimes \text{id}_Y)(\text{id}_Y \otimes \text{copy}_X \otimes \text{id}_Y) \quad (30)$$

$$(\mathbb{L} \otimes \text{id}_X \otimes \text{id}_Y)](A \times B \times C|w) \quad (31)$$

$$= [(\mathbb{K} \text{copy}_Y^3 \otimes \text{id}_W)(\text{id}_Y \otimes \mathbb{J} \otimes \text{id}_Y) \quad (32)$$

$$(\text{id}_Y \otimes \text{copy}_X \otimes \text{id}_Y) \quad (33)$$

$$(\mathbb{L} \otimes \text{id}_X \otimes \text{id}_Y)](A \times B \times C|w, w) \quad (34)$$

$$= \int_X \int_Y \mathbb{L}(C|x', w, y') \text{id}_X(A|x') \text{id}_Y(B|y') \mathbb{J}(dy'|x', w) \mathbb{K}(dx'|w) \quad (35)$$

$$= \int_A \int_B \mathbb{L}(C|x', w, y') \mathbb{J}(dy'|x', w) \mathbb{K}(dx'|w) \quad (36)$$

## 2.2 Semantics of observed and unobserved variables

We are interested in constructing *probabilistic models* which explain some part of the world. In a model, variables play the role of “pointing to the parts of the world the model is explaining”. Both observed and unobserved variables play important roles in causal modelling and we think it is worth clarifying what variables of either type refer to. Our approach is a standard one: a probabilistic

model is associated with an experiment or measurement procedure that yields values in a well-defined set. Observable variables are obtained by applying well-defined functions to the result of this total measurement. We use a richer sample space that includes “unobserved variables” that are formally treated the same way as observed variables, but aren’t associated with any real-world counterparts.

Consider Newton’s second law in the form  $\mathcal{F} = \mathcal{M}\mathcal{A}$  as a simple example of a model that relates variables  $\mathcal{F}$ ,  $\mathcal{M}$  and  $\mathcal{A}$ . As Feynman (1979) noted, this law is incomplete – in order to understand it, we must bring some pre-existing understanding of force, mass and acceleration as independent things. Furthermore, the nature of this knowledge is somewhat peculiar. Acknowledging that physicists happen to know a great deal about forces on an object, it remains true that in order to actually say what the net force on a real object is, even a highly knowledgeable physicist will still have to go and do some measurements, and the result of such measurements will be a vector representing the net forces on that object.

This suggests that we can think about “force”  $\mathcal{F}$  (or mass or acceleration) as a kind of procedure that we apply to a particular real world object and which returns a mathematical object (in this case, a vector). We will call  $\mathcal{F}$  a *procedure*. Our view of  $\mathcal{F}$  is akin to Menger (2003)’s notion of variables as “consistent classes of quantities” that consist of pairing between real-world objects and quantities of some type. Force  $\mathcal{F}$  itself is not a well-defined mathematical thing, as measurement procedures are not mathematically well-defined. At the same time, the set of values it may yield *are* well-defined mathematical things.

We will assume that any procedure will eventually yield an unambiguous value in a defined mathematical set. No actual procedure can be guaranteed to have this property – any apparatus, however robust, could suffer catastrophic failure – but we assume that we can study procedures reliable enough that we don’t lose much by making this assumption. This assumption allows us to say a procedure  $\mathcal{B}$  yields values in  $B$ .  $\mathcal{B} \bowtie x$  is the proposition that  $\mathcal{B}$ , when completed, yields the value  $x \in B$ , and by assumption exactly one of these propositions is true. For  $A \subset B$ ,  $\mathcal{B} \bowtie A$  is the proposition  $\bigvee_{x \in A} \mathcal{B} \bowtie x$ . Two procedures  $\mathcal{B}$  and  $\mathcal{C}$  are the same if  $\mathcal{B} \bowtie x \iff \mathcal{C} \bowtie x$  for all  $x \in B$ .

The notion of “yielding values” allows us to define an operation akin to function composition. If I have a procedure  $\mathcal{B}$  that takes values in some set  $B$ , and a function  $f : B \rightarrow C$ , define the “composition”  $f \circ \mathcal{B}$  to be the procedure  $\mathcal{C}$  that yields  $f(x)$  whenever  $\mathcal{B}$  yields  $x$ . For example,  $\mathcal{M}\mathcal{A}$  is the composition of  $h : (x, y) \mapsto xy$  with the procedure  $(\mathcal{M}, \mathcal{A})$  that yields the mass and acceleration of the same object. Composition is associative - for all  $x \in B$ :

$$(g \circ f) \circ \mathcal{B} \text{ yields } x \iff \mathcal{B} \text{ yields } (g \circ f)^{-1}(x) \quad (37)$$

$$\iff \mathcal{B} \text{ yields } f^{-1}(g^{-1}(x)) \quad (38)$$

$$\iff f \circ \mathcal{B} \text{ yields } g^{-1}(x) \quad (39)$$

$$\iff g \circ (f \circ \mathcal{B}) \text{ yields } x \quad (40)$$



One might wonder whether there is also some kind of “append” operation that takes a standalone  $\mathcal{M}$  and a standalone  $\mathcal{A}$  and returns a procedure  $(\mathcal{M}, \mathcal{A})$ . Unlike function composition, this would be an operation that acts on two procedures rather than a procedure and a function. Rather than attempt to define any operation of this type, we simply assume that somehow a procedure has been devised that measures everything of interest, which we will call  $\mathcal{S}$  which takes values in  $\Psi$ . We assume  $\mathcal{S}$  is such that any procedure of interest can be written as  $f \circ \mathcal{S}$  for some  $f$ .

For the model  $\mathcal{F} = \mathcal{M}\mathcal{A}$ , for example, we could assume  $\mathcal{F} = f \circ \mathcal{S}$  for some  $f$  and  $(\mathcal{M}, \mathcal{A}) = g \circ \mathcal{S}$  for some  $g$ . In this case, we can get  $\mathcal{M}\mathcal{A} = h \circ (\mathcal{M}, \mathcal{A}) = (h \circ g) \circ \mathcal{S}$ . Note that each procedure is associated with a unique function with domain  $\Psi$ .

Thus far,  $\Psi$  is a “sample space” that only contains observable variables. To include unobserved variables, we posit a richer sample space  $\Omega$  such that the measurement  $\mathcal{S}$  determines an element of some partition of  $\Omega$  rather than an element of  $\Omega$  itself. Then, by analogy to procedures defined with respect to  $\mathcal{S}$ , we identify variables in general with measurable functions defined on the domain  $\Omega$ .

Specifically, suppose  $\mathcal{S}$  takes values in  $\Psi$ . Then we can propose a sample space  $\Omega$  such that  $|\Omega| \geq |\Psi|$  and a surjective function  $S : \Omega \rightarrow \Psi$  associated with  $\mathcal{S}$ . We connect  $\Omega$ ,  $S$  and  $\mathcal{S}$  with the notion of *consistency with observation*:

$$\omega \in \Omega \text{ is consistent with observation iff the result yielded by } \mathcal{S} \text{ is equal to } S(\omega) \quad (41)$$

Thus the procedure  $\mathcal{S}$  eventually restricts the observationally consistent elements of  $\Omega$ . If  $\mathcal{S}$  yield the result  $s$ , then the consistent values of  $\Omega$  will be  $S^{-1}(s)$ .

One thing to note in this setup is that two different sets of measurement outcomes  $\Psi$  and  $\Psi'$  entail a different measurement procedures  $\mathcal{S}$  and  $\mathcal{S}'$ , but different sample spaces  $\Omega$  and  $\Omega'$  may be used to model a single procedure  $\mathcal{S}$ . We will sometimes consider different models of the same observable procedures.

As far as we know, distinguishing variables from procedures is somewhat nonstandard, but it is a useful distinction to make. While they may not be explicitly distinguished, both variables and procedures are often discussed in statistical texts. For example, Pearl (2009) offers the following two, purportedly equivalent, definitions of variables:

By a *variable* we will mean an attribute, measurement or inquiry that may take on one of several possible outcomes, or values, from a specified domain. If we have beliefs (i.e., probabilities) attached to the possible values that a variable may attain, we will call that variable a random variable.

This is a minor generalization of the textbook definition, according to which a random variable is a mapping from the sample space

(e.g., the set of elementary events) to the real line. In our definition, the mapping is from the sample space to any set of objects called “values,” which may or may not be ordered.

Our view is that the first definition is a definition of a procedure, while the second is a definition of a variable. Variables model procedures, but they are not the same thing. We can establish this by noting that, under our definition, every procedure of interest – that is, all procedures that can be written  $f \circ S$  for some  $f$  – is modeled by a variable, but there may be variables defined on  $\Omega$  that do not factorise through  $S$ , and these variables do not model procedures.

### 2.3 Events

To recap, we have a procedure  $S$  yielding values in  $\Psi$  that measures everything we are interested in, a sample space  $\Omega$  and a function  $S$  that models  $S$  in the sense of Definition 41. We assume also that  $\Psi$  has a  $\sigma$ -algebra  $\mathcal{E}$  (this may be the power set of  $\Psi$ , as measurement procedures are typically limited to finite precision).  $\Omega$  is equipped with a  $\sigma$ -algebra  $\mathcal{F}$  such that  $\sigma(S) \subset \mathcal{F}$ . If a procedure  $\mathcal{X} = f \circ S$  then we define  $X : \Omega \rightarrow X$  by  $X := f \circ S$ .

If a particular procedure  $\mathcal{X} = f \circ S$  eventually yields a value  $x$ , then the values of  $\Omega$  consistent with observation must be a subset of  $X^{-1}(x)$ . We define an *event*  $X \bowtie x \equiv X^{-1}(x)$ , which we read “the event that  $X$  yields  $x$ ”. An event  $X \bowtie x$  occurs if the consistent values of  $\Omega$  are a subset of  $X \bowtie x$ , thus “the event that  $X$  yields  $x$  occurs  $\equiv \mathcal{X}$  yields  $x$ ”. The definition of events applies to all types of variables, not just observables, but we only provide an interpretation of events “occurring” when the variable  $X$  is associated with some  $\mathcal{X}$ .

For measurable  $A \in \mathcal{X}$ ,  $X \bowtie A = \bigcup_{x \in A} X \bowtie x$ .

Given  $Y : \Omega \rightarrow X$ , we can define an append operation for variables:  $(X, Y) := \omega \mapsto (X(\omega), Y(\omega))$ .  $(X, Y)$  has the property that  $(X, Y) \bowtie (x, y) = X \bowtie x \cap Y \bowtie y$ , which supports the interpretation of  $(X, Y)$  as the values yielded by  $X$  and  $Y$  together.

It is common to use the symbol  $=$  instead of  $\bowtie$ , but we want to avoid this because  $Y = y$  already has a meaning, namely that  $Y$  is a constant function everywhere equal to  $y$ .

### 2.4 Probabilistic models for causal inference

The sample space  $(\Omega, \mathcal{F})$  along with our collection of variables is a “model skeleton” – it tells us what kind of data we might see. The process  $S$  which tells us which part of the world we’re interested in is related to the model  $\Omega$  and the observable variables by the criterion of *consistency with observation*. The kind of problem we are mainly interested in here is one where we make use of data to help make decisions under uncertainty. Probabilistic models have a long history of being used for this purpose, and our interest here is in constructing probabilistic models that can be attached to our variable “skeleton”.

For causal inference, we find we need to generalise the standard approach to constructing probability models on a given sample space  $(\Omega, \mathcal{F})$ . The key things we need to handle are *gaps* in our model. Hájek (2003) defines *probability gaps* as propositions that do not have a probability assigned to them. Our view of probability gaps is slightly different – in this work, a model with probability gaps as one that is missing some key parts. If we complete such a model with parts of the appropriate type, we get a standard probability model.

Probability gaps are useful in models used for decision making because, when I have a number of different options I could choose, a model can only help select from them if it tells me what is likely to happen for each choice I could make. Thus if we have a variable representing choices, we must have a model that can tolerate a probability gap for this variable.

As an initial example of a probability gap in causal inference, we will consider the example of truncated factorisation. For this example, we will assume that the reader is familiar enough with marginal probabilities, conditional probabilities and causal models to follow along. We will offer more careful definitions of terms later.

Suppose we have a causal Bayesian network  $(\mathbb{P}^{\mathbf{XYZ}}, \mathcal{G})$  where  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbf{Y} : \Omega \rightarrow Y$  and  $\mathbf{Z} : \Omega \rightarrow Z$  are variables,  $\mathbb{P}^{\mathbf{XYZ}}$  is a probability measure on  $X \times Y \times Z$  that we call “a probability model of  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ ” and  $\mathcal{G}$  is a Directed Acyclic Graph whose vertices we identify with  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{Z}$  that contains the edges  $\mathbf{X} \rightarrow \mathbf{Y}$  and  $\mathbf{X} \leftarrow \mathbf{Z} \rightarrow \mathbf{Y}$ . “Setting  $\mathbf{X}$  to  $x$ ” is an operation that takes as inputs  $\mathbb{P}^{\mathbf{XYZ}}$ ,  $\mathcal{G}$  and some  $x \in X$  and returns a new probability measure  $\mathbb{P}_x^{\mathbf{XYZ}}$  on  $X \times Y \times Z$  given by (Pearl, 2009, page 24):

$$\mathbb{P}_x^{\mathbf{XYZ}}(x', y, z) = \mathbb{P}^{\mathbf{Y|XZ}}(y|x, z) \mathbb{P}^{\mathbf{Z}}(z) \llbracket x = x' \rrbracket \quad (42)$$

Equation 42 embodies three assumptions about a model of the operation of “setting  $\mathbf{X}$  to  $x$ ”. First, such a model must assign probability 1 to the proposition that  $\mathbf{X}$  yields  $x$ . Second, such a model must assign the same marginal probability distribution to  $\mathbf{Z}$  as the input distribution;  $\mathbb{P}^{\mathbf{Z}} = \mathbb{P}_x^{\mathbf{Z}}$ . Finally, our model must also assign the same conditional probability to  $\mathbf{Y}$  given  $\mathbf{X}$  and  $\mathbf{Z}$ ;  $\mathbb{P}^{\mathbf{Y|XZ}} = \mathbb{P}_x^{\mathbf{Y|XZ}}$ .

Notice that the map  $x \mapsto \mathbb{P}_x^{\mathbf{XYZ}}$  itself “looks like” a conditional probability. It maps each  $x \in X$  to a probability distribution over  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ . In fact, a popular alternative notation for  $\mathbb{P}_x^{\mathbf{XYZ}}$  this map is  $\mathbb{P}^{\mathbf{XYZ}|do(\mathbf{X}=x)}$ , which is clearly suggestive of an interpretation as a kind of conditional probability. We will take this interpretation seriously: we will posit some variable  $\mathbf{U}$  (which may or may not be observable) and a probabilistic model  $\mathbb{Q}^{\mathbf{XYZ|U}} := x \mapsto \mathbb{P}_x^{\mathbf{XYZ}}$ .

We note that the thing called  $do(\mathbf{X} = x)$  or that we have called  $\mathbf{U}$  is often referred to as an *intervention*. Interventions are often things we can choose to do or not to do. Also, perhaps, we could also consider choosing to do or not to do an intervention based on the output of some random process. In this case, we will need a model that can tell us which result we are likely to see for any choice of  $\mathbb{Q}_\alpha^{\mathbf{U}}$ ; the distribution of  $\mathbf{U}$  is a *probability gap*.

$\mathbb{Q}^{\mathbf{XYZ|U}}$ , as we have defined it so far, is not quite an ideal candidate for a gappy probability model. Firstly, because conditional probabilities are arbitrary

on sets of measure zero with regard to  $\mathbb{P}^{XYZ}$ , definition 42 can be satisfied by multiple probability distributions that differ in meaningful ways. Suppose  $X$ ,  $Y$  and  $Z$  are binary,  $\mathbb{P}^Z(1) = 1$  and  $\mathbb{P}(X \bowtie Z) = 1$ . Then we can consistently choose  $\mathbb{P}^{Y|XZ}(1|0, 1) = 1$  or  $\mathbb{P}^{Y|XZ}(1|0, 1) = 0$  because  $\{0, 1\}$  is a measure zero event. However, the first choice gives us  $\mathbb{P}_0^{XYZ}(0, 1, 1) = 1$  while the second gives us  $\mathbb{P}_0^{XYZ}(0, 1, 1) = 0$ , which are very different opinions regarding “the result of setting  $X$  to 1”.

