

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

David Johnston

February 10, 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Our approach . . . . .	3
<b>2</b>	<b>Probability</b>	<b>4</b>
2.1	Section outline . . . . .	4
2.1.1	Brief outline of probability gap models . . . . .	4
2.2	Standard probability theory . . . . .	6
2.3	Not quite standard probability theory . . . . .	7
2.4	Probabilistic models for causal inference . . . . .	8
2.5	Probability sets . . . . .	9
2.6	Semidirect product and almost sure equality . . . . .	10
2.7	Maximal probability sets and valid conditionals . . . . .	12
2.7.1	Conditional independence . . . . .	13
<b>3</b>	<b>Decision problems</b>	<b>14</b>
3.1	Conditional probability models . . . . .	16
3.2	Example: invalidity . . . . .	17
3.3	Response conditionals . . . . .	18
3.4	Response conditionals and potential outcomes . . . . .	18
3.5	Randomness pushbacks . . . . .	19
3.5.1	Choices aren't always known . . . . .	20
3.6	Other decision theoretic causal models . . . . .	21
<b>4</b>	<b>When do response conditionals exist?</b>	<b>22</b>
4.1	Repeatable experiments . . . . .	23
4.2	Consequence contractibility . . . . .	24
4.3	Repeatable consequence conditionals exist iff a model is causally contractible . . . . .	27
4.4	Modelling different measurement procedures . . . . .	32
4.5	Example: commutativity of exchange in the context of treatment choices . . . . .	32

4.6	Causal consequences of non-deterministic variables . . . . .	36
4.7	Intersubjective causal consequences . . . . .	37
<b>5</b>	<b>Appendix, needs to be organised</b>	<b>38</b>
5.1	Existence of conditional probabilities . . . . .	38
5.2	Validity . . . . .	41
5.3	Conditional independence . . . . .	44
5.4	Extended conditional independence . . . . .	44

# 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).

Causal inference work undertaken using either approaches is often directed towards determining the likely effects of different actions that could be taken. This kind of application is strongly suggested by the terminology of “interventions” and “potential outcomes”. However, if we want to reason clearly about using data to inform choices of actions, suggestive terminology is not enough to underpin a sound understanding of the correspondence between causal inference models and action selection problems.

As a motivating example, Hernán and Taubman (2008) observed that many epidemiological papers have been published estimating the “causal effect” of body mass index. However, Hernán argued, because there are many different *actions* that might affect body mass index, the potential outcomes associated with body mass index themselves are ill-defined. This would not be particularly problematic if we regarded the search for treatment effects as an endeavour entirely separate from questions of choosing actions – it’s only because we want potential outcomes to tell us something about effects of actions that a many-to-one relationship between “actions” and “causal variables” becomes troublesome.

In a response to Hernán’s observation, Pearl (2018) argues that *interventions* (by which we mean the operation defined by a causal graphical model) are well defined, but by default they describe “virtual interventions” or “ideal, atomic interventions”, and more complicated intervention operations may be needed for a good models of real actions. In this view, the relationship between interventions and actions is clearly not straightforward. In particular, one might wonder what standard we can use to determine if an action is “ideal” and “atomic” (apart from the question begging standard of agreement with interventions in a given causal graphical model).

In another response, Shahar (2009) argued that a properly specified intervention on body mass index will necessarily yield a conclusion that intervention on body mass index has no effect at all on anything apart from body mass index itself. If this is accepted, then it might seem that there is a whole body of literature devoted to estimating a “causal effect” that is necessarily equal to zero! It seems that there is a need to clarify the relationship between actions and causal effects.

## 1.1 Our approach

We focus our attention on the following problem: given an experiment with sequence of variable pairs  $(X_i, Y_i)$  and a collection of decision functions  $A$  that the experimenter may choose, when is there a unique probabilistic function  $H \times X \rightarrow Y$  that defines the “causal consequence” of  $X_i$  on  $Y_i$  for all  $i$ ? We choose the setting of repeatable experiments because causal inferences are, in practice, usually drawn from sequential data generated from repeatable experiments or sampling procedures. Here the set  $H$  represents a set of hypotheses that – under some choices of decision function – becomes deterministic in the limit of infinite data.

To answer this question, we require clarity on what we mean by “variable”, and so we begin with an explanation of a theory of variables. This theory is close to a standard account, but we are somewhat more explicit than usual about the relationship between variables and measurement procedures.

We then address the problem of creating probabilistic models of variables that permit us to evaluate different choices of decision function. To this end, we introduce *probability sets*, which can be thought of as partially specified probability models, and *probability gap models* which can be thought of as probability sets along with a selection of choices for “filling the gap”. We typically specify probability sets with conditional probabilities, and the criterion of *validity* ensures that we don’t inadvertently end up specifying empty probability sets – a problem that can arise in the context of modelling interventions on body mass index.

We then prove that a condition we call *causal contractibility* is equivalent to the existence of repeatable causal consequences. This result is akin to De Finetti’s theorem showing the equivalence between exchangeability and the existence of a “unique but unknown” distribution such that all variables are independent and identically distributed according to it.

Finally, we consider the question of *when* causal contractibility might be reasonable to assume. This requires us to consider the measurement processes associated with the variables that we think may or may not be causally contractible. We suggest two sufficient sets of conditions:

1. The  $X_i$ s are deterministically equal for all experimental units
2. The evidence is symmetric for all experimental units, the order of experimentation is irrelevant and the  $X_i$  are deterministic given the choice of decision function

3. There exists some  $(D_i, Y_i)_{i \in M}$  causally contractible and  $Y_i \perp\!\!\!\perp D_i | HX_i$

The first two conditions for causal contractibility are very unlikely to hold for body mass index, both because body mass index may not be functionally related to the available actions and because body mass index cannot be deterministically controlled. The third condition might hold for body mass index with respect to some sets of actions, but this is an empirical question.

## 2 Probability

### 2.1 Section outline

following section hasn't been revised

This section introduces the mathematical foundations used throughout the rest of the paper. The first subsection briefly introduces probability theory, which is likely to be familiar to many readers, as well as how string diagrams can be used to represent probabilistic functions (or *Markov kernels*), which may be less familiar. We use string diagrams for probabilistic reasoning in a number of places, and this section is intended to help interpret mathematical statements in this form.

The second subsection discusses the interpretation of probabilistic variables. Our formalisation of probabilistic variables is standard – we define them as measurable functions on a fundamental probability set  $\Omega$ . We discuss how this formalisation can be connected to statements about the real world via *measurement processes*, and distinguishes observed variables (which are associated with measurement processes) from unobserved variables (which are not associated with measurement processes). This section is not part of the mathematical theory of probability gap models, but it is relevant when one wants to apply this theory to real problems or to understand how the theory of probability gap models relates to other theories of causal inference.

Finally, we introduce *probability gap models*. Probability gap models are a generalisation of probability models, and to understand the rest of this paper a reader needs to understand what a probability gap model is, how we define the common kinds of probability gap models used in this paper and what conditional probabilities and conditional independence statements mean for probability gap models.

#### 2.1.1 Brief outline of probability gap models

We consider a probability model to be a probability space  $(\Omega, \mathcal{F}, \mu)$  along with a collection of random variables. However, if I want to use probabilistic models to support decision making, then I need function from options to probability models. For example, suppose I have two options  $A = \{0, 1\}$ , and I want to compare these options based on what I expect to happen if I choose them. If I choose option 0, then I can (perhaps) represent my expectations about the consequences with a probability model, and if I choose option 1 I can represent my

expectations about the consequences with a different probability model. I can compare the two consequences, then decide which option seems to be better. To make this comparison, I have used a function from elements of  $A$  to probability models. A function that takes elements of some set as inputs (which may or may not be decisions) and returns probability models is a *probability gap model*, and the set of inputs it accepts is a *probability gap*.

We are particularly interested in probability gap models where the consequences of all inputs share some marginal or conditional probabilities. The simplest example of a model like this can be represented by a probability distribution  $\mathbb{P}^X$  for some variable  $X : \Omega \rightarrow X$ . Such a probability distribution is consistent with many base measures on the fundamental probability set  $\Omega$ , and so we can consider the choice of base measure to be a probability gap. Not every probability distribution over  $X$  can define a probability gap model in this way. In particular, we need  $\mathbb{P}^X$  to assign probability 0 to outcomes that are mathematically impossible according to the definition of  $X$  to ensure that there is some base measure that features  $\mathbb{P}^X$  as a marginal. We call probability gap models represented by probability distributions *order 0 probability gap models*.

Higher order probability gap models can be represented by conditional probabilities  $\mathbb{P}^{Y|X}$  or pairs of conditional probabilities  $\{\mathbb{P}^{X|W}, \mathbb{P}^{Z|WX^Y}\}$ , which we call *order 1* and *order 2* models respectively. Decision functions in data-driven decision problems correspond to probability gaps in order 2 models, as we discuss in Section ??, which makes this type of model particularly interesting for our purposes. We also require these to be valid, and we define conditions for validity and prove that they are sufficient to ensure that models represented by conditional probabilities can in fact be mapped to base measures on the fundamental probability set.

A conditional independence statement in a probability gap model means that the corresponding conditional independence statement holds for all base measures in the range of the function defined by the model. It is possible to deduce conditional independences from “independences” in the conditional probabilities that we use to represent these models, and conditional independences can imply the existence of conditional probabilities with certain independence properties.

We can consider causal Bayesian networks to represent order 2 probability gap models. That is, a causal Bayesian network represents a function  $\mathbb{P}$  that take inserts from some set  $A$  of conditional probabilities and returns a probability model, and it does so in such a way that there are a pair of conditional probabilities  $\{\mathbb{P}^{X|W}, \mathbb{P}^{Z|WX^Y}\}$  shared by all models in the codomain of  $\mathbb{P}$ . The observational distribution is the value of  $\mathbb{P}(\text{obs})$  for some *observational insert*  $\text{obs} \in A$ , and other choices of inserts yield interventional distributions. Defining causal Bayesian networks in this manner resolves two areas of difficulty with causal Bayesian networks. First, under the standard definition of causal Bayesian networks interventional probabilities may fail to exist; with our perspective we can see that this arises due to misunderstanding the domain of  $\mathbb{P}$ . Secondly, there may be multiple distributions that differ in important ways that all satisfy the standard definition of “interventional distributions”. The one-to-

many relationship between observations and interventions is a basic challenge of causal inference, the problem arises when this relationship is obscured by calling multiple different things “the interventional distribution”. If we consider causal Bayesian networks to represent order 2 probability gap models, we avoid doing this.

end previous section hasn't been revised

## 2.2 Standard probability theory

**Definition 2.1** (Probability measure). Given a measure space  $(X, \mathcal{X})$ , a probability measure is a  $\sigma$ -additive function  $\mu : \mathcal{X} \rightarrow [0, 1]$  such that  $\mu(\emptyset) = 0$  and  $\mu(X) = 1$ . We write  $\Delta(X)$  for the set of all probability measures on  $(X, \mathcal{X})$ .

**Definition 2.2** (Markov kernel). Given measure spaces  $(X, \mathcal{X})$ ,  $(Y, \mathcal{Y})$   $Y : \Omega \rightarrow Y$ , a Markov kernel  $\mathbb{Q} : X \rightarrow Y$  is a map  $Y \times \mathcal{X} \rightarrow [0, 1]$  such that

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

**Definition 2.3** (Delta measure). Given a measureable space  $(X, \mathcal{X})$  and  $x \in X$ ,  $\delta_x \in \Delta(X)$  is the measure defined by  $\delta_x(A) = \llbracket x \in A \rrbracket$ .

**Definition 2.4** (Probability space). A probability space is a triple  $(\mu, \Omega, \mathcal{F})$ , where  $\mu$  is a base measure on  $\mathcal{F}$ .

**Definition 2.5** (Variable). Given a measureable space  $(\Omega, \mathcal{F})$  and a set of values  $(X, \mathcal{X})$ , an  $X$ -valued variable is a measurable function  $X : \Omega \rightarrow X$ .

**Definition 2.6** (Sequence of variables). Given a measureable space  $(\Omega, \mathcal{F})$  and two variables  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ ,  $(X, Y) : \Omega \rightarrow X \times Y$  is the variable  $\omega \mapsto (X(\omega), Y(\omega))$ .

**Definition 2.7** (Marginal distribution with respect to a probability space). Given a probability space  $(\mu, \Omega, \mathcal{F})$  and a variable  $X : \Omega \rightarrow (X, \mathcal{X})$ , we can define the *marginal distribution* of  $X$  with respect to  $\mu$ ,  $\mu^X : \mathcal{X} \rightarrow [0, 1]$  by  $\mu^X(A) := \mu(X \bowtie A)$  for any  $A \in \mathcal{X}$ .

**Lemma 2.8** (Marginal distribution as a kernel product). *Given a probability space  $(\mu, \Omega, \mathcal{F})$  and a variable  $X : \Omega \rightarrow (X, \mathcal{X})$ , define  $\mathbb{F}_X : \Omega \rightarrow X$  by  $\mathbb{F}_X(A|\omega) = \delta_{X(\omega)}(A)$ , then*

$$\mu^X = \mu \mathbb{F}_X \tag{1}$$

*Proof.* Consider any  $A \in \mathcal{X}$ .

$$\mu \mathbb{F}_X(A) = \int_{\Omega} \delta_{X(\omega)}(A) d\mu(\omega) \tag{2}$$

$$= \int_{X^{-1}(A)} d\mu(\omega) \tag{3}$$

$$= \mu^X(A) \tag{4}$$

□

### 2.3 Not quite standard probability theory

Instead of having probability distributions and Markov kernels as two different kinds of thing, we can identify probability distributions with Markov kernels whose domain is a one element set  $\{*\}$ .

**Definition 2.9** (Probability measures as Markov kernels). Given  $(X, \mathcal{X})$  and  $\mu \in \Delta(X)$ , the Markov kernel  $\mathbb{K} : \{*\} \rightarrow X$  given by  $\mathbb{K}(A|*) = \mu(A)$  for all  $A \in \mathcal{X}$  is the Markov kernel associated with the probability measure  $\mu$ . We will use probability measures and their associated Markov kernels interchangeably, as it is transparent how to get from one to another.

**Definition 2.10** (Regular conditional distribution). Given a probability space  $(\mu, \Omega)$  and variables  $X : \Omega \rightarrow X, Y : \Omega \rightarrow Y$ , the probability of  $Y$  given  $X$  is any Markov kernel  $\mu^{Y|X} : X \rightarrow Y$  such that

$$\mu^{XY}(A \times B) = \int_A \mu^{Y|X}(B|x) d\mu^X(x) \quad \forall A \in \mathcal{X}, B \in \mathcal{Y} \quad (5)$$

$$\iff \quad (6)$$

$$\mu^{XY} = \begin{array}{c} \text{X} \\ \nearrow \\ \triangleleft \mu^X \end{array} \quad \bullet \quad \begin{array}{c} \boxed{\mu^{Y|X}} \\ \text{Y} \end{array} \quad (7)$$

We define higher order conditionals as “conditionals of conditionals”

**Definition 2.11** (Regular higher order conditionals). Given a probability space  $(\mu, \Omega)$  and variables  $X : \Omega \rightarrow X, Y : \Omega \rightarrow Y$  and  $Z : \Omega \rightarrow Z$ , a higher order conditional  $\mu^{Z|(Y|X)} : X \times Y \rightarrow Z$  is any Markov kernel such that, for some  $\mu^{Y|X}$ ,

$$\mu^{ZY|X}(B \times C|x) = \int_B \mu^{Z|(Y|X)}(C|x, y) \mu^{Y|X}(dy|x) \quad (8)$$

$$\iff \mu^{ZY|X} = \begin{array}{c} \text{Y} \\ \nearrow \\ \text{X} \longrightarrow \bullet \quad \begin{array}{c} \boxed{\mu^{Y|X}} \\ \text{Z} \end{array} \end{array} \quad (9)$$

Higher order conditionals are useful because  $\mu^{Z|(Y|X)}$  is a version of  $\mu^{X|YX}$ , so if we're given  $\mu^{ZY|X}$  and we can find some  $\mu^{Z|(Y|X)}$  then we have a version of  $\mu^{X|YX}$ . This also hold for conditional with respect to probability sets, which we will introduce later (Theorem 5.4).

Furthermore, given regular  $\mu^{XY|Z}$  and  $X, Y$  standard measurable, it has recently been proven that a regular higher order conditional  $\mu^{Z|(Y|X)}$  exists Bogachev and Malofeev (2020), Theorem 3.5. See also Theorem 5.3 for the extension of this theorem to probability sets.

## 2.4 Probabilistic models for causal inference

following section hasn't been revised

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  $\mathcal{S}$  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”.