Secondly, there may be no probability model at all that satisfies Equation 42. For example, suppose  $X = f \circ Z$  for some  $f$ . Then we must have  $\mathbb{P}_x^X(x') = \mathbb{P}_x^Z(f^{-1}(x'))$  for any  $x$ . However, we also have  $\mathbb{P}_x^X(x') = \llbracket x = x' \rrbracket$  for all  $x, x'$  and  $\mathbb{P}_x^Z = \mathbb{P}^Z$  for all  $x$ . Thus if  $X$  can more than one value, there is at least one choice of  $x$  that cannot simultaneously satisfy these requirements.

This might seem like an absurd model, but an analogous causal graph appears in Shahar (2009) where  $Z = (H, W)$ , representing a person’s height and weight, and  $X$  represents their body mass index, which is to say  $X = \frac{W}{H^2}$ . Furthermore, this paper uses this model to argue that body mass index cannot have a causal effect. Such an argument cannot be supported by Equation 42 because by that equation there is no probability model corresponding to an intervention on  $X$ .

Our first aim is to offer a more careful theory of probability models with gaps in them, addressing these two shortcomings.

## 2.5 Order 0 gaps: probability distributions

At this point, we will make a substantial simplifying assumption: all sets, including the sample space  $\Omega$  and any set of values a variable takes, are discrete sets. That is, they are at most countably infinite and the  $\sigma$ -algebra of measurable sets is the power set.

Define  $I : \Omega \rightarrow \Omega$  as the identity function  $\omega \mapsto \omega$ . A *standard probability model* on  $\Omega$  is a probability distribution  $\mathbb{P}^I$  on  $\Omega$ . This terminology is motivated by the fact that probability models are often given by a collection of random variables along with a probability space, which is a triple of the form  $(\mathbb{P}^I, \Omega, \mathcal{F})$  where  $\mathcal{F}$  is a  $\sigma$ -algebra on  $\Omega$ .

A standard probability model has no gaps – it associates a probability with every event  $X \bowtie A$  for every  $X$  and every  $A \in \mathcal{X}$ , and this probability is subject to no conditions.

Given a standard probability model  $\mathbb{P}^I$  and  $X : \Omega \rightarrow X$ , the probability of the event  $X \bowtie A$  for  $A \in \mathcal{X}$  is defined as  $\mathbb{P}^X(A) := \mathbb{P}^I(X \bowtie A)$ .  $\mathbb{P}^X$  is known as the pushforward of  $\mathbb{P}^I$  via  $X$ . We use the convention that the same base letter  $\mathbb{P}$  appearing in  $\mathbb{P}^I$  and  $\mathbb{P}^X$  indicates  $\mathbb{P}^X$  is a pushforward of  $\mathbb{P}^I$ . We say that  $\mathbb{P}^X$  is the distribution of  $X$  *under*  $\mathbb{P}^I$ .

Any *sample space valid* probability distribution  $\mathbb{Q}^X$  the distribution of  $X$  under some  $\mathbb{Q}^I$ . A sample space valid probability distribution is a probability model with a 0th-order gap. Unless  $X$  is invertible it contains gaps. The distinguishing feature of a 0th-order gap model is that specifying  $\mathbb{Q}^X$  is equivalent to specifying  $X$  and any  $\mathbb{Q}^I$  that pushes forward to  $\mathbb{Q}^X$ .



such that, given  $x \neq x' \in X$ ,  $\mathbb{P}^{Y|X}(x|x') = a > 0$ . Take  $\mathbb{P}_\alpha^X(x') = 1$  and then  $\mathbb{P}_\alpha^{XY}((x', x)) = a > 0$  also. For every probability model  $\mathbb{Q}^I \in \Delta(\Omega)$ ,

$$\mathbb{Q}^{XY}(x, x') = \mathbb{Q}^I(X \bowtie x \cap X \bowtie x') \quad (48)$$

$$= \llbracket x = x' \rrbracket \quad (49)$$

Thus  $\mathbb{P}_\alpha^{XY}$  is not sample space valid. Furthermore, it is surely absurd to suppose that a variable yields one value and then assign positive probability to the same variable yielding a different value. This is why we propose the validity condition in our definition of conditional probabilities.

The key justification for our definition of validity is that valid conditional probabilities guarantee that we can extend the conditional probability to a standard probability model. In fact, Definition 2.3 is equivalent to defining valid conditional probabilities as those that extend to valid probability distributions, given any valid probability distribution (Theorem 2.8). In addition, the definition of sample space validity itself is equivalent to Definition 2.3 under the identification of probability distributions with trivial conditional probabilities (Theorem 2.5).

**Theorem 2.5** (Agreement of validity criteria for probability distributions). *Given  $X : \Omega \rightarrow X$ , with  $\Omega$  and  $X$  discrete, a probability  $\mathbb{P}^X$  is sample space-valid if and only if the conditional probability  $\mathbb{P}^{X|*} := * \mapsto \mathbb{P}^X$  is valid.*

*Proof.*  $* \bowtie * = \Omega$  necessarily. Thus validity of  $\mathbb{P}^{X|*}$  means

$$\forall x \in X : X \bowtie (x) = \emptyset \implies \mathbb{P}^{X|*}(x|*) = 0 \quad (50)$$

$$= \mathbb{P}^X(x) \quad (51)$$

If: We refer to Ershov (1975) Theorem 2.5 for the proof that Equation 51 is necessary and sufficient for the existence of  $\mathbb{P}^I$  such that  $\mathbb{P}^I(X^{-1}(A)) = \mathbb{P}^X(A)$  for all  $A \in \mathcal{X}$  when  $(\Omega, \mathcal{F})$  and  $(X, \mathcal{X})$  are standard measurable. If  $\Omega$  and  $X$  are discrete, then they are standard measurable.

Only if: If  $X \bowtie x = \emptyset$  then  $\mathbb{P}^X(x) = \mathbb{P}^I(\emptyset) = 0$ .  $\square$

**Lemma 2.6** (Valid conditional probabilities are validly extendable). *Given  $(\Omega, \mathcal{F})$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ ,  $Z : \Omega \rightarrow Z$  and any valid conditional probabilities  $\mathbb{P}^{Y|X}$  and  $\mathbb{Q}^{Z|YX}$ ,  $\mathbb{P}^{Y|X} \odot \mathbb{Q}^{Z|YX}$  is also a valid conditional probability.*

*Proof.* Let  $\mathbb{R}^{YZ|X} := \mathbb{P}^{Y|X} \odot \mathbb{Q}^{Z|YX}$ .

We only need to check validity in  $x \in X(\Omega)$ , as it is automatically satisfied for other values of  $X$ .

For all  $x \in X(\Omega)$ ,  $X \bowtie x \cap Y \bowtie y = \emptyset$ ,  $\mathbb{P}^{Y|X}(y|x) = 0$  by validity. Thus

$$\mathbb{R}^{YZ|X}(y, z|x) = \mathbb{Q}^{Z|YX}(z|y, x) \mathbb{P}^{Y|X}(y|x) \quad (52)$$

$$\leq \mathbb{P}^{Y|X}(y|x) \quad (53)$$

$$= 0 \quad (54)$$

For all  $(x, y) \in (\mathbf{X}, \mathbf{Y})(\Omega)$ ,  $z \in Z$  such that  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \bowtie (x, y, z) = \emptyset$ ,  $\mathbb{Q}^{Z|\mathbf{Y}\mathbf{X}}(z|y, x) = 0$  by validity. Thus for any such  $(x, y, z)$ :

$$\mathbb{R}^{YZ|\mathbf{X}}(y, z|x) = \mathbb{Q}^{Z|\mathbf{Y}\mathbf{X}}(z|y, x)\mathbb{P}^{\mathbf{Y}|\mathbf{X}}(y|x) \quad (55)$$

$$= 0 \quad (56)$$

□

**Corollary 2.7** (Valid conditional probability is validly extendable to a probability distribution). *Given  $\Omega$ ,  $\mathbf{U} : \Omega \rightarrow U$ ,  $\mathbf{W} : \Omega \rightarrow W$  and a valid conditional probability  $\mathbb{T}^{\mathbf{W}|\mathbf{U}}$ , then for any valid conditional probability  $\mathbb{V}^{\mathbf{U}}$ ,  $\mathbb{V}^{\mathbf{U}} \odot \mathbb{T}^{\mathbf{W}|\mathbf{U}}$  is a valid probability distribution.*

*Proof.* Applying Lemma 2.6 choosing  $\mathbf{X} = *$ ,  $\mathbf{Y} = \mathbf{U}$ ,  $\mathbf{Z} = \mathbf{W}$  and  $\mathbb{P}^{\mathbf{Y}|\mathbf{X}} = \mathbb{V}^{\mathbf{U}|*}$  and  $\mathbb{Q}^{Z|\mathbf{Y}\mathbf{X}} = \mathbb{T}^{\mathbf{W}|\mathbf{U}*}$  we have  $\mathbb{R}^{WU|*} := \mathbb{V}^{\mathbf{U}|*} \odot \mathbb{T}^{\mathbf{W}|\mathbf{U}*}$  is a valid conditional probability. Then  $\mathbb{R}^{\mathbf{WU}} \cong \mathbb{R}^{WU|*}$  is valid by Theorem 2.5. □

**Theorem 2.8** (Validity of conditional probabilities). *Suppose we have  $\Omega$ ,  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbf{Y} : \Omega \rightarrow Y$ , with  $\Omega$ ,  $X$ ,  $Y$  discrete. A conditional probability  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is valid if and only if for all valid probability distributions  $\mathbb{V}^{\mathbf{X}}$ ,  $\mathbb{V}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is a valid probability distribution.*

*Proof.* If: this follows directly from Lemma 2.6.

Only if: suppose  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is invalid. Then there is some  $x \in X$ ,  $y \in Y$  such that  $\mathbf{X} \bowtie (x) \neq \emptyset$ ,  $(\mathbf{X}, \mathbf{Y}) \bowtie (x, y) = \emptyset$  and  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}(y|x) > 0$ . Choose  $\mathbb{V}^{\mathbf{X}}$  such that  $\mathbb{V}^{\mathbf{X}}(\{x\}) = 1$ ; this is possible due to standard measurability and valid due to  $\mathbf{X}^{-1}(x) \neq \emptyset$ . Then

$$(\mathbb{V}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}})(x, y) = \mathbb{T}^{\mathbf{Y}|\mathbf{X}}(y|x)\mathbb{V}^{\mathbf{X}}(x) \quad (57)$$

$$= \mathbb{T}^{\mathbf{Y}|\mathbf{X}}(y|x) \quad (58)$$

$$> 0 \quad (59)$$

Hence  $\mathbb{V}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is invalid. □

Given any two conditional probabilities  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}$ ,  $\mathbb{U}^{\mathbf{Y}|\mathbf{X}}$  such that there is some  $x \in \mathbf{X}(\Omega)$ ,  $A \in \mathcal{Y}$  for which  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}(A|x) \neq \mathbb{U}^{\mathbf{Y}|\mathbf{X}}(A|x)$ , there exists some  $\mathbb{Q}^{\mathbf{X}}$  such that  $\mathbb{Q}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}} \neq \mathbb{Q}^{\mathbf{X}} \odot \mathbb{U}^{\mathbf{Y}|\mathbf{X}}$ . We can see this by letting  $\mathbb{Q}^{\mathbf{X}}$  assign probability 1 to some value of  $x$  on which  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  and  $\mathbb{U}^{\mathbf{Y}|\mathbf{X}}$  differ. Thus if we are using conditional probability as a probability model with a gap, uniqueness requires that the conditional probability be represented by a Markov kernel that is unique up to the set of impossible values of the conditioning variable. In contrast, conditional probability derived from a standard probability model may be represented by a Markov kernel unique up to a measure zero set, which is a strictly weaker condition because the set of impossible values must be given measure 0 by any valid probability distribution.

## 2.7 Order 2 gaps: probability combs

So far we have order-0 and order-1 gap models. Extending an order-1 gap model with an order-0 model yields an order-0 model. We can continue to higher order models, such that an order-2 gap model can be extended with an order-1 gap model to yield an order-1 gap model. Order-2 gap models are *probability 2-combs* (Chiribella et al., 2008; Jacobs et al., 2019). We will provide a working definition of probability 2-combs as a pair of conditional probabilities, and refine this definition later.

We may be able to extend this theory to  $n$ -combs, where valid  $n$ -combs are those that can be extended to valid  $n - 1$ -combs by valid  $n - 1$ -combs. 2-combs are sufficient for our purposes, though.

**Definition 2.9** (Probability 2-comb). Given  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ ,  $Z : \Omega \rightarrow Z$ , a probability 2-comb  $\mathbb{P}^{X|W \square Z|Y} : W \times Y \rightarrow X \times Z$  is a Markov kernel such that for some  $\mathbb{P}^{X|W} : W \rightarrow X$

$$\begin{array}{c} W \\ Y \end{array} \begin{array}{|c|} \hline \mathbb{P}^{X|W \square Z|Y} \\ \hline \end{array} \begin{array}{c} X \\ * \end{array} = \begin{array}{c} W \\ Y \end{array} \begin{array}{|c|} \hline \mathbb{P}^{X|W} \\ \hline * \end{array} \begin{array}{c} X \\ * \end{array} \quad (60)$$

Coupled with the validity conditions

1.  $\mathbb{P}^{X|W}$  is valid as a conditional probability
2.  $(W, X, Y, Z) \bowtie (w, x, y, z) = \emptyset$  implies  $\mathbb{P}^{X|W \square Z|Y}(x, z|w, y) = 0$  or  $(W, X, Y) \bowtie (w, x, y) = \emptyset$

With discrete sets, and in general wherever we have kernel disintegrations, there exists some  $\bar{\mathbb{P}}^{Z|WXY} : W \times X \times Y \rightarrow Z$  (Lemma 2.15) such that

$$\mathbb{P}^{X|W \square Z|Y} = \begin{array}{c} W \\ Y \end{array} \begin{array}{|c|} \hline \mathbb{P}^{X|W} \\ \hline \end{array} \begin{array}{c} X \\ * \end{array} \begin{array}{|c|} \hline \bar{\mathbb{P}}^{Z|WXY} \\ \hline \end{array} \begin{array}{c} Z \\ * \end{array} \quad (61)$$

We define the operation insert that takes a 2-comb and a conditional probability and returns a conditional probability.

$$\text{insert}(\mathbb{P}_\alpha^{Y|XW}, \mathbb{P}^{X|W \square Z|Y}) = \mathbb{P}^{X|W} \odot \mathbb{P}_\alpha^{Y|XW} \odot \bar{\mathbb{P}}^{Z|WXY} \quad (62)$$

We can depict this operation graphically in a somewhat informal way as “inserting”  $\mathbb{P}_\alpha^{Y|XW}$  into  $\mathbb{P}^{X|W \square Z|Y}$ :



$$\text{Insert} \quad \left| \begin{array}{c} \text{ } \\ \text{ } \\ \text{ } \end{array} \right| \quad (63)$$

$$= \text{Diagram (64)} \quad (64)$$

Just as in Section 2.6, a key concern for probability 2-combs is *valid extendability*. The insert operation we define for 2-combs is uniquely defined by Equation 62 (Theorem 2.10, and yields a valid conditional probability given any valid inputs (Theorem 2.11).

**Theorem 2.10** (Equivalence of 2-comb disintegrations). *Given a 2-comb  $\mathbb{P}^{X|W \square Z|Y}$  and any two disintegrations  $\bar{\mathbb{P}}_1^{Z|WXY}$ ,  $\bar{\mathbb{P}}_2^{Z|WXY}$ , for all valid extensions  $\mathbb{P}_a^{Y|XW}$*

$$\mathbb{P}^{X|W} \odot \mathbb{P}_\alpha^{Y|XW} \odot \bar{\mathbb{P}}_1^{Z|WXY} = \mathbb{P}^{X|W} \odot \mathbb{P}_\alpha^{Y|XW} \odot \bar{\mathbb{P}}_2^{Z|WXY} \quad (65)$$

*Proof.* For any  $w, x, y, z$

$$(\mathbb{P}^{X|W} \odot \mathbb{P}_\alpha^{Y|XW} \odot \bar{\mathbb{P}}_1^{Z|WXY})(x, y, z|w) = \mathbb{P}_\alpha^{Y|WX}(y|w, x) \mathbb{P}^{X|W}(x|w) \bar{\mathbb{P}}_1^{Z|XYW}(z|x, y, w) \quad (66)$$

$$= \mathbb{P}^{X|W \sqcup Z|Y}(x, z|w, y) \quad (67)$$

$$= \mathbb{P}_\alpha^{\mathbf{Y}|\mathbf{W}\mathbf{X}}(y|w, x) \mathbb{P}^{\mathbf{X}|\mathbf{W}}(x|w) \bar{\mathbb{P}}_2^{\mathbf{Z}|\mathbf{X}\mathbf{Y}\mathbf{W}}(z|x, y, w) \quad (68)$$

☐

**Theorem 2.11** (Extension of valid probability 2-combs). *Given  $\Omega$ ,  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and  $Z : \Omega \rightarrow Z$ , a probability 2-comb  $\mathbb{P}^{X|W \square Z|Y}$  is valid if and only if  $\text{insert}(\mathbb{P}_\alpha^{Y|WX}, \mathbb{P}^{X|W \square Z|Y})$  is valid for all valid  $\mathbb{P}_\alpha^{Y|WX}$ .*

*Proof.* Only if:

Note that

$$\mathbb{P}_\alpha^{\text{XYZ|W}} := \text{insert}(\mathbb{P}_\alpha^{\text{Y|WX}}, \mathbb{P}^{\text{X|W}\square\text{Z|Y}}) \quad (69)$$

$$\mathbb{P}_\alpha^{\text{XYZ}|W}(xyz|w) = \mathbb{P}_\alpha^{\text{Y}|WX}(y|w, x) \mathbb{P}^{\text{X}|W}(x|w) \bar{\mathbb{P}}^{\text{Z}|XYW}(z|x, y, w) \quad (70)$$

$$= \mathbb{P}^{X|W \square Z|Y}(x, z|w, y) \mathbb{P}_\alpha(y|x) \quad (71)$$

Suppose  $\mathbb{P}^{X|W \square Z|Y}$  is valid. If  $(W, X, Y, Z) \bowtie (w, x, y, z) = \emptyset$  then either  $\mathbb{P}^{X|W \square Z|Y}(x, z|w, y) = 0$  and hence  $\mathbb{P}_\alpha^{XYZ|W}(xyz|w) = 0$  or  $(W, X, Y) \bowtie (w, x, y) = \emptyset$ .

If  $(W, X, Y) \bowtie (w, x, y) = \emptyset$  then either  $(W, X) \bowtie (w, x) \neq \emptyset$  and by validity  $\mathbb{P}_\alpha^{Y|WX}(y|w, x) = 0$  and so  $\mathbb{P}_\alpha^{XYZ|W}(xyz|w) = 0$  or  $(W, X) \bowtie (w, x) = \emptyset$ .

If  $(W, X) \bowtie (w, x) = \emptyset$  then either  $W \bowtie w \neq \emptyset$  and by validity  $\mathbb{P}^{X|W}(x|w) = 0$  and so  $\mathbb{P}_\alpha^{XYZ|W}(xyz|w) = 0$  or  $W \bowtie w = \emptyset$ , in which case  $\mathbb{P}_\alpha^{XYZ|W}(xyz|w)$  may take any value.

If: Suppose  $\mathbb{P}^{X|W \square Z|Y}$  is invalid. Then either  $\mathbb{P}^{X|W}$  is invalid or  $\mathbb{P}^{X|W \square Z|Y}(x, z|w, y) > 0$  on some  $(w, x, y, z)$  such that  $(W, X, Y, Z) \bowtie (w, x, y, z) = \emptyset$  and  $(W, X, Y) \bowtie (w, x, y) \neq \emptyset$ .

Suppose  $\mathbb{P}^{X|W}$  is invalid. Then

$$\mathbb{P}_\alpha^{X|W}(x|w) = \sum_{y \in Y, z \in Z} \mathbb{P}_\alpha^{XYZ|W}(xyz|w) \quad (72)$$

$$= \mathbb{P}^{X|W}(x|w) \quad (73)$$

Thus  $\mathbb{P}_\alpha^{X|W}(x|w)$  is invalid and therefore so too is  $\mathbb{P}_\alpha^{XYZ|W}$ .

Suppose we have some  $(w, x, y, z)$  such that  $(W, X, Y, Z) \bowtie (w, x, y, z) = \emptyset$ ,  $(W, X, Y) \bowtie (w, x, y) \neq \emptyset$  and  $\mathbb{P}^{X|W \square Z|Y}(x, z|w, y) > 0$ .

By supposition, there is a valid  $\mathbb{P}_\alpha^{Y|WX}$  such that  $\mathbb{P}_\alpha^{Y|WX}(y|w, x) = 1$ . Then

$$\mathbb{P}_\alpha^{XYZ|W}(xyz|w) = \mathbb{P}^{X|W \square Z|Y}(x, z|w, y) \mathbb{P}_\alpha(y|x) \quad (74)$$

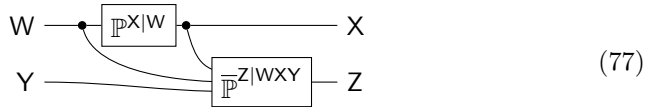
$$= \mathbb{P}^{X|W \square Z|Y}(x, z|w, y) \quad (75)$$

$$> 0 \quad (76)$$

So  $\mathbb{P}_\alpha^{XYZ|W}$  is invalid.  $\square$

It is also the case that if we combine two valid conditional probabilities to form a 2-comb, the result is a valid 2-comb.

**Theorem 2.12** (Valid conditional probabilities combine to a valid 2-comb).  
Given valid conditional probabilities  $\mathbb{P}^{X|W}$  and  $\mathbb{P}^{Z|WXY}$



Is a valid 2-comb.

*Proof.* Validity of  $\mathbb{P}^{X|W}$  is by assumption, and validity of  $\mathbb{P}^{Z|WXY}$  means  $(W, X, Y, Z) \bowtie (w, x, y, z) = \emptyset$  implies  $\mathbb{P}^{Z|WXY}(z|w, x, y) = 0$  or  $(W, X, Y) \bowtie (w, x, y) = \emptyset$ . This in turn implies  $\mathbb{P}^{Z|WXY}(z|w, x, y) = 0$  or  $(W, X, Y) \bowtie (w, x, y) = \emptyset$ .  $\square$

## 2.8 Revisiting truncated factorisation

We have established that, at least for discrete sets, probability 2-combs are a general kind of object that maps a conditional probability to a conditional probability. Thus they model choices that correspond to conditional probabilities, rather than probability distributions. We initially considered that a truncated factorisation might be a conditional probability – that is, it could be used to model choices that correspond to probability distributions of the choice variable  $U$  that we introduced. If we instead consider the possibility that the model of “interventions” that Equation 42 is suggesting is actually a probability 2-comb, we find that both of the problems we identified are automatically solved.

Suppose we have a 2-comb  $\mathbb{P}^{Z \square Y | X}$ . Then we have some  $\mathbb{P}^Z$ ,  $\bar{\mathbb{P}}^{Y | XZ}$  unique up to equivalence such that

$$\mathbb{P}^{Z \square Y | X}(z, y | x) = \bar{\mathbb{P}}^{Y | XZ}(y | x, z) \mathbb{P}^Z(z) \quad (78)$$

By Lemma 2.14, we also have

$$\mathbb{P}^{Z \square XY | X}(z, x', y | x) = \bar{\mathbb{P}}^{Y | XZ}(y | x, z) \mathbb{P}^Z(z) \llbracket x = x' \rrbracket \quad (79)$$

Compare this to Equation 42, which we reproduce here for convenience:

$$\mathbb{P}_x^{XYZ}(x', y, z) = \mathbb{P}^{Y | XZ}(y | x, z) \mathbb{P}^Z(z) \llbracket x = x' \rrbracket$$

These equations look almost identical, and suggest an alternative approach to modelling intervention. Instead of *defining* an interventional conditional probability via truncated factorisation, suppose we have an interventional model that is a probability 2-comb, and that satisfies

$$\mathbb{P}_{\text{obs}}^{XYZ} = \text{insert}(\mathbb{P}_{\text{obs}}^{X | Z}, \mathbb{P}^{Z \square XY | X}) \quad (80)$$

For some observational insert  $\mathbb{P}_{\text{obs}}^{X | Z}$ .

Under this view, we cannot necessarily derive  $\mathbb{P}^{Z \square XY | X}$  from the observational probability distribution  $\mathbb{P}_{\text{obs}}^{XYZ}$ , as we cannot guarantee that any particular  $\bar{\mathbb{P}}_{\text{obs}}^{Y | XZ}$  is a representative of  $\bar{\mathbb{P}}^{Y | XZ}$ . This is as it should be – if we have meaningfully different choices of  $\bar{\mathbb{P}}_{\text{obs}}^{Y | XZ}$  arising from measure zero sets, we cannot decide which one will model interventions without more information. The 2-comb  $\mathbb{P}^{Z \square XY | X}$  itself stands for a unique map from interventions to probability distributions, rather than the potentially large set of different maps of this type.

The second issue we raised was the nonexistence of  $\mathbb{P}_x^{XYZ}$  if, for example,  $X = f \circ Z$ . Because  $\mathbb{P}_{\text{obs}}^Z$  and  $\bar{\mathbb{P}}_{\text{obs}}^{Y | XZ}$  are both valid, Lemma 2.12 ensures that any 2-comb defined by any choice of these conditional probabilities is also valid, and so Equation 79 defines a valid object. Because we have a 2-comb which is extended

by *conditional probabilities* rather than by probability distributions, and these conditional probabilities *must be valid*, we can validly extend  $\mathbb{P}^{Z \square XY|X}$  even when  $X = f \circ Z$  – it's just that validity requires that any extending conditional must be such that  $\mathbb{P}_\alpha^{X|Z}(x|z) = \llbracket x = f(z) \rrbracket$ .

So: 2-combs do what we want truncated factorisation to do without breaking down in edge cases, except maybe when we go to uncountably infinite sets.

## 2.9 Useful results

### 2.9.1 Repeated variables

Lemmas 2.13 and 2.14 establish that models of repeated variables must connect the repetitions with a copy map.

**Lemma 2.13** (Output copies of the same variable are identical). *For any  $\Omega$ ,  $X, Y, Z$  random variables on  $\Omega$  and conditional probability  $\mathbb{K}^{YZ|X}$ , there is a conditional probability  $\mathbb{K}^{YYZ|X}$  unique up to impossible values of  $X$  such that*

$$X \text{ --- } \boxed{\mathbb{K}^{YYZ|X}} \begin{matrix} \text{---}^* \\ \text{---} \\ \text{---} \end{matrix} \begin{matrix} Y \\ Y \\ Z \end{matrix} = \mathbb{K}^{YZ|X} \quad (81)$$

and it is given by

$$\mathbb{K}^{YYZ|X} = X \text{ --- } \boxed{\mathbb{K}^{YZ|X}} \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \begin{matrix} Y \\ Y \\ Z \end{matrix} \quad (82)$$

$$\iff \quad (83)$$

$$\mathbb{K}^{YYZ|X}(y, y', z|x) = \llbracket y = y' \rrbracket \mathbb{K}^{YZ|X}(y, z|x) \quad (84)$$

$$(85)$$

*Proof.* If we have a valid  $\mathbb{K}^{YYZ|X}$ , it must be the pushforward of  $(Y, Y, Z)$  under some  $\mathbb{K}^{||X}$ . Furthermore,  $\mathbb{K}^{YZ|X}$  must be the pushforward of  $(*, Y, Z) \cong (Y, Z)$  under the same  $\mathbb{K}^{||X}$ .

For any  $x \in X(\Omega)$ , validity requires  $(X, Y, Y, Z) \bowtie (x, y, y', z) = \emptyset \implies \mathbb{K}^{YYZ|X}(y, y', z|x) = 0$ . Clearly, whenever  $y \neq y'$ ,  $\mathbb{K}^{YYZ|X}(y, y', z|x) = 0$ . Because  $\mathbb{K}^{YYZ|X}$  is a Markov kernel, there is some  $\mathbb{L} : X \rightarrow X \times Z$  such that

$$\mathbb{K}^{YYZ|X}(y, y', z|x) = \llbracket y = y' \rrbracket \mathbb{L}(y, z|x) \quad (86)$$

$$(87)$$

But then

$$\mathbb{K}^{YZ|X}(y, z|x) = \sum_{y' \in Y} \mathbb{K}^{YYZ|X}(y, y', z|x) \quad (88)$$

$$= \mathbb{L}(y, z|x) \quad (89)$$

$$(90)$$

□

**Lemma 2.14** (Copies shared between input and output are identical). *For any  $\mathbb{K} : (\mathbf{X}, \mathbf{Y}) \rightarrow (\mathbf{X}, \mathbf{Z})$ ,  $\mathbb{K}$  is a model iff there exists some  $\mathbb{L} : (\mathbf{X}, \mathbf{Y}) \rightarrow \mathbf{Z}$  such that*

$$\mathbb{K} = \begin{array}{c} \mathbf{X} \\ \mathbf{Y} \end{array} \begin{array}{c} \text{---} \bullet \text{---} \\ \text{---} \end{array} \boxed{\mathbb{K}^{\mathbf{Z}|\mathbf{X}\mathbf{Y}}} \begin{array}{c} \text{---} \mathbf{Z} \\ \text{---} \end{array} \quad (91)$$

$$\iff \quad (92)$$

$$\mathbb{K}_{x,y}^{x',z} = \llbracket x = x' \rrbracket \mathbb{L}_{x,y}^z \quad (93)$$

For any  $\Omega$ ,  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  random variables on  $\Omega$  and conditional probability  $\mathbb{K}^{\mathbf{Z}|\mathbf{X}\mathbf{Y}}$ , there is a conditional probability  $\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}$  unique up to impossible values of  $(\mathbf{X}, \mathbf{Y})$  such that

$$\begin{array}{c} \mathbf{X} \\ \mathbf{Y} \end{array} \boxed{\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}} \begin{array}{c} \text{---} * \\ \text{---} \mathbf{Z} \end{array} = \mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}} \quad (94)$$

and it is given by

$$\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}} = \mathbf{X} \text{---} \boxed{\mathbb{K}^{\mathbf{Y}\mathbf{Z}|\mathbf{X}}} \begin{array}{c} \text{---} \mathbf{Y} \\ \text{---} \mathbf{Z} \end{array} \quad (95)$$

$$\iff \quad (96)$$

$$\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}(x, z|x', y) = \llbracket x = x' \rrbracket \mathbb{K}^{\mathbf{Z}|\mathbf{X}\mathbf{Y}}(z|x', y) \quad (97)$$

$$(98)$$

*Proof.* If we have a valid  $\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}$ , it must be the pushforward of  $(\mathbf{X}, \mathbf{Z})$  under some  $\mathbb{K}^{||\mathbf{X}\mathbf{Y}}$ . Furthermore,  $\mathbb{K}^{\mathbf{Z}|\mathbf{X}\mathbf{Y}}$  must be the pushforward of  $(*, \mathbf{Z}) \cong (\mathbf{Z})$  under the same  $\mathbb{K}^{||\mathbf{X}}$ .