Given a model skeleton, a common approach to attaching a probabilistic model involves defining a base measure  $\mu$  on  $(\Omega, \mathcal{F})$  which yields a probability space  $(\Omega, \mathcal{F}, \mu)$ . For causal inference, we need a to generalise this approach, because we need to handle *choices*. If I have different options I can choose, and I want to use a model to compare the options according to some criteria, then I need a model that can accept a choice and output the expected result of that choice. According to this model, anything that we consider a “consequence of a choice” doesn’t have a definite probability, because it depends on the choice we make.

In general, we might have arbitrary sets of choices that map to probabilistic models in an arbitrary way. However, we are here interested in a simpler case: we suppose that there are a number of points at which we can act, and prior to acting we can observe some variables, and we are able to choose probabilistic maps from observations to acts. We also assume that, given the same observation and the same act, the same consequence is expected. That is, the consequences do not depend directly way on the choice of map from observations to acts.

These assumptions together imply that our model should contain a number of fixed conditional probabilities – the probabilities of consequences given observations and acts – and a number of “choosable” conditional probabilities – the probabilities of acts given observations. The fixed conditional probabilities form a probability model with *gaps*, and those gaps correspond to choices we can make. When we combine the fixed conditional probabilities and a choice of a conditional probability for each gap, we get a regular probability model. The terminology of “probability gaps” comes from Hájek (2003).

To restate our general approach: we model decision problems with a collection of fixed conditional probabilities and a collection of choosable conditional probabilities, and combine the fixed conditionals with particular choices to get a probability measure. Two issues present themselves here: firstly, what *is* a collection of conditional probabilities without a fixed underlying probability measure? Secondly, we need to ensure that our chosen collection of conditional probabilities actually does induce a probability model. We address these questions with *probability sets*. A probability set is a collection of probability measures on  $(\Omega, \mathcal{F})$ , and we identify a collection of conditional probabilities with the set of probability measures that induce those conditional probabilities. We then define



an operation  $\odot$  for combining conditional probabilities, and a criterion of *validity* such that a collection of valid conditional probabilities recursively combined using  $\odot$  is guaranteed to corresponds to a non-empty probability set.

Previous section hasn't been revised

## 2.5 Probability sets

A probability set is a set of probability measures. This section establishes a number of useful properties of conditional probability with respect to probability sets. Unlike conditional probability with respect to a probability space, conditional probabilities don't always exist for probability sets. Where they do, however, they are almost surely unique and we can marginalise and disintegrate them to obtain other conditional probabilities with respect to the same probability set.

**Definition 2.12** (Probability set). A probability set  $\mathbb{P}_{\{\}}$  on  $(\Omega, \mathcal{F})$  is a collection of probability measures on  $(\Omega, \mathcal{F})$ . In other words it is a subset of  $\mathcal{P}(\Delta(\Omega))$ , where  $\mathcal{P}$  indicates the power set.

Given a probability set  $\mathbb{P}_{\{\}}$ , we define marginal and conditional probabilities as probability measures and Markov kernels that satisfy Definitions 2.7 and 2.10 respectively for *all* base measures in  $\mathbb{P}_{\{\}}$ . There are generally multiple Markov kernels that satisfy the properties of a conditional probability with respect to a probability set, and this definition ensures that marginal and conditional probabilities are “almost surely” unique (Definition 2.18) with respect to probability sets.

**Definition 2.13** (Marginal probability with respect to a probability set). Given a sample space  $(\Omega, \mathcal{F})$ , a variable  $X : \Omega \rightarrow X$  and a probability set  $\mathbb{P}_{\{\}}$ , the marginal distribution  $\mathbb{P}_{\{\}}^X = \mathbb{P}_{\alpha}^X$  for any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$  if a distribution satisfying this condition exists. Otherwise, it is undefined.

**Definition 2.14** (Regular conditional distribution with respect to a probability set). Given a fundamental probability set  $\Omega$  variables  $X : \Omega \rightarrow X$  and  $Y : \Omega \rightarrow Y$  and a probability set  $\mathbb{P}_{\{\}}$ , a conditional  $\mathbb{P}_{\{\}}^{Y|X}$  is any Markov kernel  $X \rightarrow Y$  such that  $\mathbb{P}_{\{\}}^{Y|X}$  is an  $Y|X$  conditional probability of  $\mathbb{P}_{\alpha}$  for all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$ . If no such Markov kernel exists,  $\mathbb{P}_{\{\}}^{Y|X}$  is undefined.

**Definition 2.15** (Regular higher order conditional with respect to a probability set). Given a fundamental probability set  $\Omega$ , variables  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and  $Z : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_{\{\}}$ , if  $\mathbb{P}_{\{\}}^{ZY|X}$  exists then a higher order conditional  $\mathbb{P}_{\{\}}^{Z|(Y|X)}$  is any Markov kernel  $X \times Y \rightarrow Z$  that is a higher order conditional of some version of  $\mathbb{P}_{\{\}}^{ZY|X}$ . If no  $\mathbb{P}_{\{\}}^{ZY|X}$  exists,  $\mathbb{P}_{\{\}}^{Z|(Y|X)}$  is undefined.

Under the assumption of standard measurable spaces, the existence of a conditional probability  $\mathbb{P}_{\{\}}^{ZY|X}$  implies the existence of a higher order conditional

$\mathbb{P}_{\{\}}^{Z|(Y|X)}$  with respect to the same probability set (Theorem 5.3).  $\mathbb{P}_{\{\}}^{Z|(Y|X)}$  is in turn a version of the conditional  $\mathbb{P}_{\{\}}^{Z|YX}$  (Theorem 5.4). Thus, from the existence of  $\mathbb{P}_{\{\}}^{ZY|X}$  we can derive the existence of  $\mathbb{P}_{\{\}}^{Z|YX}$ .

## 2.6 Semidirect product and almost sure equality

The operation used in Equation 7 that combines  $\mu^X$  and  $\mu^{Y|X}$  is something we will use repeatedly, so we call it the *semidirect product* and give it the symbol

$\odot$ . We also define a notion of almost sure equality with respect to  $\odot$ :  $\mathbb{K} \stackrel{\mu^X}{\cong} \mathbb{L}$  if  $\mu^X \odot \mathbb{K} = \mu^X \odot \mathbb{L}$ . Thus if two terms are almost surely equal, they are substitutable when they both appear in a semidirect product.

**Definition 2.16** (Semidirect product). Given  $\mathbb{K} : X \rightarrow Y$  and  $\mathbb{L} : Y \times X \rightarrow Z$ , define the copy-product  $\mathbb{K} \odot \mathbb{L} : X \rightarrow Y \times Z$  as

$$\mathbb{K} \odot \mathbb{L} := \text{copy}_X(\mathbb{K} \otimes \text{id}_X)(\text{copy}_Y \otimes \text{id}_X)(\text{id}_Y \otimes \mathbb{L}) \quad (10)$$

$$= \begin{array}{c} \text{Diagram: } X \text{ connects to a node, which connects to box } \mathbb{K}. \text{ Box } \mathbb{K} \text{ connects to a node, which connects to box } \mathbb{L}. \text{ Box } \mathbb{L} \text{ connects to } Z. \text{ A curved line connects the first node to box } \mathbb{L}. \text{ Box } \mathbb{L} \text{ also has an output to } Y. \end{array} \quad (11)$$

$$\iff \quad (12)$$

$$(\mathbb{K} \odot \mathbb{L})(A \times B|x) = \int_A \mathbb{L}(B|y, x) \mathbb{K}(dy|x) \quad A \in \mathcal{Y}, B \in \mathcal{Z} \quad (13)$$

**Lemma 2.17** (Semidirect product is associative). Given  $\mathbb{K} : X \rightarrow Y$ ,  $\mathbb{L} : Y \times X \rightarrow Z$  and  $\mathbb{M} : Z \times Y \times X \rightarrow W$

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (14)$$

$$(15)$$

*Proof.*

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \begin{array}{c} \text{Diagram: } X \text{ connects to a node, which connects to box } \mathbb{K}. \text{ Box } \mathbb{K} \text{ connects to a node, which connects to box } \mathbb{L}. \text{ Box } \mathbb{L} \text{ connects to a node, which connects to box } \mathbb{M}. \text{ Box } \mathbb{M} \text{ has three outputs to } X, Y, \text{ and } Z. \end{array} \quad (16)$$

$$= \begin{array}{c} \text{Diagram: } X \text{ connects to a node, which connects to box } \mathbb{K}. \text{ Box } \mathbb{K} \text{ connects to a node, which connects to box } \mathbb{L}. \text{ Box } \mathbb{L} \text{ connects to a node, which connects to box } \mathbb{M}. \text{ Box } \mathbb{M} \text{ has three outputs to } X, Y, \text{ and } Z. \end{array} \quad (17)$$

$$= \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (18)$$

□

Two Markov kernels are almost surely equal with respect to a probability set  $\mathbb{P}_{\Omega}$  if the semidirect product  $\odot$  of all marginal probabilities of  $\mathbb{P}_{\Omega}^{\mathbf{X}}$  with each Markov kernel is identical.

**Definition 2.18** (Almost sure equality). Two Markov kernels  $\mathbb{K} : X \rightarrow Y$  and  $\mathbb{L} : X \rightarrow Y$  are almost surely equal  $\stackrel{\mathbb{P}_{\Omega}}{\cong}$  with respect to a probability set  $\mathbb{P}_{\Omega}$  and variable  $\mathbf{X} : \Omega \rightarrow X$  if for all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$ ,

$$\mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{K} = \mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{L} \quad (19)$$

**Lemma 2.19** (Conditional probabilities are almost surely equal). *If  $\mathbb{K} : X \rightarrow Y$  and  $\mathbb{L} : X \rightarrow Y$  are both versions of  $\mathbb{P}_{\Omega}^{Y|\mathbf{X}}$  then  $\mathbb{K} \stackrel{\mathbb{P}_{\Omega}}{\cong} \mathbb{L}$*

*Proof.* For all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$

$$\mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{K} = \mathbb{P}_{\alpha}^{\mathbf{X}Y} \quad (20)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{L} \quad (21)$$

□

**Lemma 2.20** (Substitution of almost surely equal Markov kernels). *Given  $\mathbb{P}_{\Omega}$ , if  $\mathbb{K} : X \times Y \rightarrow Z$  and  $\mathbb{L} : X \times Y \rightarrow Z$  are almost surely equal  $\mathbb{K} \stackrel{\mathbb{P}_{\Omega}}{\cong} \mathbb{L}$ , then for any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$*

$$\mathbb{P}_{\alpha}^{Y|\mathbf{X}} \odot \mathbb{K} \stackrel{a.s.}{\cong} \mathbb{P}_{\alpha}^{Y|\mathbf{X}} \odot \mathbb{L} \quad (22)$$

*Proof.* For any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$

$$\mathbb{P}_{\alpha}^{\mathbf{X}Y} \odot \mathbb{K} = (\mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{P}_{\Omega}^{Y|\mathbf{X}}) \odot \mathbb{K} \quad (23)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\Omega}^{Y|\mathbf{X}} \odot \mathbb{K}) \quad (24)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\Omega}^{Y|\mathbf{X}} \odot \mathbb{L}) \quad (25)$$

□

**Lemma 2.21** (Semidirect product of conditionals is a joint conditional). *Given a probability set  $\mathbb{P}_{\Omega}$  on  $(\Omega, \mathcal{F})$  along with conditional probabilities  $\mathbb{P}_{\Omega}^{Y|\mathbf{X}}$  and  $\mathbb{P}_{\Omega}^{Z|\mathbf{X}Y}$ ,  $\mathbb{P}_{\Omega}^{YZ|\mathbf{X}}$  exists and is equal to*

$$\mathbb{P}_{\Omega}^{YZ|\mathbf{X}} = \mathbb{P}_{\Omega}^{Y|\mathbf{X}} \odot \mathbb{P}_{\Omega}^{Z|\mathbf{X}Y} \quad (26)$$

$$(27)$$

*Proof.* By definition, for any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$

$$\mathbb{P}_{\alpha}^{\mathbf{X}YZ} = \mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{P}_{\alpha}^{YZ|\mathbf{X}} \quad (28)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\alpha}^{Y|\mathbf{X}} \odot \mathbb{P}_{\alpha}^{Z|\mathbf{X}Y}) \quad (29)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\Omega}^{Y|\mathbf{X}} \odot \mathbb{P}_{\Omega}^{Z|\mathbf{X}Y}) \quad (30)$$

□

## 2.7 Maximal probability sets and valid conditionals

So far we have defined probability sets and conditional probabilities as Markov kernels that can sometimes be derived from a probability set. We are often interested in working in the opposite direction: at the outset, we have a conditional probability and we want to reason about the largest probability set admitting this conditional probability. We call this a *maximal probability set*.

We need to be a little bit careful when we proceed in this fashion: we can't take an arbitrary Markov kernel  $\kappa : X \rightarrow Y$  and declare it to be a conditional probability  $\mathbb{P}_{\{\}}^{Y|X}$  for some  $X : \Omega \rightarrow X$  and  $Y : \Omega \rightarrow Y$  and a maximal probability set  $\mathbb{P}_{\{\}}$ . The reason for this is that some collections of variables cannot have arbitrary conditional probabilities, and so  $\mathbb{P}_{\{\}}$  may in fact be the empty set. We address this with the notion of validity; a *valid distribution* is a distribution associated with a particular variable that defines a nonempty set of base measures on  $\Omega$  (Theorem 5.9), and *valid conditionals* are a set of conditional probabilities closed under  $\odot$  and reducing to valid distributions when conditioning on a trivial variable (Lemma 5.12).

Consider, for example,  $\Omega = \{0, 1\}$  with  $X = (Z, Z)$  for  $Z := \text{id}_{\Omega}$  and any measure  $\kappa \in \Delta(\{0, 1\}^2)$  such that  $\kappa(\{1\} \times \{0\}) > 0$ . Note that  $X^{-1}(\{1\} \times \{0\}) = Z^{-1}(\{1\}) \cap Z^{-1}(\{0\}) = \emptyset$ . Thus for any probability measure  $\mu \in \Delta(\{0, 1\})$ ,  $\mu^X(\{1\} \times \{0\}) = \mu(\emptyset) = 0$  and so  $\kappa$  cannot be the marginal distribution of  $X$  for any base measure at all.

**Definition 2.22** (Valid distribution). Given  $(\Omega, \mathcal{F})$  and a variable  $X : \Omega \rightarrow X$ , an  $X$ -valid probability distribution is any probability measure  $\mathbb{K} \in \Delta(X)$  such that  $X^{-1}(A) = \emptyset \implies \mathbb{K}(A) = 0$  for all  $A \in \mathcal{X}$ .

**Definition 2.23** (Valid conditional). Given  $(\Omega, \mathcal{F})$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  a  $Y|X$ -valid conditional probability is a Markov kernel  $\mathbb{L} : X \rightarrow Y$  such that:

$$\forall B \in \mathcal{Y}, x \in X : (X, Y) \bowtie \{x\} \times B = \emptyset \implies (\mathbb{L}(B|x) = 0) \vee (X \bowtie \{x\} = \emptyset) \quad (31)$$

**Definition 2.24** (Maximal probability set). Given  $(\Omega, \mathcal{F})$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and a  $Y|X$ -valid conditional probability  $\mathbb{L} : X \rightarrow Y$  the maximal probability set  $\mathbb{P}_{\{\}}^{Y|X[M]}$  associated with  $\mathbb{L}$  is the probability set such that for all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$ ,  $\mathbb{L}$  is a version of  $\mathbb{P}_{\alpha}^{Y|X}$ .

We use the notation  $\mathbb{P}_{\{\}}^{Y|X[M]}$  as shorthand to refer to the probability set  $\mathbb{P}_{\{\}}$  maximal with respect to  $\mathbb{P}_{\{\}}^{Y|X}$ .

Lemma 5.12 shows that the semidirect product of any pair of valid conditional probabilities is itself a valid conditional. Suppose we have some collection of  $X_i|X_{[i-1]}$ -valid conditionals  $\{\mathbb{P}_i^{X_i|X_{[i-1]}} | i \in [n]\}$ ; then recursively taking the semidirect product  $\mathbb{M} := \mathbb{P}_1^{X_1} \odot (\mathbb{P}_2^{X_2|X_1} \odot \dots)$  yields a  $X_{[n]}$  valid distribution. Furthermore, the maximal probability set associated with  $\mathbb{M}$  is nonempty.

Collections of recursive conditional probabilities often arise in causal modelling – in particular, they are the foundation of the structural equation modelling approach Richardson and Robins (2013); Pearl (2009).

Note that validity is not a necessary condition for a conditional to define a non-empty probability set. The intuition for this is: if we have some  $\mathbb{K} : X \rightarrow Y$ ,  $\mathbb{K}$  might be an invalid  $Y|X$  conditional on all of  $X$ , but might be valid on some subset of  $X$ , and so we might have some probability model  $\mathbb{P}$  that assigns measure 0 to the bad parts of  $X$  such that  $\mathbb{K}$  is a version of  $\mathbb{P}^{Y|X}$ . On the other hand, if we want to take the product of  $\mathbb{K}$  with arbitrary valid  $X$  probabilities, then the validity of  $\mathbb{K}$  is necessary (Theorem 5.14).