For any  $(x, y) \in (\mathbf{X}, \mathbf{Y})(\Omega)$ , validity requires  $(\mathbf{X}, \mathbf{Y}, \mathbf{X}, \mathbf{Z}) \bowtie (x, y, x', z) = \emptyset \implies \mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}(x', z|x, y) = 0$ . Clearly, whenever  $x \neq x'$ ,  $\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}(x', z|x, y) = 0$ . Because  $\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}$  is a Markov kernel, there is some  $\mathbb{L} : \mathbf{X} \times \mathbf{Y} \rightarrow \mathbf{Z}$  such that

$$\mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}(x', z|x, y) = 0 = \llbracket x = x' \rrbracket \mathbb{L}(z|x, y) \quad (99)$$

$$(100)$$

But then

$$\mathbb{K}^{\mathbf{Z}|\mathbf{X}\mathbf{Y}}(y, z|x) = \sum_{x' \in \mathbf{X}} \mathbb{K}^{\mathbf{X}\mathbf{Z}|\mathbf{X}\mathbf{Y}}(x', z|x, y) \quad (101)$$

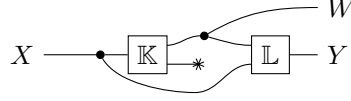
$$= \mathbb{L}(z|x, y) \quad (102)$$

$$(103)$$

□

### 2.9.2 Disintegrations

**Lemma 2.15** (Disintegration existence). *For any model  $\mathbb{K} : X \rightarrow W \times Y$  where  $X, W, Y$  are discrete, there exists  $\mathbb{L} : W \times X \rightarrow Y$  such that*



$$\mathbb{K} = \quad (104)$$

*Proof.* Consider any Markov kernel  $\mathbb{L} : W \times X \rightarrow Y$  with the property

$$\mathbb{L}(y|w, x) = \frac{\mathbb{K}(w, y|x)}{\sum_{y \in Y} \mathbb{K}(w, y|x)} \quad \forall x, w : \text{the denominator is positive} \quad (105)$$

Then

$$\sum_{y \in Y} \mathbb{K}(w, y|x) \mathbb{L}(y|w, x) = \sum_{y \in Y} \mathbb{K}(w, y|x) \frac{\mathbb{K}(w, y|x)}{\sum_{y \in Y} \mathbb{K}(w, y|x)} \quad \text{if } \sum_{y \in Y} \mathbb{K}(w, y|x) > 0 \quad (106)$$

$$= \mathbb{K}(w, y|x) \quad \text{if } \sum_{y \in Y} \mathbb{K}(w, y|x) > 0 \quad (107)$$

$$= 0 \quad \text{otherwise} \quad (108)$$

$$= \mathbb{K}(w, y|x) \quad \text{otherwise} \quad (109)$$

In general there are many indexed Markov kernels that satisfy this.  $\square$

**Lemma 2.16** (Valid disintegration choice). *Given valid  $\mathbb{K}^{WY|X}$ , there exists a valid conditional probability  $\bar{\mathbb{L}}^{Y|WX}$  such that*

$$\mathbb{K}^{WY|X} = \mathbb{K}^{W|X} \odot \bar{\mathbb{L}}^{Y|WX} \quad (110)$$

*Proof.* From Lemma 2.15, we have the existence of some  $\bar{\mathbb{L}}^{Y|WX} : W \times X \rightarrow Y$ . We need to check that  $\bar{\mathbb{L}}^{Y|WX}$  can be chosen so that it is valid. By validity of  $\mathbb{K}^{W, Y|X}$ ,  $w \in W(\Omega)$  and  $(X, W, Y) \bowtie (x, w, y) = \emptyset \implies \mathbb{K}^{W, Y|X} = 0$ , so we only need to check for  $(w, x, y)$  such that  $\mathbb{K}^{W, Y|X}(w, y|x) = 0$ . For all  $x, y$  such that  $\mathbb{K}^{Y|X}(y|x)$  is positive, we have  $\mathbb{K}^{W, Y|X}(w, y|x) = 0 \implies \bar{\mathbb{L}}^{Y|WX}(y|w, x) = 0$ . Furthermore, where  $\mathbb{K}^{W|X}(w|x) = 0$ , we either have  $(W, X) \bowtie (w, x) = \emptyset$  or we can choose some  $\omega \in (W, X) \bowtie (w, x)$  and let  $\bar{\mathbb{L}}^{Y|WX}(Y(\omega)|w, x) = 1$ .  $\square$

### 3 Decision theoretic causal inference

The first question we want to investigate is: supposing that we are happy to use the modelling approach described in the previous section, what kind of model would we want to use to help make good choices when we have to make choices?

Suppose we will be given an observation, modelled by  $X$  taking values in  $X$ , and in response to this we can select any decision, modelled by  $D$  taking values in  $D$ . The process by which we choose a decision or mixture of decisions, is called a decision rule or a *strategy*, designated  $\alpha$  and modelled by  $\mathbb{S}_\alpha : X \rightarrow \Delta(D)^2$ . We assume that the collection of strategies under consideration  $\{\mathbb{S}_\alpha\}_\alpha$  is convex. We are interested in some defined collection of things that will be determined at some point after we have taken our decision; these will be modelled by the variable  $Y$  and we will call them *consequences*.

For different observations and decisions we will generally expect different consequences. We will assume that we expect the same observations whatever strategy we choose. We will also assume that given the same observations and the same decision, we expect the same consequences regardless of the strategy. These assumptions rule out certain classes of decision problem where, for example, there is controversy over whether the strategy chosen should depend on the time at which it is chosen Weirich (2016); Lewis (1981); Paul F. Christiano (2018).

We will entertain a collection of probabilistic models to represent postulated relationships between  $X$ ,  $D$  and  $Y$  for each strategy  $\alpha$ ; to do this, we will introduce a latent variable  $H$  such that each value of  $H$  corresponds to a particular probabilistic model of  $X$ ,  $D$  and  $Y$ . Concretely, for each strategy  $\alpha$  our forecast will be represented by a probability model  $\mathbb{P}_\alpha : I \rightarrow (H, X, D, Y)$ . We assume that – holding the hypothesis fixed – the same observations are expected whatever strategy we choose:  $\mathbb{P}_\alpha^{X|H} = \mathbb{P}_\beta^{X|H}$  for all  $\alpha, \beta$ . We assume that under each hypothesis, the decision chosen is always modelled by the chosen strategy:  $\mathbb{P}_\alpha^{D|HX} = \mathbb{S}_\alpha \otimes \text{erase}_H$ . Finally, we assume that, holding the hypothesis fixed, the same consequences are expected under any strategy given the same observations and the same decision:  $\mathbb{P}_\alpha^{Y|XHD} = \mathbb{P}_\beta^{Y|XHD}$  for all  $\alpha, \beta$ .

Under these assumptions, there exists a “see-do model”  $\mathbb{T}^{XY|HD}$  such that  $X \perp\!\!\!\perp_{\mathbb{T}} D|H$  and for all  $\alpha$ ,

$$\mathbb{P}_\alpha = \begin{array}{c} \begin{array}{ccccc} & & \boxed{\mathbb{S}_\alpha^{D|X}} & & \\ D & \text{---} & & & D \\ & & & & \uparrow \\ H & \text{---} & \boxed{\mathbb{T}^{X|H}} & \text{---} & \boxed{\mathbb{T}^{Y|DHX}} & \text{---} & Y \\ & & & & \uparrow \\ & & & & X \end{array} \end{array} \quad (111)$$

The proof is given in Appendix 6. Note that  $\mathbb{T}^{X|H}$  exists by virtue of the fact  $X \perp\!\!\!\perp_{\mathbb{T}} D|H$ .

<sup>2</sup>We don’t make the strategy a variable simply because we would need an uncountable version of our theory to do it.

We will call the the see-do model along with the collection of strategies  $\{\mathbb{T}^{\text{XY}|\text{HD}}, \{\mathbb{S}_\alpha | \alpha \in \mathcal{A}\}\}$  a *standard decision problem*.

### 3.1 Combs

The conditional independence  $X \perp\!\!\!\perp_{\mathbb{T}} D | H$  of  $\mathbb{T}$  is the property that allows us to write Equation 111, but it also implies that  $\mathbb{T}$  is *not* a submodel of  $\mathbb{P}_\alpha$  for most strategies  $\alpha$ , because for most such strategies  $X$  and  $D$  are not independent. Instead,  $\mathbb{T}$  is a *comb*. This structure was introduced by Chiribella et al. (2008) in the context of quantum circuit architecture, and Jacobs et al. (2019) adapted the concept to causal modelling.

We don't formally define any special operations with combs here, but because they come up multiple times we will explain the notion a little. A comb is a Markov kernel with an “insert” operation; to obtain the probability model associated with a particular strategy, we “insert” the strategy into our see-do model.

$$\mathbb{T} = \begin{array}{c} \text{H} \rightarrow \boxed{\mathbb{T}_{\text{X|H}}} \xrightarrow{\text{X}} \boxed{\mathbb{T}_{\text{Y|XD}}} \xrightarrow{\text{Y}} \\ \quad \quad \quad \downarrow \quad \quad \quad \downarrow \\ \quad \quad \quad \text{D} \end{array} \quad (112)$$

$$= \text{H} - \boxed{\text{T} - \text{X} \text{ D}} - \text{Y} \quad (113)$$

A key feature of a comb is that a strategy can be chosen such that  $D$  is independent of any variable on the “upper arm” ( $X$  in this example) conditional on  $H$ . There is an intuitive appeal to the notion that, with access to a randomiser, we could if we wanted to choose a decision independent of all of our observations. We may wish to introduce additional variables that we do not observe, but we can nonetheless choose  $D$  independent of them. Such variables we will call *pre-choice variables*

**Definition 3.1** (Pre-choice variable). Given a see-do model  $\mathbb{T}$ ,  $W$  is a pre-choice variable iff for every other pre-choice variable  $V$ ,  $(W, V) \perp\!\!\!\perp_{\mathbb{T}} D|H$ . The hypothesis  $H$  is always a pre-choice variable, and we also assume the same is true of the observation  $X$ .

Given that  $\mathbf{H}$  is necessarily a pre-choice variable, we wonder if it may be possible to define a hypothesis  $\mathbf{H}$  such that all pre-choice variables are functions of it. This would reduce the number of different elements of our theory, as we would no longer distinguish between “hypotheses” and “pre-choice variables”. The reason why we have not done so thus far is that hypotheses are motivated by classical statistics while pre-choice variables are motivated by approaches to causal inference, and we haven’t yet investigated whether the two can be identified without losing anything important.

Lattimore and Rohde (2019a) and Lattimore and Rohde (2019b) describe a novel approach to causal inference: they consider an observational probability model and a collection of indexed interventional probability models, with the



probability model tied to the interventional models by shared parameters. In these papers, they show how such a model can reproduce inferences made using Causal Bayesian Networks. This kind of model is very close to a type of see-do model, where we identify the hypotheses  $H$  with the parameter variables in that work. The only difference is that we consider interventional maps (see-do models represent a map  $(D, H) \rightarrow Y$ ) rather than interventional probability models, and this is a superficial difference as an indexed collection of probability models is a map.

Dawid (2020) describes a different version of a decision theoretic approach to causal inference:

A fundamental feature of the DT approach is its consideration of the relationships between the various probability distributions that govern different regimes of interest. As a very simple example, suppose that we have a binary treatment variable  $T$ , and a response variable  $Y$ . We consider three different regimes [...] the first two regimes may be described as interventional, and the last as observational.

This is somewhat different to a see-do model, as it features a probabilistic model that uses the same random variables  $T$  and  $Y$  to represent both interventional and observational regimes, while a see-do model uses different random variables. This difference can be thought of as the difference between positing a sequence  $(X_1, X_2, X_3)$  distributed according to  $\mathbb{P}^X$ , or saying that the  $X_i$  are distributed according to  $\mathbb{P}$  such that they are mutually independent ( $i \notin A \subset [3] \implies X_i \perp\!\!\!\perp_{\mathbb{P}} (X_j)_{j \in A}$ ) and identically distributed ( $\mathbb{P}^{X_i} = \mathbb{P}^{X_j}$  for all  $i, j$ ). The former can be understood as a shorthand of the latter, but because in this paper we are particularly interested in problems that arise regarding the relation between the map and the territory, we favour the second approach because it is more explicit.

Jacobs et al. (2019) has used a comb decomposition theorem to prove a sufficient identification condition similar to the identification condition given by Tian and Pearl (2002). This theorem depends on the particular inductive hypotheses made by causal Bayesian networks.

### 3.2 See-do models and classical statistics

See-do models are capable of expressing the expected results of a particular choice of decision strategy, but they cannot by themselves tell us which strategies are more desirable than others. To do this, we need some measure of the desirability of our collection of results  $\{\mathbb{P}_\alpha | \alpha \in A\}$ . A common way to do this is to employ the principle of expected utility. The classic result of Von Neumann and Morgenstern (1944) shows that all preferences over a collection of probability models that obey their axioms of completeness, transitivity, continuity and independence of irrelevant alternatives must be able to be expressed via the principle of expected utility. This does not imply that anyone knows what the appropriate utility function is.

We introduced the hypothesis  $H$  as a latent variable to allow us to postulate multiple different models of observations, decisions and consequences. In general, both the hypothesis and the observation  $X$  may influence our views about the consequences  $Y$  that are likely to follow from a given decision. It is very common to model sequences of observations as independent and identically distributed given some parameter or latent variable. In such cases, we can identify  $H$  with this latent variable (our setup does not preclude introducing a prior over  $H$ , nor does it require it). Furthermore, in such cases where we have a collection of  $X_i$  such that  $X_i \perp\!\!\!\perp_{\mathbb{T}} X_j | H$ , it may be reasonable to expect that  $Y \perp\!\!\!\perp_{\mathbb{T}} X | H$  also. In fact, this is the standard view in causal modelling – given “the probability distribution over observations” (which is to say, conditional on  $H$ ), interventional distributions have no additional dependence on *particular* observations. We can find exceptions with questions like “given what actually happened, what would have happened if a different action had been taken?” (Pearl, 2009; Tian and Pearl, 2000; Mueller et al., 2021), but this is not the kind of question we are considering here.

Given these two choices – to use the principle of expected utility to evaluate strategies, and to use a see-do model  $\mathbb{T}$  with the conditional independence  $Y \perp\!\!\!\perp_{\mathbb{T}} X | H, D$  – we obtain a statistical decision problem in the form introduced by Wald (1950).

A *statistical model* (or *statistical experiment*) is a collection of probability distributions  $\{\mathbb{P}_\theta\}$  indexed by some set  $\Theta$ . A statistical decision problem gives us an observation variable  $X : \Omega \rightarrow X$  and a statistical experiment  $\{\mathbb{P}_\theta^X\}_\Theta$ , a decision set  $D$  and a loss  $l : \Theta \times D \rightarrow \mathbb{R}$ . A strategy  $\mathbb{S}_\alpha^{D|X}$  is evaluated according to the risk functional  $R(\theta, \alpha) := \sum_{x \in X} \sum_{d \in D} \mathbb{P}_\theta^X(x) S_\alpha^{D|X}(d|x) l(h, d)$ . A strategy  $\mathbb{S}_\alpha^{D|X}$  is considered more desirable than  $\mathbb{S}_\beta^{D|X}$  if  $R(\theta, \alpha) < R(\theta, \beta)$ .

Suppose we have a see-do model  $\mathbb{T}^{X|H,D}$  with  $Y \perp\!\!\!\perp_{\mathbb{T}} X | (H, D)$ , and suppose that the random variable  $Y$  is a “reverse utility” function taking values in  $\mathbb{R}$  for which low values are considered desirable. Then, defining a loss  $l : H \times D \rightarrow \mathbb{R}$  by  $l(h, d) = \sum_{y \in \mathbb{R}} y \mathbb{T}^{Y|HD}(y|h, d)$ , we have

$$\mathbb{E}_{\mathbb{P}_\alpha}[Y|h] = \sum_{x \in X} \sum_{d \in D} \sum_{y \in Y} \mathbb{T}^{X|H}(x|h) \mathbb{S}_\alpha^{D|X}(d|x) \mathbb{T}^{Y|HD}(y|h, d) \quad (114)$$

$$= \sum_{x \in X} \sum_{d \in D} \mathbb{T}^{X|H}(x|h) \mathbb{S}_\alpha^{D|X}(d|x) l(h, d) \quad (115)$$

$$= R(h, \alpha) \quad (116)$$

If we are given a see-do model where we interpret  $\mathbb{T}^{X|H}$  as a statistical experiment and  $Y$  as a reversed utility, the expectation of the utility under the strategy forecast given in equation 111 is the risk of that strategy under hypothesis  $h$ .