### 2.7.1 Conditional independence

Conditional independence has a familiar definition in probability models. We define conditional independence with respect to a probability gap model to be equivalent to conditional independence with respect to every base measure in the range of the model. This definition is closely related to the idea of *extended conditional independence* proposed by Constantinou and Dawid (2017), see Appendix 5.4.

**Definition 2.25** (Conditional independence with respect to a probability model). For a *probability model*  $\mathbb{P}_\alpha$  and variables  $A, B, Z$ , we say  $B$  is conditionally independent of  $A$  given  $C$ , written  $B \perp\!\!\!\perp_{\mathbb{P}_\alpha} A|C$ , if

$$\mathbb{P}_\alpha^{ABC} = \begin{array}{c} \triangleleft \mathbb{P}_\alpha^C \quad \begin{array}{l} \boxed{\mathbb{P}_\alpha^{A|C}} \text{---} A \\ \boxed{\mathbb{P}_\alpha^{B|C}} \text{---} B \\ \text{---} C \end{array} \end{array} \quad (32)$$

Cho and Jacobs (2019) have shown that this definition coincides with the standard notion of conditional independence for a particular probability model (Theorem 5.15).

Conditional independence satisfies the *semi-graphoid axioms*. For all standard measurable spaces  $(\Omega, \mathcal{F})$  and all probability measures  $\mathbb{P} \in \Delta(\Omega)$ :

1. Symmetry:  $A \perp\!\!\!\perp_{\mathbb{P}} B|C$  iff  $B \perp\!\!\!\perp_{\mathbb{P}} A|C$
2. Decomposition:  $A \perp\!\!\!\perp_{\mathbb{P}} (B, C)|W$  implies  $A \perp\!\!\!\perp_{\mathbb{P}} B|W$  and  $A \perp\!\!\!\perp_{\mathbb{P}_\square} C|W$
3. Weak union:  $A \perp\!\!\!\perp_{\mathbb{P}} (B, C)|W$  implies  $A \perp\!\!\!\perp_{\mathbb{P}} B|(C, W)$
4. Contraction:  $A \perp\!\!\!\perp_{\mathbb{P}} C|W$  and  $A \perp\!\!\!\perp_{\mathbb{P}} B|(C, W)$  implies  $A \perp\!\!\!\perp_{\mathbb{P}_\square} (B, C)|W$

We define *universal conditional independence* with respect to a probability set as conditional independence for every probability model in the set.

**Definition 2.26** (Universal conditional independence). For a *probability set*  $\mathbb{P}_\square$  and variables  $A, B, Z$ , we say  $B$  is universally conditionally independent of  $A$  given  $C$ , written  $B \perp\!\!\!\perp_{\mathbb{P}_\square} A|C$ , if for all  $\mathbb{P}_\alpha \in \mathbb{P}_\square$   $A \perp\!\!\!\perp_{\mathbb{P}_\alpha} C|B$ .

It is straightforward to show that universal conditional independence satisfies the semi-graphoid axioms.

**Lemma 2.27.**  $[\forall x : (f(x) \implies g(x))] \implies [(\forall x : f(x)) \implies (\forall x : g(x))]$

*Proof.*

$$\begin{array}{ll}
\forall x : f(x) \implies g(x) & \text{premise} \\
(33) & \\
\forall x : f(x) & \text{premise} \\
(34) & \\
f(a) & \text{universal instantiation on 34 substitute } a/x \\
(35) & \\
f(a) \implies g(a) & \text{universal instantiation on 33 substitute } a/x \\
(36) & \\
g(a) & \text{modus ponens 35 and 36} \\
(37) & \\
\forall x : g(x) & \text{universal generalisation on 37} \\
(38) & \\
(\forall x : f(x)) \implies (\forall x : g(x)) & \text{conditional proof 34 – 38} \\
(39) & \\
[\forall x : (f(x) \implies g(x))] \implies [(\forall x : f(x)) \implies (\forall x : g(x))] & \text{conditional proof 33–39} \\
(40) &
\end{array}$$

With thanks to (<https://math.stackexchange.com/users/252356/kylew>) for the proof.  $\square$

**Lemma 2.28.** *Given a standard measurable space  $(\Omega, \mathcal{F})$  and  $\mathbb{P}_{\{\}} on  $\Omega$ , universal conditional independence with respect to  $\mathbb{P}_{\{\}}$  satisfies the semi-graphoid axioms.$*

*Proof.* For a particular probability  $\mathbb{P}_\alpha$ , each of the semi-graphoid axioms consists of a statement of the form  $\forall \mathbb{P} : f(\mathbb{P}) \implies g(\mathbb{P})$  (in the case of the first axiom, it corresponds to two such statements).

As the axioms hold for conditional independence for any probability model, we have, for arbitrary  $\mathbb{P}_{\{\}}, \forall \mathbb{P}_\alpha \in \mathbb{P}_{\{\}} : f(\mathbb{P}_\alpha) \implies g(\mathbb{P}_\alpha)$ .

Then, by Lemma 2.27,  $(\forall \mathbb{P}_\alpha \in \mathbb{P}_{\{\}} : f(\mathbb{P}_\alpha)) \implies (\forall \mathbb{P}_\alpha \in \mathbb{P}_{\{\}} : g(\mathbb{P}_\alpha))$ .

Note that  $(\forall \mathbb{P}_\alpha \in \mathbb{P}_{\{\}} : f(\mathbb{P}_\alpha))$  is, by definition, a universal conditional independence statement with respect to  $\mathbb{P}_{\{\}}$ .  $\square$

### 3 Decision problems

We want to construct models to help make decisions. For our purposes, “making a decision” means we have some mathematically well-defined set  $C$  of choices under consideration, and some means of comparing one choice to another that

induces a partial order on choices. If we are trying to be precise with language, we might call this *formal* decision making; people may often make decisions where it is difficult to identify a set  $C$  of choices under consideration or a rule for inducing a partial order on it.

A procedure for making formal decisions is, in a sense, the opposite of a measurement procedure. With measurement, we have some unclear process that interacts with the world and leaves us with a collection of mathematical objects. Formal decision making starts with a collection of mathematical objects – a set of choices, a partial order and a tie-breaking rule which together imply a *choice*, and then on the basis of this choice some unclear procedure takes place that has consequences in the world. If a measurement is a “function” whose domain is the world, a formal decision can be thought of as a “function” whose codomain is the world.

We make the following assumptions about how choices are compared:

- Each choice is associated with a probability set over some sample space  $(\Omega, \mathcal{F})$ ; that is, we have a function  $f : C \rightarrow \mathcal{P}(\Delta(\Omega))$
- There is some method available for comparing the desirability of  $\mathbb{P}_\alpha$  and  $\mathbb{P}_{\alpha'}$  for all  $\alpha, \alpha' \in C$ . For example, we might have a utility function  $u : \Omega \rightarrow \mathbb{R}$ , and  $|\mathbb{P}_\alpha| = 1$  for all  $\alpha \in C$ ; then we can compare  $\mathbb{P}_\alpha$  with  $\mathbb{P}_{\alpha'}$  using expected utility

We can represent such a model  $f$  with probability sets. There are many different ways to do this, but the general scheme is as follows:

- There is a probability set  $\mathbb{P}_\square$  representing “properties that hold regardless of which choice is taken”
- There is a collection of probability sets  $\{\mathbb{P}_{\tilde{\alpha}}\}_A$  representing “properties that hold just for the choice  $\alpha$ ”
- $\mathbb{P}_\alpha = \mathbb{P}_\square \cap \mathbb{P}_{\tilde{\alpha}}$

It is always possible to accomplish this: take  $\mathbb{P}_\square \supset \cup_{\alpha \in C} \mathbb{P}_\alpha$  and  $\mathbb{P}_{\tilde{\alpha}} = \mathbb{P}_\alpha$ . However, we might be motivated to make different choices for  $\mathbb{P}_\square$  and the  $\mathbb{P}_{\tilde{\alpha}}$ s.

The probability set representation allows us to make use of universal conditional probabilities and universal conditional independence to reason about decision problem models. By construction  $\mathbb{P}_\alpha \subset \mathbb{P}_\square$  for all  $\alpha$ , any universal conditional independence that holds for  $\mathbb{P}_\square$  holds for all  $\mathbb{P}_\alpha$  as well, and similarly any universal conditional probability with respect to  $\mathbb{P}_\square$  is also a universal conditional probability for all  $\mathbb{P}_\alpha$ .

We call this general scheme a *probability gap model*. A probability gap model specifies some universal behaviour via  $\mathbb{P}_\square$ , but this specification is incomplete; it has some gaps. We then have a selection of probability sets  $\{\mathbb{P}_{\tilde{\alpha}}\}_A$  that specify choice-specific behaviour; this set represents different ways to fill the gap.

**Definition 3.1** (Probability gap model). Given  $(\Omega, \mathcal{F})$ , a probability gap model is a triple  $(\mathbb{P}_\square, \{\mathbb{P}_{\tilde{\alpha}}\}_A, f)$  where  $\mathbb{P}_\square$  is a probability set and  $\{\mathbb{P}_{\tilde{\alpha}}\}_A$  is a collection of probability sets and  $f : A \rightarrow \mathcal{P}(\Delta(\Omega))$  is the map

$$\mathbb{P}_\alpha := f(\alpha) \quad (41)$$

$$= \mathbb{P}_\square \cap \mathbb{P}_{\tilde{\alpha}} \quad (42)$$

### 3.1 Conditional probability models

A simple but interesting class of probability gap model is the *conditional probability model*. Conditional probability models can be thought of as describing problems in which there is a variable  $X$  whose marginal probability can be chosen and a variable  $Y$  that responds in a fixed way to  $X$ .

**Definition 3.2** (Conditional probability model). Given  $(\Omega, \mathcal{F})$ ,  $Y : \Omega \rightarrow Y$ ,  $X : \Omega \rightarrow X$ , a  $Y|X$  conditional probability model is a probability gap model  $(\mathbb{P}_\square^{Y|X[M]}, \{\mathbb{P}_\alpha^{X[M]}\}_A, f)$ .

Conditional probability models arise when we have variables that represent the choice made, and each choice is associated with a *unique* probability distribution over an outcome of interest.

**Example 3.3** (Choice variable). Suppose we have a procedure  $\mathcal{C}$  that compares the elements of a countable set  $C$  and produces a choice according to some notion of the desirability of each one. Another procedure  $\mathcal{Y}$  measures outcomes of interest. These are modelled with variables  $C : \Omega \rightarrow C$  and  $Y : \Omega \rightarrow Y$  respectively. Without loss of generality, we take  $\Omega = C \times Y$  and  $C$  and  $Y$  are projections onto the matching sets.

For each  $\alpha \in C$ , suppose that  $\mathbb{P}_\alpha^C$  exists and is equal to  $\delta_\alpha$ . This expresses the notion that, if the choice procedure actually yields  $\alpha$ , then the model should always assign probability 1 to  $C \bowtie \alpha$ .

Suppose also that  $\mathbb{P}_\alpha^Y$  exists for all  $\alpha$ . That is, the marginal probability of  $Y$  is unique given any choice  $\alpha$ .

We thus have a model  $f : C \rightarrow \Delta(\Omega)$  given by  $\alpha \mapsto \mathbb{P}_\alpha^{CY}$  where

$$\mathbb{P}_\alpha^{CY}(B \times D) = \delta_\alpha(B) \mathbb{P}_\alpha^Y(D) \quad (43)$$

Validity of  $\mathbb{P}_\alpha^{CY}$  is guaranteed by the definitions of  $C$  and  $Y$  because  $(C, Y)$  is surjective.

We can implement  $f$  with a conditional probability model  $(\mathbb{P}_\square^{Y|C[M]}, \{\mathbb{P}_\alpha^{C[M]}\}_A, f)$  where  $\mathbb{P}_\square^{Y|C} := (D|\alpha) \mapsto \mathbb{P}_\alpha^Y(D)$ .

Then

$$\mathbb{P}_\alpha^{CY} = \mathbb{P}_\alpha^C \odot \mathbb{P}_\square^{Y|C} \quad (44)$$



and so by Lemma 5.10

$$f(\alpha) = \mathbb{P}_\alpha \cap \mathbb{P}_\square \quad (45)$$

find a home for the following remarks

If the conditional probability  $\mathbb{P}_{\{\}}^{Y|X}$  and all the marginal probabilities  $\mathbb{P}_\alpha^X$  are valid, then by Lemma 5.12  $\mathbb{P}_{\{\}} \cap \alpha \neq \emptyset$  for all  $\alpha \in A$ . Thus validity of all the individual parts is enough to ensure compatibility.

We can define more complex probability gap models with a similar approach where, for example, the model is specified by an incomplete collection of conditional probabilities and the choices are each a complementary collection of conditional probabilities; we call such models *probability comb models* after Chiribella et al. (2008); Jacobs et al. (2019), but we will not address them in this paper.

### 3.2 Example: invalidity

Body mass index is defined as a person's weight divided by the square of their height. Suppose we have a measurement process  $\mathcal{S} = (\mathcal{W}, \mathcal{H})$  and  $\mathcal{B} = \frac{\mathcal{W}}{\mathcal{H}^2}$  - i.e. we figure out someone's body mass index first by measuring both their height and weight, and then passing the result through a function that divides the second by the square of the first. Thus, given the random variables  $\mathcal{W}, \mathcal{H}$  modelling  $\mathcal{W}, \mathcal{H}$ ,  $\mathcal{B}$  is the function given by  $\mathcal{B} = \frac{\mathcal{W}}{\mathcal{H}^2}$ . Given  $x \in \mathbb{R}$ , consider the conditional probability

$$\nu^{\mathcal{B}|\mathcal{WH}} = \begin{array}{c} \mathcal{H} \text{ --- } * \\ \mathcal{W} \text{ --- } * \end{array} \triangleleft_{\delta_x} \text{ --- } \mathcal{B} \quad (46)$$

Then pick some  $w, h \in \mathbb{R}$  such that  $\frac{w}{h^2} \neq x$  and  $(\mathcal{W}, \mathcal{H}) \bowtie (w, h) \neq \emptyset$  (our measurement procedure could possibly yield  $(w, h)$  for a person's height and weight). We have  $\nu^{\mathcal{B}|\mathcal{WH}}(x|w, h) = 1$ , but

$$(\mathcal{B}, \mathcal{W}, \mathcal{H}) \bowtie \{(x, w, h)\} = \{\omega | (\mathcal{W}, \mathcal{H})(\omega) = (w, h), \mathcal{B}(\omega) = \frac{w}{h^2}\} \quad (47)$$

$$= \emptyset \quad (48)$$

so  $\nu^{\mathcal{B}|\mathcal{WH}}$  is invalid, and there is some valid  $\mu^X$  such that the probability set  $\mathbb{P}_{\{\}}$  with  $\mathbb{P}_{\{\}}^{XY} = \mu^X \odot \nu^{Y|X}$  is empty.

Validity rules out conditional probabilities like 46. We guess that in many cases this condition may either be trivial or unconsciously taken into account when constructing conditional probabilities. However, if we are not cognizant of the conditional our model depends on, we may inadvertently propose a model that depends on invalid conditional probabilities. For example, the conditional probability 46 would be used to evaluate the causal effect of body mass index in the causal diagram found in Shahar (2009), presuming the author used the

term “causal effect” to depend somehow on the function  $x \mapsto P(\cdot | do(B = x))$  as is the usual convention when discussing causal Bayesian networks.

End find a home for the following remarks

### 3.3 Response conditionals

Given any probability gap model  $(\mathbb{P}_\square, \{\mathbb{P}_\alpha\}_A, f)$  and variables  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ , if we have a conditional probability  $\mathbb{P}_\square^{Y|X}$ , then this must be the conditional probability of  $Y$  given  $X$  for any  $\alpha \in A$ . It must therefore be the case that

$$\mathbb{P}_\alpha^Y = \mathbb{P}_\alpha^X \mathbb{P}_\square^{Y|X} \quad \forall \alpha \in A \quad (49)$$

If it is possible to control  $X$  to some extent (i.e.  $\mathbb{P}_\alpha^X \neq \mathbb{P}_{\alpha'}^X$  for some  $\alpha, \alpha'$ ), then  $\mathbb{P}_\square^{Y|X}$  tells us how  $\mathbb{P}_\alpha^X$  determines its effect on  $\mathbb{P}_\alpha^Y$ . This feature motivates the name *response conditional* for conditional probabilities of this type.

The motivating question we introduced at the beginning of this paper was “when are potential outcomes well-defined?”. This is not written using the potential outcomes framework, so we cannot directly address this question. However, we can ask “when do probability gap models feature response conditionals?”.

### 3.4 Response conditionals and potential outcomes

There is a connection between the question of when a particular conditional probability exists and when potential outcomes are well-defined. Given any Markov kernel there is an operation akin to function currying termed “randomness pushback” that represents the kernel as the product of a probability measure and a deterministic Markov kernel. We observe that potential outcomes models share a number of features in common with a Markov kernel represented using a randomness pushback.

Unlike function currying, there are many different randomness pushbacks that represent the same Markov kernel. The interpretation of potential outcomes models seems to require that exactly one of these possible pushbacks is a genuine potential outcomes model. Several works including Dawid (2000) and Richardson and Robins (2013) have expressed the view that some of the degrees of freedom in potential outcomes models are superfluous; both propose that the degrees of freedom that can be tested using some idealised experiment are the degrees that should be kept.

Notably, the single world intervention graphs of Richardson and Robins (2013) feature an operation that splits an “intervenable” variable  $X$  into two versions,  $X$  and  $X'$ , representing “the actual  $X$ ” and “the unobserved value  $X$  would have taken absent intervention” respectively (this interpretation is ours, not theirs). Thus Richardson and Robins might argue that they are interested in the existence of response conditionals of the form  $\mathbb{P}_\square^{Y|XX'}$  rather than  $\mathbb{P}_\square^{Y|X}$ .

We do not take a position on which degrees of freedom are good or bad in a typical potential outcomes model, nor do we explore conditionals of the form  $\mathbb{P}^{Y|XX'}$ . We can, however, distinguish two different questions:

- How are response conditionals represented?
- Which response conditionals are of interest?

The first question seems to be an inconsequential stylistic matter, while the second question may determine the direction of one's analysis.

### 3.5 Randomness pushbacks

Given a function  $f : X \times Y \rightarrow Z$ , we can obtain a curried version  $\lambda f : Y \rightarrow Z^X$ . In particular, if  $Y = \{*\}$  then  $\lambda f : \{*\} \rightarrow Y^X$ . At least for countable  $X$ , we can apply this construction to Markov kernels: given a kernel  $\mathbb{K} : X \rightarrow Y$ , define  $\mathbb{L} : \{*\} \rightarrow Y^X$  by

$$\lambda \mathbb{K}((y_i)_{i \in X}) = \prod_{i \in X} \mathbb{K}(y_i|i) \quad (50)$$

We can then define an evaluation map  $\text{ev} : Y^X \times X \rightarrow Y$  by  $\text{ev}((y_i)_{i \in X}, x) = y_x$ . Then

$$\mathbb{K} = \begin{array}{c} \triangleleft \mathbb{L} \\ \text{X} \end{array} \xrightarrow{\quad} \boxed{\mathbb{F}_{\text{ev}}} \longrightarrow \text{Y} \quad (51)$$

$$\iff \quad (52)$$

$$\mathbb{K}(A|x) = \int_{Y^X} \delta_{\text{ev}(y^X, x)}(A) \mathbb{L}(dy^X|x) \quad (53)$$

Unlike the case of function currying,  $\lambda \mathbb{K}$  is not the unique Markov kernel for which 51 holds. In fact, we can substitute any  $\mathbb{M}$  such that, for any  $i \in X$

$$\sum_{y_{\{i\}^C} \in Y^{|X|-1}} \mathbb{M}((y_i)_{i \in X}) = \mathbb{K}(y_i|i) \quad (54)$$

This representation of a Markov kernel is called a *randomness pushback* by Fritz (2020).

Randomness pushbacks have a few features in common with potential outcomes causal models. For our purposes, we will say a potential outcomes model is a probability set  $(\Omega, \mathcal{F}, \mathbb{P}_{\{\cdot\}})$  along with variables  $\text{X}, \text{Y}, \text{Y}^X$  such that

$$\mathbb{P}_{\{\cdot\}}^{Y|XY^X} = \mathbb{F}_{\text{ev}} \quad (55)$$

More commonly, this property is expressed as

$$Y \stackrel{a.s.}{=} \text{ev}(X, Y^X) \quad (56)$$

We consider a potential outcomes model to be a probability set here, but we can formally recover a “traditional” potential outcomes model by considering probability sets of size 1.

If we additionally have the existence of  $\mathbb{P}_{\{\}}^{Y^X|X}$  and  $Y^X \perp\!\!\!\perp_{\mathbb{P}_{\{\}}} X$  then

$$\mathbb{P}_{\{\}}^{Y|X} = \begin{array}{c} \triangle \mathbb{P}_{\{\}}^{Y^X} \\ \swarrow \searrow \\ X \quad \text{F}_{\text{ev}} \end{array} \longrightarrow Y \quad (57)$$

Equation 57 is clearly a version of 51. As we have established, provided  $\mathbb{P}_{\{\}}^{Y|X}$  exists, we can always introduce some variable  $Y^X$  and corresponding  $\mathbb{P}_{\{\}}^{Y^X}$  such that Equation 57 holds.

### 3.5.1 Choices aren’t always known

One area of potential difficulty with our approach to formalising causal inference from the starting point of modelling decision problems is related to the issue of unknown choice sets. While causal investigations are often concerned with helping someone to make better decisions, the kind of “decision making process” associated with them is not necessarily well modeled by the setup above. Often the identity of the decision maker and the exact choices at hand are vague. Consider Banerjee et al. (2016): a large scale experiment was conducted trialling a number of different strategies all aiming to increase the amount of learning level appropriate instruction available to students in four Indian states. It is not clear who, exactly, is going to make a decision on the basis of this information, but one can guess:

- They’re someone with interest in and authority to make large scale changes to a school system
- They consider the evidence of effectiveness of teaching at the right level relevant to their situation
- They consider the evidence regarding which strategies work to implement this approach relevant to their situation

This could describe a writer who is considering what kind of advice they can provide in a document, a grantmaker looking to direct funds, a policy maker trying to design policies with appropriate incentives a program manager trying to implement reforms or someone in a position we haven’t thought of yet. All of these people have very different choices facing them, and to some extent it is desirable that this research is relevant to all of them.

These situations are common in the field of causal inference and the question of how to formalise them may be worthy of study. The situation of known choices that we focus on here is easier to study, and it may serve as an idealisation or a limiting case of situations where choices are unknown.

### 3.6 Other decision theoretic causal models

Lattimore and Rohde (2019a) and Lattimore and Rohde (2019b) consider an observational probability model coupled to a collection of indexed interventional probability models, with the observational probability model coupled to the interventional models by shared (unobserved) parameters. In these papers, they show how such a model can reproduce inferences made using Causal Bayesian Networks. In our terms, Lattimore and Rohde’s model family can be thought of as conditional probability models  $(\mathbb{P}_{\square}^{\text{OCH}|\mathcal{D}[M]}, \{\mathbb{P}_{\alpha}^{\mathcal{D}}\}_A, f)$  where  $\mathcal{D}$  represents the choice of intervention,  $\text{O}$  observations,  $\text{I}$  interventional consequences and  $\mathbb{H}$  model parameters or hypotheses, with the properties  $\text{H} \perp\!\!\!\perp_{\mathbb{P}_{\square}} \mathcal{D}$ ,  $\text{O} \perp\!\!\!\perp_{\mathbb{P}_{\square}} \mathcal{D}|\text{H}$  and  $\text{C} \perp\!\!\!\perp_{\mathbb{P}_{\square}} \text{O}|\mathcal{D}, \text{H}$ .

Note that we may well prefer to use a model in which the intervention  $\mathcal{X}$  is chosen to depend on the observations  $\text{O}$  somehow. It’s possible to represent this using probability gap models, but doing so is beyond the scope of this paper.

The approach to decision theoretic causal inference described by Dawid (2020) is somewhat different to Rohde and Lattimore’s:

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  $\text{T}$ , and a response variable  $\text{Y}$ . We consider three different regimes [...] the first two regimes may be described as interventional, and the last as observational.

The key difference is that Rohde and Lattimore’s model employs different variables  $\text{O}$  and  $\text{C}$  to represent observations and interventional consequences, which allows us to write conditional independence statements like  $\text{C} \perp\!\!\!\perp_{\mathbb{P}_{\{\}}}\text{O}|\mathcal{D}, \text{H}$ , while Dawid’s approach uses the same variable  $\text{Y}$  to represent observations and interventional consequences, depending on the choice of regime. In this scheme there is no way we can see to express something like “the observations are independent of the interventional consequences given the choice of intervention and the parameters”.

If we require that the map from regimes to probability distributions is measurable (which can be trivially satisfied if the set of regimes is countable), then we can also characterise Dawid’s approach using a probability set  $\mathbb{P}_{\{\}}^{\text{TY}|\text{F}_T[M]}$  where  $\text{F}_T$  represents the prevailing regime, and  $\text{T}$  and  $\text{Y}$  are as in the quote above. Note that we do not consider conditional probability models here, because it is not obvious how we should say  $\text{F}_T$  is distributed given any choice of intervention; for any choice of intervention,  $\text{F}_T$  will sometimes take the value

“observational” and sometimes take a value corresponding to the chosen intervention. Moreover, it seems inappropriate to posit a sequence of independent and identically distributed copies of  $F_T$  because we are likely to know it advance exactly when it will and won’t take on the observational value, and so we can’t associate it with a unique marginal distribution.

Heckerman and Shachter (1995) also explore a decision theoretic approach to causal inference. Their approach is based on the decision theory of Savage (1954), and represents a decision problem in terms of choices  $D$ , outcome variables  $U$  and unobserved states of the world  $S$ ; each state of the world defines a deterministic map  $D \rightarrow U$ . In comparison with the theory Lattimore and Rohde Heckerman and Schachter’s comes with the added requirement of deterministic outcomes. It also features no built in distinction between observations and outcomes of interventions, though one could always consider models where  $U$  is in fact a pair of variables  $(O, I)$  such that  $O$  satisfies the independences required of observations in the Lattimore and Rohde model.

## 4 When do response conditionals exist?

The specific question we ask here is: given a conditional probability model  $(\mathbb{P}_{\square}^{Y|D^{[M]}}, \{\mathbb{P}_{\alpha}^D\}_A, f)$  where  $D = (D_i)_{i \in M}$  are choice variables,  $Y = (Y_i)_{i \in M}$  are outcome variables and  $M$  is some index set, when does there exist a response conditional  $\mathbb{P}_{\square}^{Y_0|D_0}$  that explains how each  $Y_i$  responds to each  $D_i$ ; we call these *repeatable response conditionals*. The reason why we consider a sequential problem is that in practice, causal inference almost always deals with observational data that is assumed to be appropriately modeled with a sequence of independent and identically distributed random variables. Furthermore, we can characterise precisely the kinds of models for which such response conditionals exist. Thus the sequential setting is both a theoretically tractable and widely applicable starting point.

We examine this question from two points of view: firstly, we ask *what kind of conditional probability models exhibit repeatable response conditionals?* Secondly, we ask *what kinds of experiments are these conditional probability models appropriate for?* In the course of answering the second question, we show that apparently subtle differences in the description of an experimental procedure can determine whether a particular experiment should or should not be modeled with repeatable response conditionals.

We need to make strong assumptions about an experiment to establish the existence of repeatable response conditional in the first place, once we have repeatable response conditionals with respect to some pair of variables  $(Y, D)$ , response conditionals with respect to different pairs of variables may exist due to conditional independences in the original  $\mathbb{P}_{\square}^{Y_0|D_0}$ . These conditional independences can be tested for in the observed data.

## 4.1 Repeatable experiments

Our setup is a conditional probability model  $(\mathbb{P}_{\square}^{\overline{Y|D}}, A)$  where  $Y := Y_M = (Y_i)_{i \in M}$  and  $D := D_M = (D_i)_{i \in M}$  for some index set  $M$ ; we say such a model is a model of a *sequential experiment*. We say that  $Y_i$  is the consequence corresponding to the decision  $D_i$  for all  $i \in M$  (i.e. variables with matching indices correspond). We say that  $\mathbb{P}_{\square}^{\overline{Y|D}}$  features *repeatable response conditionals* if there exists a hypothesis  $H$  such that  $\mathbb{P}_{\square}^{Y_i|HD_i} = \mathbb{P}_{\square}^{Y_j|HD_j}$  for all  $i, j \in M$ ,  $H \perp\!\!\!\perp_{\mathbb{P}_{\square}} D$  and  $Y_i \perp\!\!\!\perp_{\mathbb{P}_{\square}} Y_{M \setminus \{i\}} D_{M \setminus \{i\}} | HD_i$ . We remind the reader that in a conditional probability model, arbitrary conditional probabilities do not always exist, see Definition 2.14.

There are two assumptions relevant to the existence of repeatable response conditionals. The first is a condition of interchangeability: in particular, given a permutation of  $M$ , we get the same result from applying this permutation only to the set of actions we take, or from applying it only to the set of outcomes we observe. We call this *exchange commutativity*.

The second is a condition of *locality of consequences*; that is the assumption that  $Y_i$  is independent of  $D_j$  given  $D_i$  for any  $j$ . It is possible to have models in which commutativity to exchange holds but locality of consequences does not. Such a situation could arise in a model of stimulus payments to individuals in a nation; if exactly  $n$  payments of \$10 000 are made, we might consider that it doesn't matter much exactly who receives the payments (this is a subtle question, though, we will return to it in more detail later). However, the amount of inflation induced depends on the number of payments; making 100 such payments will have a negligible effect on inflation, while making payments to everyone in the country is likely to have a substantial effect. Dawid (2000) discusses condition of *post-treatment exchangeability* which is similar to exchange commutativity, and there he gives the example of herd immunity in vaccination campaigns as a situation where post-treatment exchangeability holds but locality of consequences does not.

As we have mentioned, exchange commutativity is similar to the condition of *post-treatment exchangeability* found in Dawid (2020). Exchange commutativity is also very similar to the exchangeability assumption of GREENLAND and ROBINS (1986), and the assumption of exchangeability found in Banerjee et al. (2017). However, in every case there is a subtle but important difference; exchange commutativity concerns exchanges of actions and outcomes, while these other exchangeability conditions concern exchange of “people” or “experimental units”. Swapping people and experimental units are actions in the real world, and so these symmetries have to be described as part of the measurement procedure, and can't be purely characterised as symmetries of a probabilistic model. On the other hand, swapping the orders of variables *can* be described purely as a symmetry of a probabilistic model, as these swaps involve only function composition. As we will discuss in detail, exchangeability of experimental units does not always imply exchange commutativity.

Locality of consequences is similar to the stable unit treatment distribution

assumption (SUTDA) in Dawid (2020). It is also related to the “no interference” part of the stable unit treatment value assumption (SUTVA). The stable unit treatment value assumption (SUTVA) is given as (Rubin, 2005):

“(SUTVA) comprises two subassumptions. First, it assumes that *there is no interference between units* (Cox 1958); that is, neither  $Y_i(1)$  nor  $Y_i(0)$  is affected by what action any other unit received. Second, it assumes that *there are no hidden versions of treatments*; no matter how unit  $i$  received treatment 1, the outcome that would be observed would be  $Y_i(1)$  and similarly for treatment 0.

Not sure if or where I want to put this, I just think it helps to illustrate the difference

Exchange commutativity is not equivalent to exchangeability in the sense of De Finetti’s well-known theorem de Finetti ([1937] 1992). The latter can be understood as expressing an indifference between conducting the experiment as normal, or conducting the experiment and then swapping some labels. However, swapping *choices* will (usually) lead to different “pieces of the experiment” receiving different treatment, which is something that can’t be achieved by swapping labels after the experiment has concluded.

The difference is illustrated by the following pair of diagrams.

Exchangeability (swapping labels):

(58)

Exchange commutativity (swapping choices  $\sim$  swapping labels):

(59)

—end not sure where to put—

## 4.2 Consequence contractibility

We offer formal definitions of exchange commutativity and locality of consequences, as well as “consequence contractibility”, which is the conjunction of both conditions.

A conditional probability model commutes with exchange if applying a permutation to the choice  $D_M$  “before” it is taken yields the same result as applying the corresponding permutation to  $Y_M$  “after” it is observed.

**Definition 4.1** (Swap map). Given  $M \subset \mathbb{N}$  a finite permutation  $\rho : M \rightarrow M$  and a variable  $X : \Omega \rightarrow X^M$  such that  $X = (X_i)_{i \in M}$ , define the Markov kernel  $\text{swap}_{\rho(X)} : X^M \rightarrow X^M$  by  $(d_i)_{i \in \mathbb{N}} \mapsto \delta_{(d_{\rho(i)})_{i \in \mathbb{N}}}$ .