## 4 Causal Bayesian Networks

When do causal relationships as defined by causal Bayesian networks exist? We will consider a simplified case where a single node may be intervened on, and find the implied see-do model. With this condition, according to Pearl (2009), a causal Bayesian network is a probability model  $\mathbb{P}$ , a collection of interventional probability models  $\{\mathbb{P}_{X=a} | a \in X_i\}$  and a directed acyclic graph  $\mathcal{G}$  whose nodes are identified with some collection of variables, which we can group into three variables  $\{W, X, Y\}$ , where  $W$  is the sequence of variables associated with the parents of  $X$  in  $\mathcal{G}$ ,  $X$  is the “intervenable” node of  $\mathcal{G}$  and  $Y$  are associated with the other nodes. The interventional probability models must all obey the truncated factorisation condition with respect to  $\mathcal{G}$ :

$$\mathbb{P}_{X=a}^{WXY}(w, x, y) = \mathbb{P}^W(w) \mathbb{P}^{Y|XW}(y|x, w) \llbracket x = a \rrbracket \quad (117)$$

A standard interpretation of the observational and interventional probability distributions is that we have a sequence of observations modeled by  $V_A := (W_i, X_i, Y_i)_{i \in A}$  mutually independent and identically distributed according to  $\mathbb{P}^{WXY}$ , and a sequence of consequences modeled by  $V_B := (W_i, X_i, Y_i)_{i \in B}$  mutually independent and identically distributed according to  $\mathbb{P}_{X=a}^{WXY}$ , and  $\mathbb{P}$  and  $\mathbb{P}_{X=a}$  are coupled by Equation 117. What it means for  $\mathbb{P}$  and  $\mathbb{P}_{X=a}$  to be coupled is: if  $\mathbb{P}$  is the “actual” distribution of observations, then  $\mathbb{P}_{X=a}$  is the “actual” distribution of consequences. This can be explicitly represented by introducing a variable  $H$  representing the “actual” distribution of observations, and we introduce a model  $\mathbb{U}^{\cdot|H}$  such that

$$\begin{aligned} \mathbb{P}^{V_i} &:= \mathbb{U}^{V_i|H}(v|h) \text{ for some } h \in H \text{ and any } i \in A, v \in W \times X \times Y \quad (118) \\ \mathbb{P}_{X_j=a}^{V_j} &:= \mathbb{U}^{V_j|HX_j}(v|h, a) \text{ for some } h \in H \text{ and any } j \in B, v \in W \times X \times Y \end{aligned} \quad (119)$$

We justify line 119 by noting that  $\mathbb{U}^{V_i|HX}$  is a Markov kernel  $H \times X \rightarrow W \times X \times Y$ , which is the same type as the map  $\mathbb{Q} := h, a \mapsto \mathbb{P}_{X=a}$ , and in addition Equation 117 ensures that defining  $\mathbb{U}^{V_i|HX} := \mathbb{Q}$  is consistent via Lemma 2.14.

Note that the assumptions of mutual independence  $V_i \perp\!\!\!\perp V_{A \cup B \setminus \{i\}} | H$  for  $i \in A$  and  $V_j \perp\!\!\!\perp V_{A \cup B \setminus \{j\}} | HX_j$  for  $j \in B$  are required for the existence of  $\mathbb{U}^{V_i|H}$  and  $\mathbb{U}^{V_j|HX_j}$  respectively.

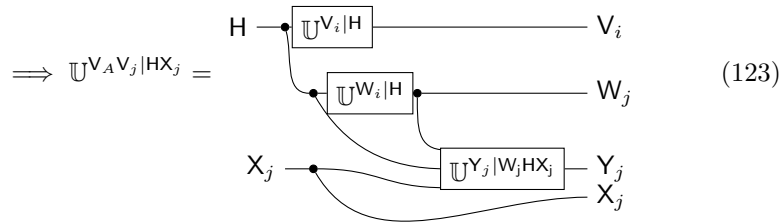
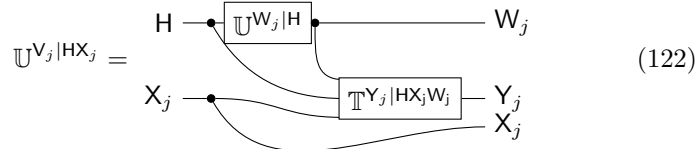
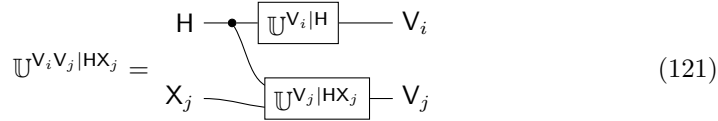
Then Equation 117 becomes

$$\mathbb{U}^{W_j X_j Y_j | HX_j}(w, x, y|h, a) = \mathbb{U}^{W_i|H}(w) \mathbb{U}^{Y_i | X_i W_i H}(y|x, w, h) \llbracket x = a \rrbracket \quad i \in A, j \in B \quad (120)$$

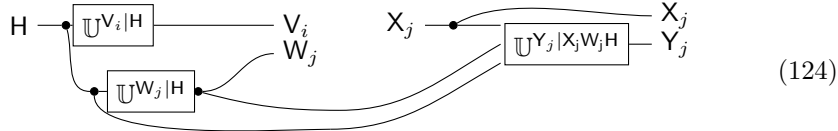
The only difference here is that the coupling between distributions of observations and consequences via  $H$  is explicit.

In most situations,  $A$  will be disjoint from  $B$ . While we don’t necessarily want to rule out considering consequences to be equal to observations, we usually want to consider consequences that may take different values from observations.

For  $i \in A, j \in B$ , we can write  $\mathbb{U}^{\mathbf{V}_i \mathbf{V}_j | \mathbf{H}\mathbf{X}_j}$  as follows



Note that we replace the single observation  $V_i$  with the full observations  $V_A$  as we will make use of them subsequently, and we can do this without issue due to the assumption of conditional independence among the  $V_k$ s. It will be sufficient to consider a single consequence  $V_j$ . Equation 123 defines a model  $\mathbb{U}^{\text{HX}_j}$  which relates observations to consequences in the manner suggested by Equation 117. We will call  $\mathbb{U}$  a “CBN model”. We note that the model in Equation 123 looks like a 2-comb:



However, we have not at this point assumed that we have a convex set of strategies. Suppose we have some standard see-do model  $\mathcal{M} := \{\mathbb{T}^{\text{OV}_B|\text{HD}}, \{\mathbb{S}_\alpha^{\text{D|O}}|\alpha \in \mathcal{A}\}\}$ . The question we want to ask is: when can we posit a see-do model  $\{\mathbb{U}^{\text{V}_A\text{V}_j|\text{HX}_j}, \{\mathbb{R}_\alpha^{\text{X}_j|\text{V}_A\text{W}_j\text{H}}|\alpha \in \mathcal{A}\}\}$  consistent with  $\mathcal{M}$  in the sense that, for all

$\alpha \in \mathcal{A}$ :

$$\mathbb{P}_\alpha^{V_B|H} := H \text{ --- } \boxed{\mathbb{T}^{V_A|H}} \text{ --- } \boxed{\mathbb{S}_\alpha^{D_j|V_A}} \text{ --- } \boxed{\mathbb{T}^{V_j|DV_A H}} \begin{matrix} X_j \\ Y_j \\ W_j \end{matrix} \quad (125)$$

$$= H \text{ --- } \boxed{\mathbb{U}^{V_A W_j|H}} \text{ --- } \boxed{\mathbb{R}_\alpha^{X_j|H V_A W_j}} \text{ --- } \boxed{\mathbb{U}^{Y_j|X_j W_j H}} \begin{matrix} X_j \\ Y_j \\ W_j \end{matrix} \quad (126)$$

$$=: \mathbb{Q}_\alpha^{V_B|H} \quad (127)$$

I think reusing the same  $H$  between  $\mathbb{U}$  and  $\mathbb{T}$  is a mistake here. Maybe not a big problem, but ideally one would check!

**Theorem 4.1.** *Given a standard see-do model  $\mathcal{M} := \{\mathbb{T}^{OV_B|HD}, \{\mathbb{S}_\alpha^{D|V_A} | \alpha \in \mathcal{A}\}\}$  and a CBN model  $\mathbb{U}^{V_A V_j|HX_j}$  as defined in Equation 123, assuming  $W_j$  is a pre-choice variable, then there exists a see-do model  $\{\mathbb{U}^{V_i V_j|HX_j}, \{\mathbb{R}_\alpha^{X_j|V_A W_j H} | \alpha \in \mathcal{A}\}\}$  consistent with  $\mathcal{M}$  if and only if*

1.  $W_j$  is a pre-choice variable, i.e.  $(V_A, W_j) \perp\!\!\!\perp_{\mathbb{T}} D|H$
2.  $\mathbb{T}^{V_A W_j|H} = \mathbb{U}^{V_A W_j|H}$
3.  $Y_j \perp\!\!\!\perp_{\mathbb{T}} D|W_j V_A H X_j$
4.  $\mathbb{T}^{Y_j|W_j V_A H X_j} = \mathbb{U}^{Y_j|W_j V_A H X_j}$

*Proof. If:* If all assumptions hold, we can write

$$\mathbb{T}^{V_A V_j|HD} = H \text{ --- } \boxed{\mathbb{U}^{W_j V_A|H}} \text{ --- } \boxed{\mathbb{T}^{X|W_j V_A HD}} \text{ --- } \boxed{\mathbb{U}^{Y_j|W_j H X_j}} \begin{matrix} V_A \\ W_j \\ Y_j \\ X_j \end{matrix} \quad (128)$$

For each  $\mathbb{S}_\alpha^{D|V_A}$ , define

$$\mathbb{R}_\alpha^{X_j|V_A W_j H} := W_j \text{ --- } \boxed{\mathbb{T}^{X|W_j V_A HD}} \text{ --- } X_j \quad (129)$$

$H$  ---  $\boxed{\mathbb{S}_\alpha^{D|V_A}}$  ---  $V_A$

Then

$$(130)$$

$$(131)$$

$$(132)$$

**Only if:** Suppose assumption 1 does not hold. Then there exists some  $d, d' \in D$ ,  $w \in W$ ,  $h \in H$  such that  $\mathbb{T}^{W_j|HD}(w_j|h, d) \neq \mathbb{T}^{W_j|HD}(w_j|h, d')$ . Then choose  $\mathbb{S}_d^{D|V_A} : v_A \mapsto \delta_d$  and  $\mathbb{S}_{d'}^{D|V_A} : v \mapsto \delta_{d'}$  for all  $v \in V^{[A]}$ . Then define

$$\mathbb{P}_d^{W_j|H}(w|h) = \mathbb{T}^{W_j|HD}(w_j|h, d) \quad (133)$$

$$\neq \mathbb{T}^{W_j|HD}(w_j|h, d') \quad (134)$$

$$= \mathbb{P}_{d'}^{W_j|H}(w|h) \quad (135)$$

But for any  $\alpha, \alpha'$ ,  $\mathbb{Q}_{\alpha}^{W_j|H} = \mathbb{Q}_{\alpha'}^{W_j|H}$  as  $W_j \perp_{\mathcal{U}} X_j|H$ , so  $\mathbb{Q} \neq \mathbb{P}$ . Suppose assumption 1 holds but assumption 2 does not. Then for any  $\alpha$

$$\mathbb{P}_{\alpha}^{V_A W_j|H} = \mathbb{T}^{V_A W_j|H} \quad (136)$$

$$\neq \mathbb{U}^{V_A W_j|H} \quad (137)$$

$$= \mathbb{Q}_{\alpha}^{V_A W_j|H} \quad (138)$$

Suppose assumption 3 does not hold. Then there is some  $d, d' \in D$ ,  $w \in W$ ,  $h \in H$ ,  $v \in V^{[A]}$ ,  $x \in X$  and  $y \in Y$  such that

$$\mathbb{T}^{Y_j|W_j V_A H X_j D}(y|w, v, h, x, d) \neq \mathbb{T}^{Y_j|W_j V_A H X_j D}(y|w, v, h, x, d') \quad (139)$$

$$\text{and } \mathbb{T}^{X_j W_j V_A|HD}(x, w, v|h, d) > 0 \quad (140)$$

$$\text{and } \mathbb{T}^{X_j W_j V_A|HD}(x, w, v|h, d') > 0 \quad (141)$$

$$(142)$$

The latter conditions hold as if Equation 139 only held on sets of measure 0 then we could choose versions of the conditional probabilities such that the independence held.

Then

$$\mathbb{P}_d^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) = \mathbb{T}^{Y_j|W_jV_AHX_jD}(y|w, v, h, x, d) \quad (143)$$

$$\neq \mathbb{T}^{Y_j|W_jV_AHX_jD}(y|w, v, h, x, d') \quad (144)$$

$$= \mathbb{P}_{d'}^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \quad (145)$$

$$\implies \mathbb{P}_d^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \neq \mathbb{Q}_d^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \quad (146)$$

$$\text{or } \mathbb{P}_{d'}^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \neq \mathbb{Q}_{d'}^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \quad (147)$$

As the conditional probabilities disagree on a positive measure set,  $\mathbb{P} \neq \mathbb{Q}$ .

Suppose assumption 3 holds but assumption 4 does not. Then for some  $h \in H$ , some  $w \in W$ ,  $v \in V^{|A|}$ ,  $x \in X$  with positive measure and some  $y \in Y$

$$\mathbb{P}_d^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) = \mathbb{T}^{Y_j|W_jV_AHX_j}(y|w, v, h, x) \quad (148)$$

$$\neq \mathbb{U}^{Y_j|W_jV_AHX_j}(y|w, v, h, x) \quad (149)$$

$$\neq \text{model}Q_d^{Y_j|W_jV_AHX_jD}(y|w, v, h, x) \quad (150)$$

□

Conditional independences like  $(V_A, W_j) \perp\!\!\!\perp_{\mathbb{T}} D|H$  and  $Y_j \perp\!\!\!\perp_{\mathbb{T}} D|W_jV_AHX_j$  bear some resemblance to the condition of “limited unresponsiveness” proposed by Heckerman and Shachter (1995). They are conceptually similar in that they indicate that a particular variable does not “depend on” a decision  $D$  in some sense. As Heckerman points out, however, limited unresponsiveness is not equivalent to conditional independence. We tentatively speculate that there may be a relation between our “pre-choice variables”  $(W_j, V_A, H)$  and the “state” in Heckerman’s work crucial for defining limited unresponsiveness.

## 4.1 Proxy control

We say that  $(V_A, W_j) \perp\!\!\!\perp_{\mathbb{T}} D|H$  expresses the notion that  $W_j$  is a *pre-choice variable* and  $(W_j, V_A, X_j)$  are *proxies for*  $D$  with respect to  $Y$  under conditions of full information. To justify this terminology, we note that under a strong assumption of identifiability  $Y_j \perp\!\!\!\perp H|W_jV_AX_j$  (i.e. the observed data allow us

$$\begin{aligned}
\mathbb{T}^{V_A V_B | HD} &= \text{Diagram (151)} \\
&= \text{Diagram (152)}
\end{aligned}$$

Smoking cannot be stopped by any legal or educational means available to us today; cigarette advertising can. That does not stop researchers from aiming to estimate “the effect of smoking on cancer,” and doing so from experiments in which they vary the instrument—cigarette advertisement—not smoking. The reason they would be interested in the atomic intervention  $P(\text{cancer}|\text{do}(\text{smoking}))$  rather than (or in addition to)  $P(\text{cancer}|\text{do}(\text{advertising}))$  is that the former represents a stable biological characteristic of the population, uncontaminated by social factors that affect susceptibility to advertisement, thus rendering it transportable across cultures and environments. With the help of this stable characteristic, one can assess the effects of a wide variety of practical policies, each employing a different smoking-reduction instrument.

Like causal Bayesian networks, causal models in the potential outcomes framework typically do not include any variables representing what we call “consequences”. A potential outcomes model features a sequence of observable variables  $(Y_i, X_i, Z_i)_{i \in [n]}$  and a collection of potential outcomes  $(Y_i^x)_{x \in X, i \in [n]}$ . Also like causal Bayesian networks, we think that introducing the idea of consequences clarifies the meaning of potential outcomes models.

32



**Definition 5.1** (Selector). Given variables  $X : \Omega \rightarrow X$  and  $\{Y^x : \Omega \rightarrow Y \mid x \in X\}$ , define  $Y^X : (\mathcal{Y}^x)_{x \in X}$ . The selector  $\pi : X \times Y^X \rightarrow Y$  is the function that sends  $(x, y^1, \dots, y^{|X|}) \rightarrow y^x$ .

**Definition 5.2** (Potential outcomes: formal requirement). Given variables  $Y : \Omega \rightarrow Y$  and  $X : \Omega \rightarrow X$ , we introduce a collection of latent variables called *potential outcomes*  $Y^X := (\mathcal{Y}^x)_{x \in X}$  such that  $Y = \pi \circ (X, Y^X)$ . A *potential outcomes model* is any consistent model of  $Y$ ,  $X$  and  $Y^X$ .

Lemma 5.3 shows we can always define trivial potential outcomes of  $Y$  with respect to  $X$  by taking the product of  $|X|$  copies of  $Y$ . We need some other constraint on the values of potential outcomes besides the formal definition 5.2 if we want them to be informative.