**Definition 4.2** (Exchange commutativity). Suppose we have a sample space  $(\Omega, \mathcal{F})$  and a conditional probability model  $(\mathbb{P}_{\square}^{Y|D[M]}, \{\mathbb{P}_{\alpha}^D\}_A, f)$  with  $Y := Y_M := (Y_i)_M$ ,  $D := D_M := (D_i)_M$ ,  $M \subseteq \mathbb{N}$ . If, for any decision rule  $\alpha \in A$ ,

$$\mathbb{P}_{\alpha}^D \odot \text{swap}_{\rho(D)} \mathbb{P}_{\square}^{Y|D} = \mathbb{P}_{\alpha}^D \odot \mathbb{P}_{\square}^{Y|D} \text{swap}_{\rho(Y)} \quad (60)$$

Then  $\mathbb{P}_{\square}$  *commutes with exchange*.

A conditional probability model exhibits locality of consequences if, given two different choices that agree on an subsequence of indices, the model yields identical outcomes if we restrict our attention to the subsequence on which the different choices match. For example, if we have  $D = (D_1, D_2, D_3)$  and  $Y = (Y_1, Y_2, Y_3)$  and  $\mathbb{P}_{\alpha}^{D_1 D_3} = \mathbb{P}_{\beta}^{D_1 D_3}$  then  $\mathbb{P}_{\alpha}^{Y_1 D_1 Y_3 D_3} = \mathbb{P}_{\beta}^{Y_1 D_1 Y_3 D_3}$ .

**Definition 4.3** (locality of consequences). Suppose we have a sample space  $(\Omega, \mathcal{F})$  and a conditional probability model  $(\mathbb{P}_{\square}^{Y|D[M]}, \{\mathbb{P}_{\alpha}^D\}_A, f)$  with  $Y := Y_M := (Y_i)_M$ ,  $D := D_M := (D_i)_M$ ,  $M \subseteq \mathbb{N}$ . For any ordered sequence  $S = (s_i)_{i \in Q}$  where  $Q \subset M$  and  $i < j \implies s_i < s_j$ , let  $D_S := (D_i)_{i \in S}$  and  $D_T := (D_i)_{i \in T}$ . If for any  $\alpha, \beta \in R$

$$\mathbb{P}_{\alpha}^{D_S} = \mathbb{P}_{\beta}^{D_S} \quad (61)$$

$$\implies \mathbb{P}_{\alpha}^{(D_i, Y_i)_{i \in S}} = \mathbb{P}_{\beta}^{(D_i, Y_i)_{i \in S}} \quad (62)$$

then  $\mathbb{P}_{\square}$  exhibits *locality of consequences*.

Neither condition implies the other.

**Lemma 4.4.** *Exchange commutativity does not imply locality of consequences or vice versa.*

*Proof.* A conditional probability model that exhibits exchange commutativity but some choices have non-local consequences:

Suppose  $D = Y = \{0, 1\}$  and we have a conditional probability model  $(\mathbb{P}_{\square}^{Y|D[M]}, \{\mathbb{P}_{\alpha}^D\}_A, f)$  where  $D = (D_1, D_2)$ ,  $Y = (Y_1, Y_2)$  and  $A$  contains all deterministic probability measures in  $\Delta(D^2)$ . If

$$\mathbb{P}_{\square}^{Y_1 Y_2 | D_1 D_2}(y_1, y_2 | d_1, d_2) = \llbracket (y_1, y_2) = (d_1 + d_2, d_1 + d_2) \rrbracket \quad (63)$$

Then  $\mathbb{P}_{\delta_{00}}^{Y_1 D_1}(y_1) = \llbracket y_1 = 0 \rrbracket$  while  $\mathbb{P}_{\delta_{01}}^{Y_1} = \llbracket y_1 = 1 \rrbracket$ . However,  $\delta_{00}^{D_1} = \delta_{01}^{D_1} = \delta_0^{D_1}$  so  $\mathbb{P}_{\square}$  exhibits non-local consequences. However, taking  $(d_i, d_j) := \delta_{d_i d_j} \in A$ ,

$$\mathbb{P}_{d_2, d_1}^{Y_1 D_1 Y_2 D_2}(y_1, d_1, y_2, d_2) = \llbracket (y_1, y_2) = (d_2 + d_1, d_2 + d_1) \rrbracket \quad (64)$$

$$= \llbracket (y_2, y_1) = (d_1 + d_2, d_1 + d_2) \rrbracket \quad (65)$$

$$= \mathbb{P}_{d_1, d_2}^{Y_1 D_1 Y_2 D_2}(y_2, d_2, y_1, d_1) \quad (66)$$

so  $\mathbb{P}_\square$  commutes with exchange.

A conditional probability model that exhibits locality of consequences but does not commute with exchange:

Alternatively, suppose the same setup, but define  $\mathbb{P}_\square$  instead by

$$\mathbb{P}_\square Y_1 Y_2 | D_1 D_2 (y_1, y_2 | d_1, d_2) = \llbracket (y_1, y_2) = (0, 1) \rrbracket \quad (67)$$

for all  $\alpha \in A$ .

Then  $\mathbb{P}_\square$  exhibits locality of consequences. If  $\mathbb{P}_\alpha^{D_S} = \mathbb{P}_\beta^{D_S}$  for  $S \subset \{0, 1\}$  then:

$$\mathbb{P}_\alpha^{Y_S D_S}(y_s, d_s) = \sum_{y'_2 \in \{0, 1\}^{S^C}} \llbracket (y_1, y_2) = (0, 1) \rrbracket \mathbb{P}_\alpha^{D_S}(d_s) \quad (68)$$

$$= \mathbb{P}_\beta^{Y_S D_S}(y_s, d_s) \quad (69)$$

However,  $\mathbb{P}_\square$  does not commute with exchange. For all  $\alpha, \beta \in A$ :

$$\mathbb{P}_\alpha Y_1 Y_2(y_1, y_2) = \llbracket (y_1, y_2) = (0, 1) \rrbracket \quad (70)$$

$$\neq \mathbb{P}_\beta Y_1 Y_2(y_2, y_1) \quad (71)$$

□

Although locality of consequences has a lot in common with an assumption non-interference, it still allows for some models in which exhibit certain kinds of interference between actions and outcomes of different indices. For example: I have an experiment where I first flip a coin and record the results of this flip as the outcome of the first step of the experiment, but I can choose either to record this same outcome as the provisional result of the second step (this is the choice  $D_1 = 0$ ), or choose to flip a second coin and record the result of that as the provisional result of the second step of the experiment (this is the choice  $D_1 = 1$ ). At the second step, I may further choose to copy the provisional results ( $D_2 = 0$ ) or invert them ( $D_2 = 1$ ). Then

- The marginal distribution of both experiments in isolation is Bernoulli(0.5) no matter what choices I make, so a model of this experiment would satisfies Definition 4.3
- Nevertheless, the choice for the first experiment affects the result of the second experiment

Note that this example would not satisfy exchange commutativity.

We call the conjunction of exchange commutativity and consequence localilty *causal contractibility*.

**Definition 4.5** (Causal contractibility). A conditional probability model  $(\mathbb{P}_\square^{\overline{Y|D}}, A)$  is causally contractible if it is both commutative with exchange and commutative with marginalisation.

### 4.3 Repeatable consequence conditionals exist iff a model is causally contractible

The main result in this section is Theorem 4.7 which shows that a conditional probability model  $\mathbb{P}_\square$  is causally contractible if and only if it can be represented as the product of a distribution over hypotheses  $\mathbb{P}_\square^H$  and a collection of identical conditional probabilities  $\mathbb{P}_\square^{Y_1|D_1H}$ . Note the hypothesis  $H$  that appears in this conditional; it can be given the interpretation of a random variable that expresses the “true but initially unknown”  $Y_1|D_1$  conditional probability.

**Lemma 4.6** (Exchangeable randomness pushback). *A conditional probability model  $(\mathbb{P}_\square^{Y|D}, A)$  such that  $D := (D_i)_{i \in \mathbb{N}}$  and  $Y := (Y_i)_{i \in \mathbb{N}}$ .  $\mathbb{P}_\square$  is causally contractible if and only if there exists a column exchangeable probability distribution  $\mathbb{P}^{Y^D}$  such that*

$$\mathbb{P}_\square^{Y|D} = \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_\square^{Y^D} \text{ with } D \text{ as input and } Y \text{ as output, connected by a box labeled } \mathbb{F}_{ev}. \end{array} \quad (72)$$

$$\iff \quad (73)$$

$$\mathbb{P}_\square^{Y|D}(y|d) = \mathbb{P}^{(Y_{d_i i})_{i \in \mathbb{N}}}(y) \quad (74)$$

Where  $\mathbb{F}_{ev}$  is the Markov kernel associated with the evaluation map

$$ev : D^{\mathbb{N}} \times Y^{D \times \mathbb{N}} \rightarrow Y \quad (75)$$

$$((d_i)_{i \in \mathbb{N}}, (y_{ij})_{i \in D, j \in \mathbb{N}}) \mapsto (y_{d_i i})_{i \in \mathbb{N}} \quad (76)$$

*Proof.* Only if: Choose  $e := (e_i)_{i \in \mathbb{N}}$  such that  $e_{|D|+j}$  is the  $j$ th element of  $D$  for all  $j \in \mathbb{N}$ . Abusing notation, write  $e$  also for the decision function that chooses  $e$  deterministically.

Define

$$\mathbb{P}^{Y^D}((y_{ij})_{D \times \mathbb{N}}) := \mathbb{P}_e^Y((y_{|D|+j})_{j \in \mathbb{N}}) \quad (77)$$

Now consider any  $d := (d_i)_{i \in \mathbb{N}} \in D^{\mathbb{N}}$ . By definition of  $e$ ,  $e_{|D|+i} = d_i$  for any  $i \in \mathbb{N}$ .

$$\mathbb{Q} : D \rightarrow Y \quad (78)$$

$$\mathbb{Q} := \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_\square^{Y^D} \text{ with } D \text{ as input and } Y \text{ as output, connected by a box labeled } \mathbb{F}_{ev}. \end{array} \quad (79)$$

and consider some ordered sequence  $A \subset \mathbb{N}$  and  $B := ((|D|d_i + i))_{i \in A}$ . Note that  $e_B := (e_{|D|d_i + i})_{i \in B} = d_A = (d_i)_{i \in A}$ . Then

$$\sum_{y \in Y^{-1}(y_A)} \mathbb{Q}(y|d) = \sum_{y \in Y^{-1}(y_A)} \mathbb{P}^{(Y_{d_i i}^D)^A}(y) \quad (80)$$

$$= \sum_{y \in Y^{-1}(y_A)} \mathbb{P}_e^{(Y_{|D|d_i+i}^D)^A}(y) \quad (81)$$

$$= \mathbb{P}_e^{Y_B}(y_A) \quad (82)$$

$$= \mathbb{P}_d^{Y_A}(y_A) \quad \text{by causal contractibility} \quad (83)$$

Because this holds for all  $A \subset \mathbb{N}$ , by the Kolmogorov extension theorem

$$\mathbb{Q}(y|d) = \mathbb{P}_d^Y(y) \quad (84)$$

Because  $d$  is the decision function that deterministically chooses  $d$ , for all  $d \in D$

$$\mathbb{Q}(y|d) = \mathbb{P}_d^{Y|D}(y|d) \quad (85)$$

And because  $\mathbb{P}_d^{Y|D}(y|d)$  is unique for all  $d \in D^{\mathbb{N}}$  and  $\mathbb{P}^{Y|D}$  exists by assumption

$$\mathbb{P}^{Y|D} = \mathbb{Q} \quad (86)$$

Next we will show  $\mathbb{P}^{Y^D}$  is contractible. Consider any subsequences  $Y_S^D$  and  $Y_T^D$  of  $Y^D$  with  $|S| = |T|$ . Let  $\rho(S)$  be the “expansion” of the indices  $S$ , i.e.  $\rho(S) = (|D|i+j)_{i \in S, j \in D}$ . Then by construction of  $e$ ,  $e_{\rho(S)} = e_{\rho(T)}$  and therefore

$$\mathbb{P}^{Y_S^D} = \mathbb{P}_e^{Y_{\rho(S)}} \quad (87)$$

$$= \mathbb{P}_e^{Y_{\rho(T)}} \quad \text{by contractibility of } \mathbb{P} \text{ and the equality } e_{\rho(S)} = e_{\rho(T)} \quad (88)$$

$$= \mathbb{P}^{Y_T^D} \quad (89)$$

If: Suppose

$$\mathbb{P}^{Y|D} = \begin{array}{c} \triangle \\ \mathbb{P}_{\square}^{Y^D} \\ \text{D} \end{array} \begin{array}{c} \text{---} \\ \text{---} \end{array} \begin{array}{c} \square \\ \mathbb{F}_{\text{ev}} \end{array} \text{---} Y \quad (90)$$

and consider any two deterministic decision functions  $d, d' \in D^{\mathbb{N}}$  such that some subsequences are equal  $d_S = d'_T$ .

Let  $Y^{d_S} = (Y_{d_i i})_{i \in S}$ .

By definition,

$$\mathbb{P}^{Y_S|D}(y_S|d) = \sum_{y_S^D \in Y^{D|S}} \mathbb{P}^{Y_S^D}(y_S^D) \mathbb{F}_{\text{ev}}(y_S|d, y_S^D) \quad (91)$$

$$= \sum_{y_S^D \in Y^{D|S|T}} \mathbb{P}^{Y_T^D}(y_S^D) \mathbb{F}_{\text{ev}}(y_S|d, y_S^D) \quad \text{by contractibility of } \mathbb{P}^{Y_T^D} \quad (92)$$

$$= \mathbb{P}^{Y_T|D}(y_S|d) \quad (93)$$

□

As we pointed out, there are similarities between tabular distributions like  $\mathbb{P}^{Y^D}$  that appears in Lemma 4.6 and potential outcomes causal models. However, the  $\mathbb{P}^{Y^D}$  that appears in this lemma usually can't be interpreted as a distribution of potential outcomes. For example, consider a series of bets on fair coinflips. Model the consequence  $Y_i$  as uniform on  $\{0, 1\}$  for any decision  $D_i$ , for all  $i$ . Specifically,  $D = Y = \{0, 1\}$  and  $\mathbb{P}_{\alpha^n}(y) = \prod_{i \in [n]} 0.5$  for all  $n$ ,  $y \in Y^n$ ,  $\alpha \in R$ . Then the construction of  $\mathbb{P}^{Y^D}$  following the method in Lemma 4.6 yields  $\mathbb{P}^{Y_i^D}(y_i^D) = \prod_{j \in D} 0.5$  for all  $y_i^D \in Y^D$ . In this model  $Y_i^0$  and  $Y_i^1$  are independent and uniformly distributed. However, if we wanted  $Y_i^0$  to be interpretable as “what would happen if I bet on outcome 0 on turn  $i$ ” and  $Y^1$  to represent “what would happen if I bet on outcome 1 on turn  $i$ ”, then we ought to have  $Y_i^0 = 1 - Y_i^1$ .

Lemma 4.6 also does not establish that causal contractibility is necessary for the existence of a potential outcomes. A counterexample is any potential outcomes model with potential outcomes  $Z^D$  where the distribution  $\mathbb{P}^{Z^D}$  is not column exchangeable. Such a model is not causally contractible.

The tabular distribution  $\mathbb{P}^{Y^D}$  along with the evaluation function  $\mathbb{F}_{\text{ev}}$  is a randomness pushback of the conditional probability  $\mathbb{P}^{Y|D}$ . Because  $\mathbb{P}^{Y^D}$  is a column exchangeable probability distribution we can apply De Finetti's theorem to show  $\mathbb{P}^{Y^D}$  is representable as a product of identical parallel copies of  $\mathbb{P}^{Y^D|H}$  and a common prior  $\mathbb{P}^H$ . This in turn can be used to show that  $\mathbb{P}_{\square}^{Y|D}$  can be represented as a product of identical parallel copies of  $\mathbb{P}_{\square}^{Y_1|D_1H}$  and the same common prior  $\mathbb{P}_{\square}^H$ . This is the main result: the copies of  $\mathbb{P}_{\square}^{Y_1|D_1H}$  are the repeatable response conditionals.

**Theorem 4.7.** *Suppose we have a sample space  $(\Omega, \mathcal{F})$  and a conditional probability model  $(\mathbb{P}_{\square}^{Y|D}, A)$  such that  $D := (D_i)_{i \in \mathbb{N}}$  and  $Y := (Y_i)_{i \in \mathbb{N}}$ .  $\mathbb{P}_{\square}$  is causally contractible if and only if there exists some  $H : \Omega \rightarrow H$  such that  $\mathbb{P}^{Y_i|HD_i}$  exists*

for all  $i \in \mathbb{N}$  and

$$\mathbb{P}^{Y|HD} = \begin{array}{c} \text{H} \\ \text{D} \end{array} \begin{array}{c} \boxed{\begin{array}{c} \boxed{\Pi_{D,i}} \quad \boxed{\mathbb{P}_{\square}^{Y_0|HD_0}} \end{array}} \end{array} \begin{array}{c} \text{Y}_i \\ i \in \mathbb{N} \end{array} \quad (94)$$

$$\iff \quad (95)$$

$$Y_i \perp\!\!\!\perp Y_{\mathbb{N} \setminus i}, D_{\mathbb{N} \setminus i} | HD_i \quad \forall i \in \mathbb{N} \quad (96)$$

$$\wedge \mathbb{P}^{Y_i|HD_i} = \mathbb{P}^{Y_0|HD_0} \quad \forall i \in \mathbb{N} \quad (97)$$

Where  $\Pi_{D,i} : D^{\mathbb{N}} \rightarrow D$  is the  $i$ th projection map.

*Proof.* We make use of Lemma 4.6 to show that we can represent the conditional probability  $\mathbb{P}_{\square}^{Y|D}$  as

$$\mathbb{P}_{\square}^{Y|D} = \begin{array}{c} \triangle \\ \text{D} \end{array} \begin{array}{c} \mathbb{P}_{\square}^{Y^D} \\ \text{F}_{\text{ev}} \end{array} \text{Y} \quad (98)$$

$$(99)$$

As a preliminary, we will show

$$\mathbb{F}_{\text{ev}} = \begin{array}{c} \text{H} \\ \text{D} \end{array} \begin{array}{c} \boxed{\begin{array}{c} \Pi_{Y^D,i} \\ \Pi_{D,i} \end{array}} \quad \boxed{\mathbb{F}_{\text{ev},i}} \end{array} \begin{array}{c} \text{Y}_i \\ i \in \mathbb{N} \end{array} \quad (100)$$

Where  $\Pi_{Y^D,i} : Y^{D \times \mathbb{N}} \rightarrow Y^D$  is the  $i$ th column projection map on  $Y^{D \times \mathbb{N}}$  and  $\text{ev}_{Y^D \times D} : Y^D \times D \rightarrow Y$  is the evaluation function

$$((y_i)_{i \in D}, d) \mapsto y_d \quad (101)$$

Recall that  $\text{ev}$  is the function

$$((d_i)_{i \in \mathbb{N}}, (y_{ij})_{i \in D, j \in \mathbb{N}}) \mapsto (y_{d_i i})_{i \in \mathbb{N}} \quad (102)$$

By definition, for any  $\{A_i \in \mathcal{Y} | i \in \mathbb{N}\}$

$$\mathbb{F}_{\text{ev}}\left(\prod_{i \in \mathbb{N}} A_i | (d_i)_{i \in \mathbb{N}}, (y_{ij})_{i \in D, j \in \mathbb{N}}\right) = \delta_{(y_{d_i i})_{i \in \mathbb{N}}} \left(\prod_{i \in \mathbb{N}} A_i\right) \quad (103)$$

$$= \prod_{i \in \mathbb{N}} \delta_{y_{d_i i}}(A_i) \quad (104)$$

$$= \text{copy}^{\mathbb{N}} \prod_{i \in \mathbb{N}} (\Pi_{D,i} \otimes \Pi_{Y,i}) \mathbb{F}_{\text{ev}_{Y^D \times D}} \quad (105)$$

Which is what we wanted to show.

Only if: With  $\mathbb{P}_{\square}^{Y^D}$  column exchangeable. That is, letting  $Y^D = (Y_i^D)_{i \in \mathbb{N}}$ , the  $Y_i^D$  are exchangeable with respect to  $\mathbb{P}_{\square}^{Y^D}$ . From kal (2005) we have a directing random measure  $H$  such that

$$\mathbb{P}_{\square}^{Y^D|H} = H \longrightarrow \boxed{\begin{array}{c} \boxed{\mathbb{P}_{\square}^{Y^D|H}} - Y_i \\ i \in \mathbb{N} \end{array}} \quad (106)$$

$$\iff \quad (107)$$

$$\mathbb{P}_{\square}^{Y^D|H}(\prod_{i \in \mathbb{N}} A_i | h) = \prod_{i \in \mathbb{N}} \mathbb{P}_{\square}^{Y_i^D|H}(A_i | h) \quad (108)$$

Furthermore, because  $Y$  is a deterministic function of  $D$  and  $Y^D$ ,  $Y \perp\!\!\!\perp_{\mathbb{P}_{\square}} H|(D, Y^D)$  and by definition of  $Y^D$ ,  $Y^D \perp\!\!\!\perp_{\mathbb{P}_{\square}} D$  and so

$$\mathbb{P}_{\square}^{Y|HD} = \mathbb{P}_{\square}^{Y^D|HD} \odot \mathbb{P}_{\square}^{Y|Y^DHD} \quad (109)$$

$$\begin{array}{c} H \longrightarrow \boxed{\mathbb{P}_{\square}^{Y^D|H}} \\ D \longrightarrow \boxed{\mathbb{P}_{\square}^{Y|Y^DHD}} \longrightarrow Y \end{array} = \boxed{\begin{array}{c} H \longrightarrow \boxed{\Pi_{D,i}} \boxed{\mathbb{P}_{\square}^{Y_0|HD_0}} - Y_i \\ D \longrightarrow \boxed{\Pi_{D,i}} \\ i \in \mathbb{N} \end{array}} \quad (110)$$

If: By assumption

$$\mathbb{P}_{\square}^{Y|D}(\prod_{i \in \mathbb{N}} A_i | h, (d_i)_{i \in \mathbb{N}}) = \int_H \prod_{i \in \mathbb{N}} \mathbb{P}_{\square}^{Y_1|HD_1}(A_i | h, d_i) \mathbb{P}_{\square}^H(dh) \quad (111)$$

Consider  $\alpha, \alpha'$  such that  $\mathbb{P}_{\alpha}^{D_M} = \mathbb{P}_{\alpha'}^{D_L}$  for  $L, M \subset \mathbb{N}$  with  $|M| = |L|$ , both finite. Then

$$\mathbb{P}_\alpha^{Y_M}(A) = \int_{D^N} \mathbb{P}_\alpha^{Y_M|D}(A|d) \mathbb{P}_\alpha^D(dd) \quad (112)$$

$$= \int_H \int_{D^N} \prod_{i \in M} \mathbb{P}_\square^{Y_1|HD_1}(A_i|h, d_i) \mathbb{P}_\alpha^D(dd) \mathbb{P}_\square^H(dh) \quad (113)$$

$$= \int_H \int_{D^{|M|}} \prod_{i \in M} \mathbb{P}_\square^{Y_1|HD_1}(A_i|h, d_i) \mathbb{P}_\alpha^{D_M}(dd_M) \mathbb{P}_\square^H(dh) \quad (114)$$

$$= \int_H \int_{D^{|M|}} \prod_{i \in M} \mathbb{P}_\square^{Y_1|HD_1}(A_i|h, d_i) \mathbb{P}_{\alpha'}^{D_N}(dd_N) \mathbb{P}_\square^H(dh) \quad (115)$$

$$= \int_H \int_{D^N} \prod_{i \in M} \mathbb{P}_\square^{Y_1|HD_1}(A_i|h, d_i) \mathbb{P}_{\alpha'}^D(dd) \mathbb{P}_\square^H(dh) \quad (116)$$

$$= \mathbb{P}_{\alpha'}^{Y_M}(A) \quad (117)$$

□

#### 4.4 Modelling different measurement procedures

An important question is: when is it reasonable to assume causal contractibility? We're going to focus just on the assumption of commutativity of exchange because we have more interesting things to say about it. There is a tempting but false line of argument one could adopt:  $(\mathbb{P}_\square^{Y_M|D_M}, A)$  is a model of  $|M|$  indistinguishable “experimental units”, because they are indistinguishable they can be interchanged without altering the appropriate model, and so commutativity of exchange holds.

The problem with this line of reasoning is that interchangeability of “experimental units” doesn't imply commutativity of exchange. The problem is, roughly speaking, we may have indistinguishable experimental units when a decision function is chosen, but the decision function might leave some uncertainty over the actual decisions, which means the experimental units may be distinguishable when the actual decisions are made. If the decision function is deterministic, this possibility is ruled out. We'll explain this in more detail with an example, and in the next section we'll discuss randomisation.

#### 4.5 Example: commutativity of exchange in the context of treatment choices

To justify an assumption of commutativity of exchange, we will argue as follows:

- Two measurement procedures should be considered equivalent in the sense that the same model is appropriate for both
- The models associated with the two procedures are related to one another by composition with the relevant swap maps



- Therefore the model associated with the first experiment is equivalent to the same model composed with the relevant swap maps

First, we want to spell out in detail how composing a model of one measurement procedure with a swap map can result in a model applicable to a different measurement procedure. Recall that we assume that a single master measurement procedure  $\mathcal{S}$  taking values in  $\Psi$ , and observables are all functions of  $\mathcal{S}$ . Given a model  $(\mathbb{P}_{\square}, A)$  associated with  $\mathcal{S}$ , the model does not in general apply to an alternative measurement procedure  $\mathcal{S}'$ .

However, it is also a principle of measurement procedures that a measurement procedure followed by the application of a function is itself a measurement procedure. Thus a model  $(\mathbb{P}_{\square}, A)$  associated with  $\mathcal{S}$  may also be informative about a procedure  $f \circ \mathcal{S}$  for any  $f : \Psi \rightarrow X$ .

In particular, consider measurement procedures related by *swaps*. For example, suppose we have  $(\mathcal{D}_1, \mathcal{D}_2)$  and  $(\mathcal{D}_1^{\text{swap}}, \mathcal{D}_2^{\text{swap}}) := (\mathcal{D}_2, \mathcal{D}_1)$ . Then, given any probability model  $\mathbb{P}_{\alpha}^{\mathcal{D}_1 \mathcal{D}_2}$  we have  $\mathbb{P}_{\alpha}^{\mathcal{D}_1^{\text{swap}} \mathcal{D}_2^{\text{swap}}} = \mathbb{P}_{\alpha}^{\mathcal{D}_1 \mathcal{D}_2}$ . In this way,  $\mathbb{P}_{\alpha}^{\mathcal{D}_1 \mathcal{D}_2}$  is a model of  $(\mathcal{D}_1, \mathcal{D}_2)$  and induces a unique model of  $(\mathcal{D}_1^{\text{swap}}, \mathcal{D}_2^{\text{swap}})$  via composition with a swap map.

Technically, this requires an assumption: if  $X$  is associated with  $\mathcal{X}$  then  $f \circ X$  is associated with  $f \circ \mathcal{X}$  (roughly: the abstract mathematical idea of composing a function with something and the actual process of applying a function to something and obtaining a result are treated as the same thing)

Concretely, commutativity of exchange can be justified if we suppose that the same model  $(\mathbb{P}_{\square}^{Y_M | D}{}^M, A)$  should describe

- A measurement procedure  $\mathcal{S}$  that yields  $|M|$  outcomes  $\mathcal{Y}_M$  and  $|M|$  decisions  $\mathcal{D}_M$
- Any other  $|M|$  outcomes  $\mathcal{Y}_M^{\text{swap}}$  and  $|M|$  decisions  $\mathcal{D}_M^{\text{swap}}$ , related to the originals by a swap.

Consider the following two scenarios:

1. Dr Alice is going to see two patients who are both complaining of lower back pain and are otherwise unknown to Alice. Prior to seeing them, she settles on a decision function  $\alpha$  which deterministically sets her treatment choices according to a function decisions( $\alpha$ )
2. As before, but  $\alpha$  is a “decision inclination” and  $\mathbb{P}_{\alpha}^{\mathcal{D}_1 \mathcal{D}_2}$  nondeterministic

Alice could model both situations with a sequential conditional probability model  $(\mathbb{P}_{\square}^{Y_1 Y_2 | D_1 D_2}, A)$  with the elements of  $A$  identified with probability models of the form  $\mathbb{P}_{\alpha}^{\mathcal{D}_1 \mathcal{D}_2}$ . Might she, in one or both situations, consider this conditional probability model to be causally contractible?

We will assume that both satisfy commutativity of marginalisation – that is, the first patient’s outcomes are expected to be the same no matter what is planned for the second patient and vice versa. We want to know if they satisfy commutativity of exchange.

The argument we want to make (if it can be supported) is:

- We can describe two measurement procedures that should share the same model
- The first is a measurement procedure for  $(D_1, D_2, Y_1, Y_2)$
- The second is a measurement procedure for  $(D_1^{\text{swap}}, D_2^{\text{swap}}, Y_1^{\text{swap}^{-1}}, Y_2^{\text{swap}^{-1}})$

At the outset, Alice does not know any features that might distinguish the two patients, so it is reasonable to think that she should adopt the same model for a) the original experiment and b) the same experiment, except with the patients interchanged. Note that interchanging *patients* does not correspond directly to any operation on the model  $(\mathbb{P}_{\square}^{Y_1 Y_2 | D_1 D_2}, A)$  which describes decisions and, not patients.

We will define measurement procedures using pseudocode, because we find it a lot easier to keep track of operations like swaps in this format. This presentation has the unintended effect of suggesting that measurement procedures are like computer programs. We're not sure if this is a helpful way to think about things – one of the key points of this example is that precise and imprecise measurement procedures may need quite different models, but thinking of measurement procedures as computer programs suggests that all measurement procedures are precise, which is not the case. Some steps may be precise, and we can express these steps with pseudocode, while other steps may be less precise.

Suppose the first scenario corresponds to the following procedure  $\mathcal{S}$  which yields values in  $A \times D^2 \times Y^2$ .  $D_i$  is the projection  $(\alpha, d_1, d_2, y_1, y_2) \mapsto d_i$  composed with  $\mathcal{S}$  and  $Y_i$  is the projection  $(\alpha, d_1, d_2, y_1, y_2) \mapsto y_i$  composed with  $\mathcal{S}$ .

**procedure  $\mathcal{S}$**

```

assert(patient A knowledge=patient B knowledge)
 $\alpha \leftarrow \text{choose\_alpha}$ 
 $(\mathcal{D}_1, \mathcal{D}_2) \leftarrow \text{decisions}(\alpha)$ 
 $\mathcal{Y}_1 \leftarrow \text{apply}(\mathcal{D}_1, \text{patient A})$ 
 $\mathcal{Y}_2 \leftarrow \text{apply}(\mathcal{D}_2, \text{patient B})$ 
return  $(\alpha, \mathcal{D}_1, \mathcal{D}_2, \mathcal{Y}_1, \mathcal{Y}_2)$ 

```

**end procedure**

Make the assumption that, on the basis that the patients are indistinguishable to Alice at the time of model construction, the same model is appropriate for the original measurement procedure and a modified measurement procedure in which the patients are swapped (we say the measurement procedures are “equivalent”). Assume also that swapping the order of treatment and swapping the order in which outcomes are recorded yields an equivalent measurement procedure (in Walley (1991)’s language, the first assumption is based on “symmetry of evidence” and the second on “evidence of symmetry”). Putting these two assumptions together, the following procedure  $\mathcal{S}'$  is equivalent to the original:

**procedure  $\mathcal{S}'$**

```

assert(patient A knowledge=patient B knowledge)
 $\alpha \leftarrow \text{choose\_alpha}$ 
 $(\mathcal{D}_1, \mathcal{D}_2) \leftarrow \text{decisions}(\alpha)$ 

```

```

 $y_2 \leftarrow \text{apply}(\mathcal{D}_2, \text{patient A})$ 
 $y_1 \leftarrow \text{apply}(\mathcal{D}_1, \text{patient B})$ 
return  $(\alpha, \mathcal{D}_1, \mathcal{D}_2, y_1, y_2)$ 
end procedure

```

Consider another measurement procedure  $\mathcal{S}''$ , which is a modified version of  $\mathcal{S}$  where steps are added to swap decisions after they are chosen, then outcomes are swapped back once they have been observed:

```

procedure  $\mathcal{S}''$ 
  assert(patient A knowledge=patient B knowledge)
   $\alpha \leftarrow \text{choose\_alpha}$ 
   $(\mathcal{D}_1, \mathcal{D}_2) \leftarrow \text{decisions}(\alpha)$ 
   $(\mathcal{D}_1^{\text{swap}}, \mathcal{D}_2^{\text{swap}}) \leftarrow (\mathcal{D}_2, \mathcal{D}_1)$ 
   $y_1^{\text{swap}} \leftarrow \text{apply}(\mathcal{D}_1^{\text{swap}}, \text{patient A})$ 
   $y_2^{\text{swap}} \leftarrow \text{apply}(\mathcal{D}_2^{\text{swap}}, \text{patient B})$ 
   $(y_1, y_2) \leftarrow (y_2^{\text{swap}}, y_1^{\text{swap}})$ 
  return  $(\alpha, \mathcal{D}_1, \mathcal{D}_2, y_1, y_2)$ 
end procedure

```

Instead of explicitly performing the swaps, we can substitute  $\mathcal{D}_2$  for  $\mathcal{D}_1^{\text{swap}}$ ,  $y_2$  for  $y_1^{\text{swap}}$  and so on. The result is a procedure identical to  $\mathcal{S}'$

```

procedure  $\mathcal{S}''$ 
  assert(patient A knowledge=patient B knowledge)
   $\alpha \leftarrow \text{choose\_alpha}$ 
   $(\mathcal{D}_1, \mathcal{D}_2) \leftarrow \text{decisions}(\alpha)$ 
   $y_2 \leftarrow \text{apply}(\mathcal{D}_2, \text{patient A})$ 
   $y_1 \leftarrow \text{apply}(\mathcal{D}_1, \text{patient B})$ 
  return  $(\alpha, \mathcal{D}_1, \mathcal{D}_2, y_1, y_2)$ 
end procedure

```

Thus  $\mathcal{S}''$  is exactly the same as  $\mathcal{S}'$ , which by assumption is equivalent to the original  $\mathcal{S}$ , and so the assumptions of interchangeable patients and reversible order of treatment application imply the model should commute with exchange. Thus, if we could extend this example to an infinite sequence of patients, there would exist a Markov kernel  $\mathbb{P}_{\square}^{Y|\text{DH}} : D \times H \rightarrow Y$  representing a “definite but unknown causal consequence” shared by all experimental units.