**Lemma 5.3** (Trivial formal potential outcomes). *For any variables  $Y : \Omega \rightarrow Y$ ,  $X : \Omega \rightarrow X$  and  $W : \Omega \rightarrow W$ , we can always define potential outcomes  $Y_X$  such that any consistent model  $\mathbb{K}^{YX|W}$  can be extended to a consistent model of  $\mathbb{K}^{YXY^X|W}$ .*

*Proof.* Define  $Y^X := (Y)_{x \in X}$ . Then we can consistently extend  $\mathbb{K}^{YX|W}$  to  $\mathbb{K}^{YXY^X|W}$  by repeated application of Lemma 2.13.  $\square$

The trivial potential outcomes of Lemma 5.3 are in many cases unsatisfactory for what we want potential outcomes to represent. Thus Definition 5.2 is incomplete. In common with observable variables, the definition of potential outcomes involves both the formal requirement of Definition 5.2, and an indication of the parts of the real world that they model. Unlike observable variables, the “part of the world” that potential outcomes model will not at any point resolve to a canonical value. We say the potential outcome  $Y^x := \pi(x, Y)$  is “the value that  $Y$  would take if  $X$  were  $x$ , whether or not  $X$  actually takes the value  $x$ ”. We will call this additional element of the definition of potential outcomes the *counterfactual extension*.

**Definition 5.4** (What potential outcomes model: counterfactual extension). Given observables  $X$ ,  $Y$  and  $Y^X$ ,  $Y^X$  are potential outcomes if they satisfy Definition 5.2 and for all  $x \in X$ , the individual potential outcome  $Y^x := \pi(x, Y)$  models the value  $Y$  would take if  $X$  took the value  $x$ .

Because observables resolve to a single canonical value, the conditional in Definition 5.4 is eventually satisfied for exactly one  $x \in X$ , at which point  $Y^{x'}$  for all  $x' \neq x$  are guaranteed not to resolve. Nevertheless, we can maybe draw some conclusions about  $Y^X$  from Definition 5.4. For example, it seems unreasonable in light of this definition to assert that  $Y^x$  is *necessarily* identical to  $Y$  for all  $x \in X$ , which rules out the strictly trivial potential outcomes of Lemma 5.3.

We will note at this point that if  $X$  refers to a person’s body mass index and  $Y$  to an indicator of whether or not they experience heart disease, it is metaphysically subtle to say whether  $Y^X$  is well-defined with regard to Definitions 5.2 and 5.4 together. Recall that there are multiple ways that a given level of body

mass index ( $X$ ) could be achieved. One might say that, when there are multiple possible paths, there is no unique way to choose a path. However, a very similar argument can be made that whenever there are multiple possible values of  $Y^x$  (which is whenever  $X$  does not take the value  $x$ ), then there is no unique choice of  $Y^x$ , which implies that the full set of potential outcomes  $Y^X$  is *almost never well-defined*. Alternatively, if there is some method of making a canonical choice of  $Y^x$ , then perhaps this same method can also make a canonical choice of which path was taken to achieve this value of  $X$ .

We will set Definition 5.4 aside and propose an alternative decision-theoretic extension of the definition of potential outcomes. To motivate this proposal, we first note that, if we are using potential outcomes  $Y^X$  to model an observation of  $X$  and  $Y$  only conditional on some hypothesis (or parameter)  $H$ , then by repeated application of Lemma ??, we can represent the model  $\mathbb{P}^{XY^X|H}$  of these variables as

$$\mathbb{P}^{XY^X|H} = H \quad (153)$$

For any collection of representative kernels  $\mathbb{T}^{Y^X|H}$ ,  $\mathbb{T}^{X|Y^X H}$  and  $\mathbb{T}^{Y|HY^X X}$ . We can simplify Equation 153 somewhat. Firstly,  $\mathbb{P}^{Y|HY^X X}$  must always be represented a *selector kernel*  $\Pi : X \times Y^{|X|} \rightarrow Y$ , as shown by Lemma 5.5.

**Lemma 5.5** (Selector kernel). *Let the selector kernel  $\Pi : X \times Y^X \rightarrow Y$  be defined by  $\Pi_{(x,y^X)}^y = \llbracket \pi(x, y^X) = y \rrbracket$ . Given  $X$ ,  $Y$ , potential outcomes  $Y^X$  and arbitrary  $W$ , defining  $\mathbb{Q} : X \times Y^X \times W \rightarrow Y$  by*

$$\mathbb{Q} := \begin{array}{c} Y^X \\ X \\ W \end{array} \begin{array}{c} \diagup \\ \diagdown \\ \longrightarrow \end{array} \Pi \longrightarrow Y \quad (154)$$

$$\iff \quad (155)$$

$$\mathbb{Q}_{(y^X, x, w)}^y = \Pi_{(x, y^X)}^y \quad \forall y, y^X, x, w \quad (156)$$

Then any potential outcomes model  $\mathbb{T}^{YY^X X|W}$  must have the property that, for all  $x, w, y^X$  and  $y$ ,  $\mathbb{Q}$  is a representative of  $\mathbb{T}^{Y|Y^X X W}$ .

*Proof.* Recall  $Y = \pi \circ (X, Y^X)$ . Thus consistency implies that  $Y \stackrel{a.s.}{=} \pi \circ (X, Y^X)$  for all  $(x, y^X, w) \in \text{Range}(X) \times \text{Range}(Y) \times \text{Range}(W)$  such that  $X^{-1}(x) \cap (Y^X)^{-1}(y^X) \cap W^{-1}(w) \neq \emptyset$ . However, wherever  $X^{-1}(x) \cap (Y^X)^{-1}(y^X) \cap W^{-1}(w) = \emptyset$ , consistency implies  $\mathbb{T}^{YY^X X|W}(y, y^X, x|w) = 0$  and so  $\mathbb{T}^{Y|Y^X X W}$  is arbitrary on this collection of values. Equations 154 and 156 are equivalent to the statement  $Y \stackrel{a.s.}{=} \pi \circ (X, Y^X)$ .  $\square$

Thus we can without loss of generality choose  $\Pi$  to represent  $\mathbb{T}^{Y|Y^X \times W}$ . We observe that when Rubin (2005) describes a potential outcomes model, he calls  $\mathbb{T}^{Y^X|H}$  “the science” and  $\mathbb{T}^{X|HY^X}$  the “selection function”. He goes on to explain that the science “is not affected by how or whether we try to learn about it”.

We propose a definition of potential outcomes that enshrines the stability of “the science”.

**Definition 5.6.** Potential outcomes: decision theoretic extension Given a standard decision problem  $\{\mathbb{T}^{WZ|HD}, \{\mathbb{S}_\alpha\}_{\alpha \in A}\}$ ,  $Y^X$  is a potential outcome for  $Y$  with respect to  $X$  if it satisfies Definition 5.2 and is a prechoice variable; that is,  $(Y^X, W) \perp\!\!\!\perp_{\mathbb{T}} D|H$ .

Owing to the subtlety of interpreting Definition 5.4, we don’t know a straightforward argument to the effect that Definition 5.6 is implied by it. Besides the fact that it seems to formalise the idea that the distribution of potential outcomes is unaffected by our actions, we will point out that a key feature of prechoice variables – decisions can be chosen so that they are random with respect to all prechoice variables – is used in practice to justify the assumption of ignorability in randomised experiments.

Definition 5.6 can sometimes (but not always) rule out potential outcomes if there is more than one way to achieve a given value of  $X$ . Recall that Hernán and Taubman (2008) argued potential outcomes are “ill-defined” in the presence of multiple treatments.

**Example 5.7.** Suppose we have a standard decision problem  $\{\mathbb{T}^{WZ|HD}, \{\mathbb{S}_\alpha\}_{\alpha \in A}\}$  where observations are  $W$ , consequences  $Z$ , hypotheses  $H$  and decisions  $D \in \{0, 1, 2, 3\}$ . Suppose we also have some  $X \in \{0, 1\}$ ,  $Y$  such that  $\mathbb{T}^{X|HWD}(x|h, w, d) = \mathbb{I}[x = d \bmod 2]$  for all  $h, w$  and, for some  $y$

$$\mathbb{T}^{Y|HWD}(y|h, w, 0, 0) \neq \mathbb{T}^{Y|HWD}(y|h, w, 0, 2) \quad (157)$$

Then we can consider strategies  $\mathbb{S}_0^{D|W} := w \mapsto \delta_0$  and  $\mathbb{S}_2^{D|W} := w \mapsto \delta_2$ . By assumption,

$$\mathbb{P}_0^{Y|HD}(y|h, 0) = \sum_{x \in \{0, 1\}, w \in W} \mathbb{T}^{W|H}(w|h) \mathbb{S}_0^{D|W}(0|w) \mathbb{T}^{X|HWD}(x|h, w, 0) \mathbb{T}^{Y|HWD}(y|h, w, x, 0) \quad (158)$$

$$= \mathbb{T}^{Y|HWD}(y|h, w, 0, 0) \quad (159)$$

$$\neq \mathbb{P}_2^{Y|HD} \quad (160)$$

Suppose we had some potential outcomes  $Y^X$  for  $Y$  with respect to  $X$ . Then, by

assumption

$$\mathbb{P}_0^{Y|HD}(y|h, 0) = \sum_{y^X \in Y^2, x \in \{0,1\}} \mathbb{T}^{Y^X|H}(y^X|h) \mathbb{T}^{X|HDY^X}(x|h, 0, y^X) \Pi(y|x, y^X) \quad (161)$$

$$= \sum_{y^X} \mathbb{T}^{Y^X|H}(y^X|h) \Pi(y|0, y^X) \quad (162)$$

$$= \sum_{y^X \in Y^2, x \in \{0,1\}} \mathbb{T}^{Y^X|H}(y^X|h) \mathbb{T}^{X|HDY^X}(x|h, 2, y^X) \Pi(y|x, y^X) \quad (163)$$

$$= \mathbb{P}_2^{Y|HD} \quad (164)$$

Here we use the property  $Y^X \perp\!\!\!\perp_{\mathbb{T}} D|H$ , implied by the assumption that  $Y^X$  is a prechoice variable. Equations 160 and 164 are clearly contradictory, thus there can be no potential outcomes  $Y^X$  in this example.

I think I asked the wrong question here – should’ve asked when I can extend a see-do model with additional pre-choice variables. I think it’s possible to always choose some deterministic potential outcomes.

**Theorem 5.8** (Existence of potential outcomes). *Suppose we have a standard decision problem  $\{\mathbb{T}^{WZ|HD}, \{\mathbb{S}_\alpha\}_{\alpha \in A}\}$ , and let  $U$  be the sequence of all prechoice variables. For some  $Y$  and  $X$ , there exist potential outcomes  $Y^X$  in the sense of Definition 5.6 if and only if  $\mathbb{T}^{Y|UX}$  exists and is deterministic.*

*Proof.* If: If  $\mathbb{T}^{Y|UX}$  exists and is deterministic then there exists some  $f : U \times X \rightarrow Y$  such that  $Y \stackrel{a.s.}{=} f \circ (U, X)$ . Let  $Y^X := (f(U, x))_{x \in X}$ . Then  $\pi \circ (X, Y^X) = f(U, X) \stackrel{a.s.}{=} Y$ .

Only if: By definition,  $Y^X = g \circ U$ . From Lemma 5.5,  $\mathbb{T}^{Y|XY^X}$  exists and is deterministic. Thus  $\mathbb{T}^{Y|XW}$  also exists and is also deterministic.  $\square$

**Corollary 5.9.** *Potential outcomes  $Y^X$  in the sense of Definition 5.6 exist only if*

$$Y \perp\!\!\!\perp_{\mathbb{T}} D|WX \quad (165)$$

*Proof.*  $\mathbb{T}^{Y|UX}$  exists only if  $Y \perp\!\!\!\perp_{\mathbb{T}} D|UX$ .  $\square$

Note the similarity between Equation 165 and the condition for proxy control in the previous section. Indeed, the two are identical if we identify  $U$  with  $(W_j, V_A, X_j)$ .

## 6 Appendix:see-do model representation

Update notation

**Theorem 6.1** (See-do model representation). *Suppose we have a decision problem that provides us with an observation  $x \in X$ , and in response to this we can select any decision or stochastic mixture of decisions from a set  $D$ ; that is we can choose a “strategy” as any Markov kernel  $\mathbb{S} : X \rightarrow \Delta(D)$ . We have a utility function  $u : Y \rightarrow \mathbb{R}$  that models preferences over the potential consequences of our choice. Furthermore, suppose that we maintain a denumerable set of hypotheses  $H$ , and under each hypothesis  $h \in H$  we model the result of choosing some strategy  $\mathbb{S}$  as a joint probability over observations, decisions and consequences  $\mathbb{P}_{h,\mathbb{S}} \in \Delta(X \times D \times Y)$ .*

*Define  $\mathbf{X}, \mathbf{Y}$  and  $\mathbf{D}$  such that  $\mathbf{X}_{x\mathbf{d}y} = x$ ,  $\mathbf{Y}_{x\mathbf{d}y} = y$  and  $\mathbf{D}_{x\mathbf{d}y} = d$ . Then making the following additional assumptions:*

1. *Holding the hypothesis  $h$  fixed the observations as have the same distribution under any strategy:  $\mathbb{P}_{h,\mathbb{S}}[\mathbf{X}] = \mathbb{P}_{h,\mathbb{S}'}[\mathbf{X}]$  for all  $h, \mathbb{S}, \mathbb{S}'$  (observations are given “before” our strategy has any effect)*
2. *The chosen strategy is a version of the conditional probability of decisions given observations:  $\mathbb{S} = \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]$*
3. *There exists some strategy  $\mathbb{S}$  that is strictly positive*
4. *For any  $h \in H$  and any two strategies  $\mathbb{Q}$  and  $\mathbb{S}$ , we can find versions of each disintegration such that  $\mathbb{P}_{h,\mathbb{Q}}[\mathbf{Y}|\mathbf{D}\mathbf{X}] = \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{D}\mathbf{X}]$  (our strategy tells us nothing about the consequences that we don’t already know from the observations and decisions)*

*Then there exists a unique see-do model  $(\mathbb{T}, \mathbf{H}', \mathbf{D}', \mathbf{X}', \mathbf{Y}')$  such that  $\mathbb{P}_{h,\mathbb{S}}[\mathbf{XDY}]^{ijk} = \mathbb{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbb{S}_i^j \mathbb{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ijk}^k$ .*

*Proof.* Consider some probability  $\mathbb{P} \in \Delta(X \times D \times Y)$ . By the definition of disintegration (section ??), we can write

$$\mathbb{P}[\mathbf{XDY}]^{ijk} = \mathbb{P}[\mathbf{X}]^i \mathbb{P}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{P}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (166)$$

Fix some  $h \in H$  and some strictly positive strategy  $\mathbb{S}$  and define  $\mathbb{T} : H \times D \rightarrow \Delta(X \times Y)$  by

$$\mathbb{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^k \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{kj}^l \quad (167)$$

Note that because  $\mathbb{S}$  is strictly positive and by assumption  $\mathbb{S} = \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]$ ,  $\mathbb{P}_{h,\mathbb{S}}[\mathbf{D}]$  is also strictly positive. Therefore  $\mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{D}]$  is unique and therefore  $\mathbb{T}$  is also unique.

Define  $\mathbf{X}'$  and  $\mathbf{Y}'$  by  $\mathbf{X}'_{xy} = x$  and  $\mathbf{Y}'_{xy} = y$ . Define  $\mathbf{H}'$  and  $\mathbf{D}'$  by  $\mathbf{H}'_{hd} = h$  and  $\mathbf{D}'_{hd} = d$ .

We then have

$$\mathbb{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj}^k = \mathbb{T}\mathbf{X}'_{hj}^k \quad (168)$$

$$= \sum_l \mathbb{T}_{hj}^{kl} \quad (169)$$

$$= \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^k \quad (170)$$

$$= \mathbb{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj'}^k \quad (171)$$

Thus  $\mathbf{X}' \perp\!\!\!\perp_{\mathbb{T}} \mathbf{D}'|\mathbf{H}'$  and so  $\mathbb{T}[\mathbf{X}'|\mathbf{H}']$  exists (section ??) and  $(\mathbb{T}, \mathbf{H}', \mathbf{D}', \mathbf{X}', \mathbf{Y}')$  is a see-do model.