This argument does *not* hold for scenario 2. In the absence of a deterministic function  $\text{decisions}(\alpha)$  which defines the procedure for obtaining  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , there is some flexibility for how exactly these variables are measured (or chosen). In particular, we can posit measurement procedures such that permuting patients is not equivalent to permuting decisions and then applying the reverse permutation to outcomes.

For example, procedure  $\mathcal{T}$  is compatible with scenario 2 (note that there are many procedures compatible with the given description)

```

procedure  $\mathcal{T}$ 
  assert(patient A knowledge=patient B knowledge)
   $\alpha \leftarrow \text{choose\_alpha}$ 

```

```

    patient A knowledge ← inspect(patient A)
    patient B knowledge ← inspect(patient B)
    ( $\mathcal{D}_1, \mathcal{D}_2$ ) ← vagueDecisions( $\alpha$ , patient A knowledge, patient B knowledge)
     $\mathcal{Y}_1$  ← apply( $\mathcal{D}_1$ , patient A)
     $\mathcal{Y}_2$  ← apply( $\mathcal{D}_2$ , patient B)
    return ( $\alpha, \mathcal{D}_1, \mathcal{D}_2, \mathcal{Y}_1, \mathcal{Y}_2$ )
end procedure

    Permutation of patients and treatment order now yields
procedure  $\mathcal{T}'$ 
    assert(patient A knowledge=patient B knowledge)
     $\alpha$  ← choose_α
    patient B knowledge ← inspect(patient B)
    patient A knowledge ← inspect(patient A)
    ( $\mathcal{D}_1, \mathcal{D}_2$ ) ← vagueDecisions( $\alpha$ , patient B knowledge, patient A knowledge)
     $\mathcal{Y}_2$  ← apply( $\mathcal{D}_2$ , patient A)
     $\mathcal{Y}_1$  ← apply( $\mathcal{D}_1$ , patient B)
    return ( $\alpha, \mathcal{D}_1, \mathcal{D}_2, \mathcal{Y}_1, \mathcal{Y}_2$ )
end procedure

    While paired permutation of decisions and outcomes yields
procedure  $\mathcal{T}''$ 
    assert(patient A knowledge=patient B knowledge)
     $\alpha$  ← choose_α
    patient A knowledge ← inspect(patient A)
    patient B knowledge ← inspect(patient B)
    ( $\mathcal{D}_1, \mathcal{D}_2$ ) ← vagueDecisions( $\alpha$ , patient A knowledge, patient B knowledge)
     $\mathcal{Y}_2$  ← apply( $\mathcal{D}_2$ , patient A)
     $\mathcal{Y}_1$  ← apply( $\mathcal{D}_1$ , patient B)
    return ( $\alpha, \mathcal{D}_1, \mathcal{D}_2, \mathcal{Y}_1, \mathcal{Y}_2$ )
end procedure

```

$\mathcal{T}'$  is not the same as  $\mathcal{T}''$ . In scenario 1, because decisions were deterministic on  $\alpha$ , there was no room to pick anything different once  $\alpha$  was chosen, so it doesn't matter if we add patient inspection steps or not. In scenario 2, decisions are not deterministic and there is vagueness in the procedure, so it is possible to describe compatible procedures where decisions depend on patient characteristics, and this dependence is not “undone” by swapping decisions.

## 4.6 Causal consequences of non-deterministic variables

In the previous section we gave an example of how commutativity of exchange can hold when we have a sequence of decisions such that we accept the following:

- Reordering the time at which decisions are made is held to be of no consequence
- The available information relevant to each decision is symmetric at the time the decision function is adopted

- The decision function deterministically prescribes which decisions are taken

We also discussed how the absence of determinism undermines the argument for exchange commutativity.

The determinism assumption rules out choosing decisions randomly. However, if we have causal consequences for deterministic decision variables, it is sometimes possible to extend them to indeterministic variables.

**Lemma 4.8.** *Given  $(\mathbb{P}_\square, A)$  with decisions  $D_M$  and consequences  $Y_M$ , if  $\mathbb{P}_\square^{Y_M|D_M}$  is causally contractible with consequence map  $\mathbb{P}_\square^{Y_0|D_0H}$  and there exists  $X_i = f \circ Y_i$  for some  $f : Y \rightarrow X$  such that  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} D_i | HX_i$  for all  $i \in M$ , then a causally contractible conditional probability  $\mathbb{P}_\square^{Y_M|X_M}$  exists.*

*Proof.* We want to show  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} Y_{\{i\}^c} X_{\{i\}^c} | HX_i$  for all  $i \in M$ ,  $\mathbb{P}_\square^{Y_i|HX_i}$  exists for all  $i \in M$  and  $\mathbb{P}_\square^{Y_i|HX_i} = \mathbb{P}_\square^{Y_j|HX_j}$ .

Because  $X_i$  is a function of  $Y_i$ , and  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} Y_{\{i\}^c} D_{\{i\}^c} | HD_i$ , we also have  $YX_i \perp\!\!\!\perp_{\mathbb{P}_\square} Y_{\{i\}^c} X_{\{i\}^c} | HD_i$ , and by weak union  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} Y_{\{i\}^c} X_{\{i\}^c} | HD_i X_i$ .

Thus by contraction,  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} Y_{\{i\}^c} D_M | HX_i$ .

By Corollary 5.22 and the existence of  $\mathbb{P}_\square^{Y_i X_i | HD_i}$  for all  $i \in M$ ,  $\mathbb{P}_\square^{Y_i | HX_i}$  exists for all  $i$ . Furthermore, because  $\mathbb{P}_\square^{Y_i X_i | HD_i} = \mathbb{P}_\square^{Y_j X_j | HD_j}$  for all  $i, j \in M$ ,  $\mathbb{P}_\square^{Y_i | HX_i} = \mathbb{P}_\square^{Y_j | HX_j}$  for all  $i, j \in M$ .  $\square$

If the condition  $Y_i \perp\!\!\!\perp_{\mathbb{P}_\square} D_i | HX_i$  for all  $i \in M$ , we can say  $X_i$  is a proxy for controlling  $Y_i$ .

As an example of this, suppose  $X : \Omega \rightarrow X$  is a source of random numbers, the set of decisions  $D$  is a set of functions  $X \rightarrow T$  for treatments  $T : \Omega \rightarrow T$  and  $W : \Omega \rightarrow W$  are the ultimate patient outcomes, with  $Y_i = (W_i, T_i)$ . Then it may be reasonable to assume that  $W_i \perp\!\!\!\perp (D_i, X_i) | T_i H$  (where conditioning on  $H$  can be thought of as saying that this independence holds under infinite sample size). In this case,  $T_i$  is a proxy for controlling  $Y_i$ , and there exists a causal consequence  $\mathbb{P}_\square^{Y_0 | T_0 H}$ .

A “causal consequence of body mass index” is unlikely to exist on the basis of symmetric information and deterministic decisions because there are no actions available to set body mass index deterministically. However, given an underlying problem where we have symmetric information over a collection of patients and some kind of decision that can be made deterministically, causal consequences of body mass index may exist if body mass index is a proxy for controlling the outcomes of interest.

## 4.7 Intersubjective causal consequences

While the assumption of causal contractibility itself does not depend on any notion of subjectivity, our discussion of the applicability of this assumption assumed that a conditional probability model was being used to model Dr Alice’s subjective uncertain knowledge. Crucially, the justification hinged on an assumption of the symmetry of Alice’s information regarding different patients.

Causal inference is often performed in an intersubjective setting, where Ben might perform the experiment, Carmel might do the analysis and Dr Alice make the ultimate decisions. This complicates the question of when the assumption of causal contractibility is applicable. We leave the appropriate way to generalise this theory to such a setting open.

## 5 Appendix, needs to be organised

### 5.1 Existence of conditional probabilities

**Lemma 5.1** (Conditional pushforward). *Suppose we have a sample space  $(\Omega, \mathcal{F})$ , variables  $X : \Omega \rightarrow X$  and  $Y : \Omega \rightarrow Y$ ,  $Z : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_{\{\}}^X$  with conditional  $\mathbb{P}_{\{\}}^{X|Y}$  such that  $Z = f \circ Y$  for some  $f : Y \rightarrow Z$ . Then there exists a conditional probability  $\mathbb{P}_{\{\}}^{Z|X} = \mathbb{P}_{\{\}}^{Y|X} \mathbb{F}_f$ .*

*Proof.* Note that  $(X, Z) = (\text{id}_X \otimes f) \circ (X, Y)$ . Thus, by Lemma 2.8, for any  $\mathbb{P}_\alpha \in \mathbb{P}_{\{\}}$

$$\mathbb{P}_\alpha^{XZ} = \mathbb{P}_\alpha^{XY} \mathbb{F}_{\text{id}_X \otimes f} \quad (118)$$

Note also that for all  $A \in \mathcal{X}$ ,  $B \in \mathcal{Z}$ ,  $x \in X$ ,  $y \in Y$ :

$$\mathbb{F}_{\text{id}_X \otimes f}(A \times B|x, y) = \delta_x(A) \delta_{f(y)}(B) \quad (119)$$

$$= \mathbb{F}_{\text{id}_X}(A|x) \otimes \mathbb{F}_f(B|y) \quad (120)$$

$$\implies \mathbb{F}_{\text{id}_X \otimes f} = \mathbb{F}_{\text{id}_X} \otimes \mathbb{F}_f \quad (121)$$

Thus

$$\mathbb{P}_\alpha^{XZ} = (\mathbb{P}_\alpha^X \odot \mathbb{P}_{\{\}}^{Y|X}) \mathbb{F}_{\text{id}_X} \otimes \mathbb{F}_f \quad (122)$$

$$= \begin{array}{c} \text{X} \\ \curvearrowright \\ \begin{array}{c} \triangleleft \mathbb{P}_\alpha^X \\ \bullet \\ \boxed{\mathbb{P}_{\{\}}^{Y|X}} \end{array} \xrightarrow{\quad} \boxed{\mathbb{F}_f} \xrightarrow{\quad} Z \end{array} \quad (123)$$

Which implies  $\mathbb{P}_{\{\}}^{Y|X} \mathbb{F}_f$  is a version of  $\mathbb{P}_\alpha^{Z|X}$ . Because this holds for all  $\alpha$ , it is therefore also a version of  $\mathbb{P}_{\{\}}^{Z|X}$ .  $\square$

**Theorem 5.2** (Existence of regular conditionals). *Suppose we have a sample space  $(\Omega, \mathcal{F})$ , variables  $X : \Omega \rightarrow X$  and  $Y : \Omega \rightarrow Y$  with  $Y$  standard measurable and a probability model  $\mathbb{P}_\alpha$  on  $(\Omega, \mathcal{F})$ . Then there exists a conditional  $\mathbb{P}_\alpha^{Y|X}$ .*

*Proof.* This is a standard result, see for example Çinlar (2011) Theorem 2.18.  $\square$

**Theorem 5.3** (Existence of higher order valid conditionals with respect to probability sets). *Suppose we have a sample space  $(\Omega, \mathcal{F})$ , variables  $\mathbf{X} : \Omega \rightarrow X$  and  $\mathbf{Y} : \Omega \rightarrow Y$ ,  $\mathbf{Z} : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  with regular conditional  $\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  and  $Y$  and  $Z$  standard measurable. Then there exists a regular  $\mathbb{P}_{\{\}}^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$ .*

*Proof.* Given a Borel measurable map  $m : X \rightarrow Y \times Z$  let  $f : Y \times Z \rightarrow Y$  be the projection onto  $Y$ . Then  $f \circ (\mathbf{Y}, \mathbf{Z}) = \mathbf{Y}$ . Bogachev and Malofeev (2020), Theorem 3.5 proves that there exists a Borel measurable map  $n : X \times Y \rightarrow Y \times Z$  such that

$$n(f^{-1}(y)|x, y) = 1 \quad (124)$$

$$m(\mathbf{Y}^{-1}(A) \cap B|x) = \int_A n(B|x, y) m\mathbb{F}_f(dy|x) \forall A \in \mathcal{Y}, B \in \mathcal{Y} \times \mathcal{Z} \quad (125)$$

In particular,  $\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  is a Borel measurable map  $X \rightarrow Y \times Z$ . Thus equation 125 implies for all  $A \in \mathcal{Y}, B \in \mathcal{Y} \times \mathcal{Z}$

$$\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}(\mathbf{Y}^{-1}(A) \cap B|x) = \int_A n(B|x, y) \mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}} \mathbb{F}_f(dy|x) \quad (126)$$

$$= \int_A n(B|x, y) \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (127)$$

Where Equation 127 follows from Lemma 5.1.

Then, for any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$

$$\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}(\mathbf{Y}^{-1}(A) \cap B|x) = \int_A n(B|x, y) \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (128)$$

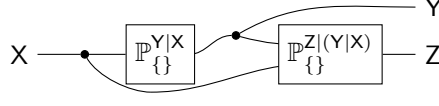
which implies  $n$  is a version of  $\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{Y}|\mathbf{X}}$ . By Lemma 5.1,  $n\mathbb{F}_f$  is a version of  $\mathbb{P}_{\{\}}^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$ .  $\square$

We might be motivated to ask whether the higher order conditionals in Theorem 5.3 can be chosen to be valid. Despite Lemma 5.8 showing that the existence of proper conditional probabilities implies the existence of valid ones, we cannot make use of this in the above theorem because Equation 124 makes  $n$  proper with respect to the “wrong” sample space  $(Y \times Z, \mathcal{Y} \otimes \mathcal{Z})$  while what we would need is a proper conditional probability with respect to  $(\Omega, \mathcal{F})$ .

We can choose higher order conditionals to be valid in the case of discrete sets, and whether we can choose them to be valid in more general measurable spaces is an open question.

**Theorem 5.4** (Higher order conditionals). *Suppose we have a sample space  $(\Omega, \mathcal{F})$ , variables  $\mathbf{X} : \Omega \rightarrow X$  and  $\mathbf{Y} : \Omega \rightarrow Y$ ,  $\mathbf{Z} : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_{\{\}}$  with conditional  $\mathbb{P}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$ . Then  $\mathbb{P}_{\{\}}^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$  is a version of  $\mathbb{P}_{\{\}}^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$ .*

*Proof.* For arbitrary  $\mathbb{P}_\alpha \in \mathbb{P}_\{\}$



$$\mathbb{P}_\alpha^{YZ|X} = \quad (129)$$

$$\Rightarrow \mathbb{P}_\alpha^{XYZ} = \triangleleft \mathbb{P}_\alpha^X \begin{array}{c} \text{---} X \\ \text{---} Y \\ \text{---} Z \end{array} \quad (130)$$

$$= \triangleleft \mathbb{P}_\alpha^X \begin{array}{c} \text{---} X \\ \text{---} Y \\ \text{---} Z \end{array} \quad (131)$$

$$= \triangleleft \mathbb{P}_\alpha^{XY} \begin{array}{c} \text{---} X \\ \text{---} Y \\ \text{---} Z \end{array} \quad (132)$$

Thus  $\mathbb{P}_\{\}^{Z|(Y|X)}$  is a version of  $\mathbb{P}_\alpha^{Z|YX}$  for all  $\alpha$  and hence also a version of  $\mathbb{P}_\{\}^{Z|YX}$ .  $\square$

**Theorem 5.5.** *Given probability gap model  $\mathbb{P}_\{\}$ ,  $X, Y, Z$  such that  $\mathbb{P}_\{\}^{Z|YX}$  exists,  $\mathbb{P}_\{\}^{Z|Y}$  exists iff  $Z \perp\!\!\!\perp_{\mathbb{P}_\{\}} X|Y$ .*

*Proof.* If: If  $Z \perp\!\!\!\perp_{\mathbb{P}_\{\}} X|Y$  then by Theorem 5.15, for each  $\mathbb{P}_\alpha \in \mathbb{P}_\{\}$  there exists  $\mathbb{P}_\alpha^{Z|Y}$  such that

$$\mathbb{P}_\alpha^{Y|WX} = \begin{array}{c} W \text{---} \boxed{\mathbb{P}_{\square}^{Y|W}} \text{---} Y \\ X \text{---} * \end{array} \quad (133)$$

$\square$

**Theorem 5.6** (Valid higher order conditionals). *Suppose we have a sample space  $(\Omega, \mathcal{F})$ , variables  $X : \Omega \rightarrow X$  and  $Y : \Omega \rightarrow Y$ ,  $Z : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_\{\}$  with regular conditional  $\mathbb{P}_\{\}^{YZ|X}$ ,  $Y$  discrete and  $Z$  standard measurable. Then there exists a valid regular  $\mathbb{P}_\{\}^{Z|XY}$ .*

*Proof.* By Theorem 5.3, we have a higher order conditional  $\mathbb{P}_\{\}^{Z|(Y|X)}$  which, by Theorem 5.4 is also a version of  $\mathbb{P}_\{\}^{Z|XY}$ .



We will show that there is a Markov kernel  $\mathbb{Q}$  almost surely equal to  $\mathbb{P}_{\{\}}^{Z|XY}$  which is also valid. For all  $x, y \in X \times Y$ ,  $A \in \mathcal{Z}$  such that  $(X, Y, Z) \bowtie \{(x, y)\} \times A = \emptyset$ , let  $\mathbb{Q}(A|x, y) = \mathbb{P}_{\{\}}^{Z|XY}(A|x, y)$ .

By validity of  $\mathbb{P}_{\{\}}^{YZ|X}$ ,  $x \in X(\Omega)$  and  $(X, Y, Z) \bowtie \{(x, y)\} \times A = \emptyset$  implies  $\mathbb{P}_{\{\}}^{YZ|X}(\{y\} \times A|x) = 0$ . Thus we need to show

$$\forall A \in \mathcal{Z}, x \in X, y \in Y : \mathbb{P}_{\{\}}^{YZ|X}(\{y\} \times A|x) = 0 \implies (\mathbb{Q}(A|x, y) = 0) \vee ((X, Y) \bowtie \{(x, y)\} = \emptyset) \quad (134)$$

For all  $x, y$  such that  $\mathbb{P}_{\{\}}^{Y|X}(\{y\}|x)$  is positive, we have  $\mathbb{P}_{\{\}}^{YZ|X}(\{y\} \times A|x) = 0 \implies \mathbb{P}_{\square}^{Z|XY}(A|x, y) = 0 =: \mathbb{Q}(A|x, y)$ .

Furthermore, where  $\mathbb{P}_{\{\}}^{Y|X}(\{y\}|x) = 0$ , we either have  $(X, Y, Z) \bowtie \{(x, y)\} \times A = \emptyset$  or can choose some  $\omega \in (X, Y, Z) \bowtie \{(x, y)\} \times A$  and let  $\mathbb{Q}(Z(\omega)|x, y) = 1$ . This is an arbitrary choice, and may differ from the original  $\mathbb{P}_{\{\}}^{Z|XY}$ . However, because  $Y$  is discrete the union of all points  $y$  where  $\mathbb{P}_{\{\}}^{Y|X}(\{y\}|x) = 0$  is a measure zero set, and so  $\mathbb{Q}$  differs from  $\mathbb{P}_{\{\}}^{Y|X}$  on a measure zero set.  $\square$

## 5.2 Validity

Validity is related to *proper* conditional probabilities. In particular, valid conditional probabilities exist when regular proper conditional probabilities exist.

**Definition 5.7** (Regular proper conditional probability). Given a probability space  $(\mu, \Omega, \mathcal{F})$  and a variable  $X : \Omega \rightarrow X$ , a regular proper conditional probability  $\mu^{|\mathbf{X}} : X \rightarrow \Omega$  is Markov kernel such that

$$\mu(A \cap X^{-1}(B)) = \int_B \mu^{|\mathbf{X}}(A|x) \mu^X(dx) \quad \forall A \in \mathcal{X}, B \in \mathcal{F} \quad (135)$$

$$\iff \quad (136)$$

$$\mu = \triangleleft \mu^X \quad \begin{array}{c} \text{---} \text{X} \\ \text{---} \mu^{Y|X} \text{---} Y \end{array} \quad (137)$$

and

$$\mu^{|\mathbf{X}}(X^{-1}(A)|x) = \delta_x(A) \quad (138)$$

**Lemma 5.8.** Given a probability space  $(\mu, \Omega, \mathcal{F})$  and variables  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ , if there is a regular proper conditional probability  $\mu^{|\mathbf{X}} : X \rightarrow \Omega$  then there is a valid conditional distribution  $\mu^{Y|X}$ .

*Proof.* Take  $\mathbb{K} = \mu^{|\mathbf{X}} \mathbb{F}_Y$ . We will show that  $\mathbb{K}$  is valid, and a version of  $\mu^{Y|X}$ .

Defining  $\mathbf{O} := \text{id}_\Omega$  (the identity function  $\Omega \rightarrow \Omega$ ),  $\mu^{|\mathbf{X}}$  is a version of  $\mu^{\mathbf{O}|\mathbf{X}}$ . Note also that  $\mathbf{Y} = \mathbf{Y} \circ \mathbf{O}$ . Thus by Lemma 5.1,  $\mathbb{K}$  is a version of  $\mu^{\mathbf{Y}|\mathbf{X}}$ .

It remains to be shown that  $\mathbb{K}$  is valid. Consider some  $x \in X$ ,  $A \in \mathcal{Y}$  such that  $\mathbf{X}^{-1}(\{x\}) \cap \mathbf{Y}^{-1}(A) = \emptyset$ . Then by the assumption  $\mu^{|\mathbf{X}}$  is proper

$$\mathbb{K}(\mathbf{Y} \bowtie A | x) = \delta_x(\mathbf{Y}^{-1}(A)) \quad (139)$$

$$= 0 \quad (140)$$

Thus  $\mathbb{K}$  is valid.  $\square$

**Theorem 5.9** (Validity). *Given  $(\Omega, \mathcal{F})$ ,  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbb{J} \in \Delta(X)$  with  $\Omega$  and  $X$  standard measurable, there exists some  $\mu \in \Delta(\Omega)$  such that  $\mu^{\mathbf{X}} = \mathbb{J}$  if and only if  $\mathbb{J}$  is a valid distribution.*

*Proof.* If: This is a Theorem 2.5 of Ershov (1975). Only if: This is also found in Ershov (1975), but is simple enough to reproduce here. Suppose  $\mathbb{J}$  is not a valid probability distribution. Then there is some  $x \in X$  such that  $\mathbf{X} \bowtie x = \emptyset$  but  $\mathbb{J}(x) > 0$ . Then

$$\mu^{\mathbf{X}}(x) = \mu(\mathbf{X} \bowtie x) \quad (141)$$

$$= \sum_{x' \in X} \mathbb{J}(x') \mathbb{K}(\mathbf{X} \bowtie x | x') \quad (142)$$

$$= 0 \quad (143)$$

$$\neq \mathbb{J}(x) \quad (144)$$

$\square$

**Lemma 5.10** (Semidirect product defines an intersection of probability sets). *Given  $(\Omega, \mathcal{F})$ ,  $\mathbf{X} : \Omega \rightarrow (X, \mathcal{X})$ ,  $\mathbf{Y} : \Omega \rightarrow (Y, \mathcal{Y})$ ,  $\mathbf{Z} : \Omega \rightarrow (Z, \mathcal{Z})$  all standard measurable and maximal probability sets  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}[M]}$  and  $\mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}[M]}$  then defining*

$$\mathbb{R}_{\{\}}^{\mathbf{YZ}|\mathbf{X}} := \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (145)$$

*we have*

$$\mathbb{R}_{\{\}}^{\mathbf{YZ}|\mathbf{X}[M]} = \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}[M]} \cap \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}[M]} \quad (146)$$

*Proof.* For any  $\mathbb{R}_a \in \mathbb{R}_{\{\}}$

$$\mathbb{R}_a^{\mathbf{XYZ}} = \mathbb{R}_a^{\mathbf{X}} \odot \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (147)$$

$$\implies \mathbb{R}_a^{\mathbf{XY}} = \mathbb{R}_a^{\mathbf{X}} \odot \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \quad (148)$$

$$\wedge \mathbb{R}_a^{\mathbf{XYZ}} = \mathbb{R}_a^{\mathbf{XY}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (149)$$

Thus  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  is a version of  $\mathbb{R}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  and  $\mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}}$  is a version of  $\mathbb{R}_{\{\}}^{\mathbf{Z}|\mathbf{YX}}$  so  $\mathbb{R}_{\{\}} \subset \mathbb{P}_{\{\}} \cap \mathbb{Q}_{\{\}}$ .

Suppose there's an element  $\mathbb{S}$  of  $\mathbb{P}_{\{\}} \cap \mathbb{Q}_{\{\}}$  not in  $\mathbb{R}_{\{\}}$ . Then by definition of  $\mathbb{R}_{\{\}}$ ,  $\mathbb{R}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  is not a version of  $\mathbb{S}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$ . But by construction of  $\mathbb{S}$ ,  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  is a version of  $\mathbb{S}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  and  $\mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}}$  is a version of  $\mathbb{S}_{\{\}}^{\mathbf{Z}|\mathbf{YX}}$ . But then by the definition of disintegration,  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}}$  is a version of  $\mathbb{S}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  and so  $\mathbb{R}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$  is a version of  $\mathbb{S}_{\{\}}^{\mathbf{YZ}|\mathbf{X}}$ , a contradiction.  $\square$

**Lemma 5.11** (Equivalence of validity definitions). *Given  $\mathbf{X} : \Omega \rightarrow X$ , with  $\Omega$  and  $X$  standard measurable, a probability measure  $\mathbb{P}^{\mathbf{X}} \in \Delta(X)$  is valid if and only if the conditional  $\mathbb{P}^{\mathbf{X}|\ast} := \ast \mapsto \mathbb{P}^{\mathbf{X}}$  is valid.*

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

$$\forall A \in \mathcal{X} : \mathbf{X} \bowtie A = \emptyset \implies \mathbb{P}^{\mathbf{X}|\ast}(A|\ast) = 0 \quad (150)$$

But  $\mathbb{P}^{\mathbf{X}|\ast}(A|\ast) = \mathbb{P}^{\mathbf{X}}(A)$  by definition, so this is equivalent to

$$\forall A \in \mathcal{X} : \mathbf{X} \bowtie A = \emptyset \implies \mathbb{P}^{\mathbf{X}}(A) = 0 \quad (151)$$

$\square$

**Lemma 5.12** (Semidirect product of valid candidate conditionals is valid). *Given  $(\Omega, \mathcal{F})$ ,  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbf{Y} : \Omega \rightarrow Y$ ,  $\mathbf{Z} : \Omega \rightarrow Z$  (all spaces standard measurable) and any valid candidate conditional  $\mathbb{P}^{\mathbf{Y}|\mathbf{X}}$  and  $\mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}$ ,  $\mathbb{P}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}$  is also a valid candidate conditional.*

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

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

For all  $x \in \mathbf{X}(\Omega)$ ,  $B \in \mathcal{Y}$  such that  $\mathbf{X} \bowtie \{x\} \cap \mathbf{Y} \bowtie B = \emptyset$ ,  $\mathbb{P}^{\mathbf{Y}|\mathbf{X}}(B|x) = 0$  by validity. Thus for arbitrary  $C \in \mathcal{Z}$

$$\mathbb{R}^{\mathbf{YZ}|\mathbf{X}}(B \times C|x) = \int_B \mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}(C|y, x) \mathbb{P}^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (152)$$

$$\leq \mathbb{P}^{\mathbf{Y}|\mathbf{X}}(B|x) \quad (153)$$

$$= 0 \quad (154)$$

For all  $\{x\} \times B$  such that  $\mathbf{X} \bowtie \{x\} \cap \mathbf{Y} \bowtie B \neq \emptyset$  and  $C \in \mathcal{Z}$  such that  $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \bowtie \{x\} \times B \times C = \emptyset$ ,  $\mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}(C|y, x) = 0$  for all  $y \in B$  by validity. Thus:

$$\mathbb{R}^{\mathbf{YZ}|\mathbf{X}}(B \times C|x) = \int_B \mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}(C|y, x) \mathbb{P}^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (155)$$

$$= 0 \quad (156)$$

$\square$

**Corollary 5.13** (Valid conditionals are validly extendable to valid distributions). *Given  $\Omega$ ,  $U : \Omega \rightarrow U$ ,  $W : \Omega \rightarrow W$  and a valid conditional  $\mathbb{T}^{W|U}$ , then for any valid conditional  $\mathbb{V}^U$ ,  $\mathbb{V}^U \odot \mathbb{T}^{W|U}$  is a valid probability.*

*Proof.* Applying Lemma 5.12 choosing  $X = *$ ,  $Y = U$ ,  $Z = W$  and  $\mathbb{P}^{Y|X} = \mathbb{V}^{U|*}$  and  $\mathbb{Q}^{Z|YX} = \mathbb{T}^{W|U*}$  we have  $\mathbb{R}^{WU|*} := \mathbb{V}^{U|*} \odot \mathbb{T}^{W|U*}$  is a valid conditional probability. Then  $\mathbb{R}^{WU} \cong \mathbb{R}^{WU|*}$  is valid by Theorem 5.11.  $\square$

**Theorem 5.14** (Validity of conditional probabilities). *Suppose we have  $\Omega$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ , with  $\Omega$ ,  $X$ ,  $Y$  discrete. A conditional  $\mathbb{T}^{Y|X}$  is valid if and only if for all valid candidate distributions  $\mathbb{V}^X$ ,  $\mathbb{V}^X \odot \mathbb{T}^{Y|X}$  is also a valid candidate distribution.*

*Proof.* If: this follows directly from Corollary 5.13.

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

$$(\mathbb{V}^X \odot \mathbb{T}^{Y|X})(x, y) = \mathbb{T}^{Y|X}(y|x) \mathbb{V}^X(x) \quad (157)$$

$$= \mathbb{T}^{Y|X}(y|x) \quad (158)$$

$$> 0 \quad (159)$$

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

### 5.3 Conditional independence

**Theorem 5.15.** *Given standard measurable  $\Omega$ , a probability model  $\mathbb{P}$  and variables  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ ,  $Y \perp\!\!\!\perp_{\mathbb{P}} X|W$  if and only if there exists some version of  $\mathbb{P}^{Y|WX}$  and  $\mathbb{P}^{Y|W}$  such that*

$$\mathbb{P}^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{P}^{Y|W}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (160)$$

$$\iff \mathbb{P}^{Y|WX}(y|w, x) = \mathbb{P}^{Y|W}(y|w) \quad (161)$$

*Proof.* See Cho and Jacobs (2019).  $\square$

### 5.4 Extended conditional independence

Constantinou and Dawid (2017) introduced the idea of *extended conditional independence*, which is a notion of conditional independence with respect to a parametrised collection of probability measures. It is motivated in part by the observation that such parametrised collections can be used to model causal questions. Furthermore, probability sets are closely related to parametrised probability sets – one can get the former from the latter by simply dropping the parameters.

This needs major revision, and is not a top priority right now

In the case of a probability gap model  $(\mathbb{P}_{\square}^{V|W}, A)$  where there is some  $\alpha \in A$  dominating  $A$ , we can relate conditional independence with respect to  $\mathbb{P}_{\square}$  to what , which is a notion they define with respect to a Markov kernel. These concepts may differ if  $A$  is not dominated. Theorem 4.4 of Constantinou and Dawid (2017) proves the following claim:

**Definition 5.16** (Extended conditional independence). adf

**Theorem 5.17.** *Let  $A^* = A \circ V$ ,  $B^* = B \circ V$ ,  $C^* = C \circ V$  ( $(A, B, C)$  are  $\mathcal{V}$ -measurable) and  $D^* = D \circ W$ ,  $E^* = E \circ W$  where  $W$  is discrete and  $W = (D^*, E^*)$ . In addition, let  $\mathbb{P}_{\alpha}^W$  be some probability distribution on  $W$  such that  $w \in W(\Omega) \implies \mathbb{P}_{\alpha}^W(w) > 0$ . Then, denoting extended conditional independence with  $\perp\!\!\!\perp_{\mathbb{P}, ext}$  and  $\mathbb{P}_{\alpha}^{VW} := \mathbb{P}_{\alpha}^W \odot \mathbb{P}^{V|W}$*

$$A \perp\!\!