Applying Equation 166 to  $\mathbb{P}_{h,\mathbb{S}}$ :

$$\mathbb{P}_{h,\mathbb{S}}[\mathbf{XDY}]^{ijk} = \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (172)$$

$$= \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (173)$$

$$= \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{T}[\mathbf{X}'\mathbf{Y}'|\mathbf{H}'\mathbf{D}']_{hj}^{ik} \quad (174)$$

$$= \mathbb{S}_i^j \mathbb{T}[\mathbf{X}'\mathbf{Y}'|\mathbf{H}'\mathbf{D}']_{hj}^{ik} \quad (175)$$

$$= \mathbb{S}_i^j \mathbb{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj}^i \mathbb{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (176)$$

$$= \mathbb{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbb{S}_i^j \mathbb{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (177)$$

Consider some arbitrary alternative strategy  $\mathbb{Q}$ . By assumption

$$\mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^i = \mathbb{P}_{h,\mathbb{Q}}[\mathbf{X}]^i \quad (178)$$

$$\mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k = \mathbb{P}_{h,\mathbb{Q}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \text{ for some version of } \mathbb{P}_{h,\mathbb{Q}}[\mathbf{Y}|\mathbf{XD}] \quad (179)$$

It follows that, for some version of  $\mathbb{P}_{h,\mathbb{Q}}[\mathbf{Y}|\mathbf{XD}]$ ,

$$\mathbb{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbb{Q}}[\mathbf{X}]^k \mathbb{P}_{h,\mathbb{Q}}[\mathbf{Y}|\mathbf{XD}]_{kj}^l \quad (180)$$

Then by substitution of  $\mathbb{Q}$  for  $\mathbb{S}$  in Equation 172 and working through the same steps

$$\mathbb{P}_{h,\mathbb{S}}[\mathbf{XDY}]^{ijk} = \mathbb{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbb{Q}_i^j \mathbb{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (181)$$

As  $\mathbb{Q}$  was arbitrary, this holds for all strategies.  $\square$

## 7 Appendix: Counterfactual representation

**Definition 7.1** (Parallel potential outcomes). Given a Markov kernel space  $(\mathbb{K}, E, F)$ , a collection of variables  $\{\mathbf{Y}_i, \mathbf{Y}(W), \mathbf{W}_i\}$ ,  $i \in [n]$ , where  $\mathbf{Y}_i$  and  $\mathbf{Y}(W)$  are random variables and  $\mathbf{W}_i$  could be either a state or random variables is a *parallel potential outcome submodel* if  $\mathbb{K}[\mathbf{Y}_i|\mathbf{W}_i\mathbf{Y}(W)]$  exists and  $\mathbb{K}[\mathbf{Y}_i|\mathbf{W}_i\mathbf{Y}(W)]_{kj_1j_2\dots j_{|W|}} = \delta[j_k]$ .

How this will change: a parallel potential outcomes model is a comb  
 $\mathbb{K}[Y(W)|H] \Rightarrow \mathbb{K}[Y_i|W_i Y(W)]$ .

A parallel potential outcomes model features a sequence of  $n$  “parallel” outcome variables  $Y_i$  and  $n$  “regime proposals”  $W_i$ , with the property that if the regime proposal  $W_i = w_i$  then the corresponding outcome  $Y_i \stackrel{a.s.}{=} Y(w_i)$ . We can identify a particular index, say  $n = 1$ , with the actual world and the rest of the indices with supposed worlds. Thus  $Y_1$  represents the value of TYT in the actual world and  $Y_i$   $i \neq 1$  represents TYT under a supposed regime  $W_i$ . Given such an interpretation, the fact that  $Y_i \stackrel{a.s.}{=} Y(w_i)$  can be interpreted as assuming “for all  $w$ , if the supposed regime  $W_i$  is  $w$  then the corresponding outcome will be almost surely equal to  $Y(w)$ , regardless of the value of the actual regime  $W_1$ ”, which is our original counterfactual assumption.

We do not intend to defend this as the only way that counterfactuals can be modeled, or even that it is appropriate to capture the idea of counterfactuals at all. It is simply a way that we can model the counterfactual assumption typically associated with potential outcomes. We will show that parallel potential outcome submodels correspond precisely to *extendably exchangeable* and *deterministically reproducible* submodels of Markov kernel spaces.

## 7.1 Parallel potential outcomes representation theorem

Exchangeable sequences of random variables are sequences whose joint distribution is unchanged by permutation. Independent and identically distributed random variables are one example: if  $X_1$  is the result of the first flip of a coin that we know to be fair and  $X_2$  is the second flip then  $\mathbb{P}[X_1 X_2] = \mathbb{P}[X_2 X_1]$ . There are also many examples of exchangeable sequences that are not mutually independent and identically distributed – for example, if we want to use random variables  $Y_1$  and  $Y_2$  to model our subjective uncertainty regarding two flips of a coin of unknown fairness, we regard our initial uncertainty for each flip to be equal  $\mathbb{P}[Y_1] = \mathbb{P}[Y_2]$  and we our state of knowledge of the second flip after observing only the first will be the same as our state of knowledge of the first flip after observing only the second  $\mathbb{P}[Y_2|Y_1] = \mathbb{P}[Y_1|Y_2]$ , then our model of subjective uncertainty is exchangeable.

De Finetti’s representation theorem establishes the fact that any infinite exchangeable sequence  $Y_1, Y_2, \dots$  can be modeled by the product of a *prior* probability  $\mathbb{P}[J]$  with  $J$  taking values in the set of marginal probabilities  $\Delta(Y)$  and a conditionally independent and identically distributed Markov kernel  $\mathbb{P}[Y_A|J]_j^{y_A} = \prod_{i \in A} \mathbb{P}[Y_i|J]_j^{y_i}$ .

We extend the idea of exchangeable sequences to cover both random variables and state variables, and we show that a similar representation theorem holds for potential outcomes. De Finetti’s original theorem introduced the variable  $J$  that took values in the set of marginal distributions over a single observation; the set of potential outcome variables plays an analogous role taking values in the set of functions from propositions to outcomes.

The representation theorem for potential outcomes is somewhat simpler than

De Finetti's original theorem due to the fact that potential outcomes are usually assumed to be *deterministically reproducible*; in the parallel potential outcomes model, this means that for  $j \neq i$ , if  $W_j$  and  $W_i$  are equal then  $Y_j$  and  $Y_i$  will be almost surely equal. This assumption of determinism means that we can avoid appeal to a law of large numbers in the proof of our theorem.

An interesting question is whether there is a similar representation theorem for potential outcomes without the assumption of deterministic reproducibility. I'm reasonably confident that this is a straightforward corollary of the representation theorem proved in my thesis. However, this requires maths not introduced in this draft of the paper.

Extendably exchangeable sequences can be permuted without changing their conditional probabilities, and can be extended to arbitrarily long sequences while maintaining this property. We consider here sequences that are exchangeable conditional on some variable; this corresponds to regular exchangeability if the conditioning variable is  $*$  where  $*_i = 1$ .

**Definition 7.2** (Exchangeability). Given a Markov kernel space  $(\mathbb{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $Y_i$  random variables is *exchangeable* conditional on  $Z$  if, defining  $Y_{[n]} = (Y_i)_{i \in [n]}$  and  $D_{[n]} = (D_i)_{i \in [n]}$ ,  $\mathbb{K}[Y_{[n]}|D_{[n]}Z]$  exists and for any bijection  $\pi : [n] \rightarrow [n]$   $\mathbb{K}[Y_{\pi([n])}|D_{\pi([n])}Z] = \mathbb{K}[Y_{[n]}|D_{[n]}Z]$ .

**Definition 7.3** (Extension). Given a Markov kernel space  $(\mathbb{K}, E, F)$ ,  $(\mathbb{K}', E', F')$  is an *extension* of  $(\mathbb{K}, E, F)$  if there is some random variable  $X$  and some state variable  $U$  such that  $\mathbb{K}'[X|U]$  exists and  $\mathbb{K}'[X|U] = \mathbb{K}$ .

If  $(\mathbb{K}', E', F')$  is an extension of  $(\mathbb{K}, E, F)$  we can identify any random variable  $Y$  on  $(\mathbb{K}, E, F)$  with  $Y \circ X$  on  $(\mathbb{K}', E', F')$  and any state variable  $D$  with  $D \circ U$  on  $(\mathbb{K}', E', F')$  and under this identification  $\mathbb{K}'[Y \circ X|D \circ U]$  exists iff  $\mathbb{K}[Y|D]$  exists and  $\mathbb{K}'[Y \circ X|D \circ U] = \mathbb{K}[Y|D]$ . To avoid proliferation of notation, if we propose  $(\mathbb{K}, E, F)$  and later an extension  $(\mathbb{K}', E', F')$ , we will redefine  $\mathbb{K} := \mathbb{K}'$  and  $Y := Y \circ X$  and  $D := D \circ U$ .

I think this is a very standard thing to do – propose some  $X$  and  $\mathbb{P}(X)$  then introduce some random variable  $Y$  and  $\mathbb{P}(XY)$  as if the sample space contained both  $X$  and  $Y$  all along.

**Definition 7.4** (Extendably exchangeable). Given a Markov kernel space  $(\mathbb{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  and a state variable  $Z$  with  $Y_i$  random variables is *extendably exchangeable* if there exists an extension of  $\mathbb{K}$  with respect to which  $((D_i, Y_i))_{i \in \mathbb{N}}$  is exchangeable conditional on  $Z$ .

Here that we identify  $Z$  and  $((D_i, Y_i))_{i \in [n]}$  defined on the extension with the original variables defined on  $(\mathbb{K}, E, F)$  while  $((D_i, Y_i))_{i \in \mathbb{N} \setminus [n]}$  may be defined only on the extension.

Deterministically reproducible sequences have the property that repeating the same decision gets the same response with probability 1. This could be a model of an experiment that exhibits no variation in results (e.g. every time I



put green paint on the page, the page appears green), or an assumption about collections of “what-ifs” (e.g. if I went for a walk an hour ago, just as I actually did, then I definitely would have stubbed my toe, just like I actually did). Incidentally, many consider that this assumption is false concerning what-if questions about things that exhibit quantum behaviour.

**Definition 7.5** (Deterministically reproducible). Given a Markov kernel space  $(\mathbb{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $Y_i$  random variables is *deterministically reproducible* conditional on  $Z$  if  $n \geq 2$ ,  $\mathbb{K}[Y_{[n]}|D_{[n]}Z]$  exists and  $\mathbb{K}[Y_{\{i,j\}}|D_{\{i,j\}}Z]_{kk}^{lm} = \llbracket l = m \rrbracket \mathbb{K}[Y_i|D_iZ]_k^l$  for all  $i, j, k, l, m$ .

**Theorem 7.6** (Potential outcomes representation). *Given a Markov kernel space  $(\mathbb{K}, E, F)$  along with a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $n \geq 2$  and a conditioning variable  $Z$ ,  $(\mathbb{K}, E, F)$  can be extended with a set of variables  $Y(D) := (Y(i))_{i \in D}$  such that  $\{Y_i, Y(D), D_i\}$  is a parallel potential outcome submodel if and only if  $((D_i, Y_i))_{i \in [n]}$  is extendably exchangeable and deterministically reproducible conditional on  $Z$ .*

*Proof.* If: Because  $((D_i, Y_i))_{i \in [n]}$  is extendably exchangeable, we can without loss of generality assume  $n \geq |D|$ .

Let  $e = (e_i)_{i \in [|D|]}$ . Introduce the variable  $Y(i)$  for  $i \in D$  such that  $\mathbb{K}[Y(D)|D_{[D]}Z]_{ez} = \mathbb{K}[Y_D|D_DZ]_{ez}$  and introduce  $X_i$ ,  $i \in D$  such that  $\mathbb{K}[X_i|D_iZY(D)]_{e_i z j_1 \dots j_{|D|}}^{x_i} = \delta[j_{e_i}]^{x_i}$ . Clearly  $\{X_{[n]}, D_{[n]}, Y(D)\}$  is a parallel potential outcome submodel. We aim to show that  $\mathbb{K}[Y_{[n]}|D_{[n]}Z] = \mathbb{K}[X_{[n]}|D_{[n]}Z]$ .

Let  $y := (y_i)_{i \in |D|} \in Y^{|D|}$ ,  $d := (d_i)_{i \in [n]} \in D^{[n]}$ ,  $x := (x_i)_{i \in [n]} \in Y^{[n]}$ .

$$\mathbb{K}[X_n|D_nZ]_{dz}^x = \sum_{y \in Y^{|D|}} \mathbb{K}[X_{[n]}|D_nZY(D)]_{dzy}^x \mathbb{K}[Y(D)|D_{[n]}Z]_{dz}^y \quad (182)$$

$$= \sum_{y \in Y^{|D|}} \prod_{i \in [n]} \delta[y_{d_i}]^{x_i} \mathbb{K}[Y(D)|D_nZ]_{dz}^y \quad (183)$$

Wherever  $d_i = d_j := \alpha$ , every term in the above expression will contain the product  $\delta[\alpha]^{x_i} \delta[\alpha]^{x_j}$ . If  $x_i \neq x_j$ , this will always be zero. By deterministic reproducibility,  $d_i = d_j$  and  $x_i \neq x_j$  implies  $\mathbb{K}[Y_{[n]}|D_{[n]}Z]_{dz}^x = 0$  also. We need to check for equality for sequences  $x$  and  $d$  such that wherever  $d_i = d_j$ ,  $x_i = x_j$ . In this case,  $\delta[\alpha]^{x_i} \delta[\alpha]^{x_j} = \delta[\alpha]^{x_i}$ . Let  $Q_d \subset [n] := \{i \mid \nexists i \in [n] : j < i \text{ \& } d_j = d_i\}$ , i.e.  $Q$  is the set of all indices such that  $d_i$  is the first time this value appears in  $d$ . Note that  $Q_d$  is of size at most  $|D|$ . Let  $Q_d^C = [n] \setminus Q_d$ , let  $R_d \subset D : \{d_i \mid i \in Q_d\}$  i.e. all the elements of  $D$  that appear at least once in the sequence  $d$  and let  $R_d^C = D \setminus R_d$ .

Let  $y' = (y_i)_{i \in Q_d^C}$ ,  $x_{Q_d} = (x_i)_{i \in Q_d}$ ,  $Y(R_d) = (Y_d)_{d \in R_d}$  and  $Y(S_d) = (Y_d)_{d \in S_d}$ .

$$\mathbb{K}[X_{[n]}|D_{[n]}Z]_{dz}^x = \sum_{y \in Y^{|\mathcal{D}|}} \prod_{i \in Q_d} \delta[y_{d_i}]^{x_i} \mathbb{K}[Y(D)|D_{[n]}Z]_{dz}^y \quad (184)$$

$$= \sum_{y' \in Y^{|\mathcal{R}_d^C|}} \mathbb{K}[Y(R_d)Y(R_d^C)|D_{Q_d}D_{Q_d^C}Z]_{d_{Q_d}d_{Q_d^C}z}^{x_{Q_d}y'} \quad (185)$$

$$= \sum_{y' \in Y^{|\mathcal{R}_d^C|}} \mathbb{K}[Y_{R_d}Y_{R_d^C}|D_{Q_d}D_{Q_d^C}Z]_{dz}^{x_{Q_d}y'} \quad (186)$$

$$= \sum_{y' \in Y^{|\mathcal{R}_d^C|}} \mathbb{K}[Y_{[n]}|D_{[n]}Z]_{dz}^{x_{Q_d}y'} \quad (\text{using exchangeability}) \quad (187)$$

Note that

Only if: We aim to show that the sequences  $Y_{[n]}$  and  $D_{[n]}$  in a parallel potential outcomes submodel are exchangeable and deterministically reproducible.  $\square$

## 8 Appendix: Connection is associative

This will be proven with string diagrams, and consequently generalises to the operation defined by Equation ?? in other Markov kernel categories.

Define

$$I_{K..} := I_K \setminus I_L \setminus I_J \quad (188)$$

$$I_{KL.} := I_K \cap I_L \setminus I_J \quad (189)$$

$$I_{K..J} := I_K \cap I_J \setminus I_L \quad (190)$$

$$I_{KLJ} := I_K \cap I_L \cap I_J \quad (191)$$

$$I_{..L} := I_L \setminus I_K \setminus I_J \quad (192)$$

$$I_{..LJ} := I_L \cap I_J \setminus I_K \quad (193)$$

$$I_{..J} := I_J \setminus I_K \setminus I_L \quad (194)$$

$$O_{K..} := O_K \setminus I_N \setminus I_J \quad (195)$$

$$O_{KL.} := O_K \cap I_L \setminus I_J \quad (196)$$

$$O_{K..J} := O_K \cap I_J \setminus I_L \quad (197)$$

$$O_{KLJ} := O_K \cap I_L \cap I_J \quad (198)$$

$$O_{..L} := O_L \setminus I_J \quad (199)$$

$$O_{..LJ} := O_L \cap I_J \quad (200)$$

Also define

$$(\mathbb{P}, l_P, O_P) := \mathbb{K} \Rightarrow \mathbb{L} \quad (201)$$

$$(\mathbb{Q}, l_Q, O_Q) := \mathbb{L} \Rightarrow \mathbb{J} \quad (202)$$

Then

$$(\mathbb{K} \Rightarrow \mathbb{L}) \Rightarrow \mathbb{J} = \mathbb{P} \Rightarrow \mathbb{J} \quad (203)$$

$$= \begin{array}{c} l_{P.} \text{---} \boxed{\mathbb{P}} \text{---} O_{P.} \\ l_{PJ} \text{---} \bullet \text{---} O_{PJ} \\ l_{.J} \text{---} \bullet \text{---} O_J \end{array} \quad (204)$$

$$= \begin{array}{c} l_{K..} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{K..} \\ l_{KL.} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{KL.} \\ l_{.L} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{K..J} \\ l_{K.J} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{KLJ} \\ l_{KLJ} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{L.} \\ l_{.LJ} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{LJ} \\ l_{..J} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_J \end{array} \quad (205)$$

$$\stackrel{perm}{=} \begin{array}{c} l_{K..} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{K..} \\ l_{KL.} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{KL.} \\ l_{K.J} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{K..J} \\ l_{KLJ} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{KLJ} \\ l_{.L} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{L.} \\ l_{.LJ} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_{LJ} \\ l_{..J} \text{---} \bullet \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} O_J \end{array} \quad (206)$$

$$= \begin{array}{c} l_{K.} \text{---} \boxed{\mathbb{K}} \text{---} O_{K.} \\ l_{KQ} \text{---} \bullet \text{---} O_{KQ} \\ l_{.Q} \text{---} \bullet \text{---} O_Q \end{array} \quad (207)$$

$$= \mathbb{K} \Rightarrow (\mathbb{L} \Rightarrow \mathbb{J}) \quad (208)$$

## 9 Appendix: String Diagram Examples

Recall the definition of *connection*:

**Definition 9.1** (Connection).

$$\mathbb{K} \Rightarrow \mathbb{L} := \begin{array}{c} \text{I}_{F^{\cdot}} \text{---} \boxed{\mathbb{K}} \text{---} \text{O}_{F^{\cdot}} \\ \text{I}_{FS} \text{---} \bullet \text{---} \text{O}_{FS} \\ \text{I}_S \text{---} \bullet \text{---} \text{O}_S \end{array} \quad (209)$$

$$:= \mathbb{J} \quad (210)$$

$$\mathbb{J}_{yqr}^{zxw} = \mathbb{K}_{yq}^{zx} \mathbb{L}_{xr}^w \quad (211)$$

Equation 209 can be broken down to the product of four Markov kernels, each of which is itself a tensor product of a number of other Markov kernels:

$$(\mathbb{J}, (\text{I}_{F^{\cdot}}, \text{I}_{FS}, \text{I}_S), (\text{O}_{F^{\cdot}}, \text{O}_{FS}, \text{O}_S)) = \left[ \begin{array}{c} \text{I}_{F^{\cdot}} \text{---} \\ \text{I}_{FS} \text{---} \bullet \\ \text{I}_S \text{---} \end{array} \right] \left[ \begin{array}{c} \boxed{\mathbb{K}} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \bullet \text{---} \\ \bullet \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{O}_{F^{\cdot}} \\ \text{O}_{FS} \\ \boxed{\mathbb{L}} \text{---} \text{O}_S \end{array} \right] \quad (212)$$

$$(213)$$

## 10 Markov variable maps and variables form a Markov category

In the following, given *arbitrary measurable sets*  $(X, \mathcal{X})$  and  $(Y, \mathcal{Y})$ , a Markov kernel is a function  $\mathbb{K} : X \times \mathcal{Y} \rightarrow [0, 1]$  such that

- For every  $A \in \mathcal{Y}$ , the function  $x \mapsto \mathbb{K}(x, A)$  is  $\mathcal{X}$ -measurable
- For every  $x \in X$ , the function  $A \mapsto \mathbb{K}(x, A)$  is a probability measure on  $(Y, \mathcal{Y})$

Note that this is a more general definition than the one used in the main paper; the version in the main paper is the restriction of this definition to finite sets.

The *delta function*  $\delta : X \rightarrow \Delta(\mathcal{X})$  is the Markov kernel defined by

$$\delta(x, A) = \begin{cases} 1 & x \in A \\ 0 & \text{otherwise} \end{cases} \quad (214)$$

Fritz (2020) defines Markov categories in the following way:

**Definition 10.1.** A Markov category  $\mathcal{C}$  is a symmetric monoidal category in which every object  $X \in \mathcal{C}$  is equipped with a commutative comonoid structure given by a comultiplication  $\text{copy}_X : X \rightarrow X \otimes X$  and a counit  $\text{del}_X : X \rightarrow I$ , depicted in string diagrams as

$$\text{del}_X := \text{---} * \text{copy}_X \quad := \text{---} \bullet \quad (215)$$

and satisfying the commutative comonoid equations

$$\text{Diagram 1} = \text{Diagram 2} \quad (216)$$

$$\begin{array}{c} \text{---} \bullet \begin{array}{l} \nearrow \\ \searrow \end{array} \begin{array}{l} \text{---} \\ * \end{array} \\ \text{---} \end{array} = \text{---} = \begin{array}{c} \text{---} \bullet \begin{array}{l} \nearrow * \\ \searrow \end{array} \\ \text{---} \end{array} \quad (217)$$

$$\begin{array}{c} \text{---} \bullet \text{---} \end{array} = \begin{array}{c} \text{---} \bullet \text{---} \end{array} \quad (218)$$

as well as compatibility with the monoidal structure

$$\begin{array}{ccc} X \otimes Y & \longrightarrow & * \\ & & \downarrow \\ & & X \longrightarrow * \end{array} \quad (219)$$

$$X \otimes Y \longrightarrow \begin{array}{c} \text{---} X \otimes Y \\ \text{---} X \otimes Y \end{array} = \begin{array}{c} X \text{---} \bullet \text{---} X \\ Y \text{---} \bullet \text{---} Y \end{array} \quad (220)$$

and the naturality of  $del$ , which means that

$$\text{---} \boxed{f} \text{---} * \text{---} * = \text{---} * \text{---} * \quad (221)$$

for every morphism  $f$ .

The category of labeled Markov kernels is the category consisting of labeled measurable sets as objects and labeled Markov kernels as morphisms. Given  $\mathbb{K} : \mathbf{X} \rightarrow \Delta(\mathbf{Y})$  and  $\mathbb{L} : \mathbf{Y} \rightarrow \Delta(\mathbf{Z})$ , sequential composition is given by

$$\mathbb{KL} : \mathbf{X} \rightarrow \Delta(\mathbf{Z}) \quad (222)$$

$$\text{defined by } (\mathbb{KL})(x, A) = \int_Y \mathbb{L}(y, A) \mathbb{K}(x, dy) \quad (223)$$

For  $\mathbb{K} : \mathbf{X} \rightarrow \Delta(\mathbf{Y})$  and  $\mathbb{L} : \mathbf{W} \rightarrow \Delta(\mathbf{Z})$ , parallel composition is given by

$$\mathbb{K} \otimes \mathbb{L} : (\mathbf{X}, \mathbf{W}) \rightarrow \Delta(\mathbf{Y}, \mathbf{Z}) \quad (224)$$

$$\text{defined by } \mathbb{K} \otimes \mathbb{L}(x, w, A \times B) = \mathbb{K}(x, A)\mathbb{L}(w, B) \quad (225)$$

The identity map is

$$\text{Id}_X : X \rightarrow \Delta(X) \quad (226)$$

$$\text{defined by } (\text{Id}_X)(x, A) = \delta(x, A) \quad (227)$$

We take an arbitrary single element labeled set  $I = (*, \{*\})$  to be the unit, which we note satisfies  $I \otimes X = X \otimes I = X$  by Lemma ??.

The swap map is given by

$$\text{swap}_{X,Y} : (X, Y) \rightarrow \Delta(Y, X) \quad (228)$$

$$\text{defined by } (\text{swap}_{X,Y})(x, y, A \times B) = \delta(x, B)\delta(y, A) \quad (229)$$

And we use the standard associativity isomorphisms for Cartesian products such that  $(A \times B) \times C \cong A \times (B \times C)$ , which in turn implies  $(X, (Y, Z)) \cong ((X, Y), Z)$ .

The copy map is given by

$$\text{copy}_X : X \rightarrow \Delta(X, X) \quad (230)$$

$$\text{defined by } (\text{copy}_X)(x, A \times B) = \delta_x(A)\delta_x(B) \quad (231)$$

and the erase map by

$$\text{del}_X : X \rightarrow \Delta(*) \quad (232)$$

$$\text{defined by } (\text{del}_X)(x, A) = \delta(*, A) \quad (233)$$

$$(234)$$

Note that the category formed by taking the underlying unlabeled sets and the underlying unlabeled morphisms is identical to the category of measurable sets and Markov kernels described in Fong (2013); Cho and Jacobs (2019); Fritz (2020).

**Theorem 10.2** (The category of labeled Markov kernels and labeled measurable sets is a Markov category). *The category described above is a Markov category.*

*Proof.*

I'm not sure how to formally argue that it is monoidal and symmetric as the relevant texts I've checked all gloss over the functors with respect to which the relevant isomorphisms should be natural, but labels with products were intentionally made to act just like sets with cartesian products which are symmetric monoidal

Equations 216 to 221 are known to be satisfied for the underlying unlabeled Markov kernels. We need to show is that they hold given our stricter criterion of labeled Markov kernel equality; that the underlying kernels *and the label sets* match. It is sufficient to check the label sets only.

□

## References

- G. Chiribella, Giacomo D’Ariano, and P. Perinotti. Quantum Circuit Architecture. *Physical review letters*, 101:060401, September 2008. doi: 10.1103/PhysRevLett.101.060401.
- Kenta Cho and Bart Jacobs. Disintegration and Bayesian inversion via string diagrams. *Mathematical Structures in Computer Science*, 29(7):938–971, August 2019. ISSN 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488.
- A. Philip Dawid. Decision-theoretic foundations for statistical causality. *arXiv:2004.12493 [math, stat]*, April 2020. URL <http://arxiv.org/abs/2004.12493>. arXiv: 2004.12493.
- M. P. Ershov. Extension of Measures and Stochastic Equations. *Theory of Probability & Its Applications*, 19(3):431–444, June 1975. ISSN 0040-585X. doi: 10.1137/1119053. URL <https://epubs.siam.org/doi/abs/10.1137/1119053>. Publisher: Society for Industrial and Applied Mathematics.
- R.P. Feynman. *The Feynman lectures on physics*. Le cours de physique de Feynman. Interditions, France, 1979. ISBN 978-2-7296-0030-3. INIS Reference Number: 13648743.
- Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv:1301.6201 [math]*, January 2013. URL <http://arxiv.org/abs/1301.6201>. arXiv: 1301.6201.
- Tobias Fritz. A synthetic approach to Markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, August 2020. ISSN 0001-8708. doi: 10.1016/j.aim.2020.107239. URL <https://www.sciencedirect.com/science/article/pii/S0001870820302656>.
- D. Heckerman and R. Shachter. Decision-Theoretic Foundations for Causal Reasoning. *Journal of Artificial Intelligence Research*, 3:405–430, December 1995. ISSN 1076-9757. doi: 10.1613/jair.202. URL <https://www.jair.org/index.php/jair/article/view/10151>.
- M. A. Hernán and S. L. Taubman. Does obesity shorten life? The importance of well-defined interventions to answer causal questions. *International Journal of Obesity*, 32(S3):S8–S14, August 2008. ISSN 1476-5497. doi: 10.1038/ijo.2008.82. URL <https://www.nature.com/articles/ijo200882>.
- Alan Hájek. What Conditional Probability Could Not Be. *Synthese*, 137(3): 273–323, December 2003. ISSN 1573-0964. doi: 10.1023/B:SYNT.0000004904.91112.16. URL <https://doi.org/10.1023/B:SYNT.0000004904.91112.16>.
- Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In Mikoaj Bojaczek and Alex Simpson, editors, *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing, 2019. ISBN 978-3-030-17127-8.

- Alfred Korzybski. *Science and sanity; an introduction to Non-Aristotelian systems and general semantics*. Lancaster, Pa., New York City, The International Non-Aristotelian Library Publishing Company, The Science Press Printing Company, distributors, 1933. URL <http://archive.org/details/sciencesanityint00korz>.
- Finnian Lattimore and David Rohde. Causal inference with Bayes rule. *arXiv:1910.01510 [cs, stat]*, October 2019a. URL <http://arxiv.org/abs/1910.01510>. arXiv: 1910.01510.
- Finnian Lattimore and David Rohde. Replacing the do-calculus with Bayes rule. *arXiv:1906.07125 [cs, stat]*, December 2019b. URL <http://arxiv.org/abs/1906.07125>. arXiv: 1906.07125.
- David Lewis. Causal decision theory. *Australasian Journal of Philosophy*, 59(1): 5–30, March 1981. ISSN 0004-8402. doi: 10.1080/00048408112340011. URL <https://doi.org/10.1080/00048408112340011>.
- Karl Menger. Random Variables from the Point of View of a General Theory of Variables. In Karl Menger, Bert Schweizer, Abe Sklar, Karl Sigmund, Peter Gruber, Edmund Hlawka, Ludwig Reich, and Leopold Schmetterer, editors, *Selecta Mathematica: Volume 2*, pages 367–381. Springer, Vienna, 2003. ISBN 978-3-7091-6045-9. doi: 10.1007/978-3-7091-6045-9\_31. URL [https://doi.org/10.1007/978-3-7091-6045-9\\_31](https://doi.org/10.1007/978-3-7091-6045-9_31).
- Scott Mueller, Ang Li, and Judea Pearl. Causes of Effects: Learning individual responses from population data. *arXiv:2104.13730 [cs, stat]*, May 2021. URL <http://arxiv.org/abs/2104.13730>. arXiv: 2104.13730.
- Paul F. Christiano. EDT vs CDT, September 2018. URL <https://sideways-view.com/2018/09/19/edt-vs-cdt/>.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2 edition, 2009.
- Judea Pearl. Does Obesity Shorten Life? Or is it the Soda? On Non-manipulable Causes. *Journal of Causal Inference*, 6(2), 2018. doi: 10.1515/jci-2018-2001. URL <https://www.degruyter.com/view/j/jci.2018.6.issue-2/jci-2018-2001/jci-2018-2001.xml>.
- Donald B. Rubin. Causal Inference Using Potential Outcomes. *Journal of the American Statistical Association*, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880. URL <https://doi.org/10.1198/016214504000001880>.
- Peter Selinger. A survey of graphical languages for monoidal categories. *arXiv:0908.3347 [math]*, 813:289–355, 2010. doi: 10.1007/978-3-642-12821-9\_4. URL <http://arxiv.org/abs/0908.3347>. arXiv: 0908.3347.



- Eyal Shoham. The association of body mass index with health outcomes: causal, inconsistent, or confounded? *American Journal of Epidemiology*, 170(8):957–958, October 2009. ISSN 1476-6256. doi: 10.1093/aje/kwp292.
- Jin Tian and Judea Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics and Artificial Intelligence*, 28(1):287–313, October 2000. ISSN 1573-7470. doi: 10.1023/A:1018912507879. URL <https://doi.org/10.1023/A:1018912507879>.
- Jin Tian and Judea Pearl. A general identification condition for causal effects. In *Aaai/iaai*, pages 567–573, July 2002.
- J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Theory of games and economic behavior. Princeton University Press, Princeton, NJ, US, 1944.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England, 1950.
- Paul Weirich. Causal Decision Theory. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2016 edition, 2016. URL <https://plato.stanford.edu/archives/win2016/entries/decision-causal/>.

## Appendix: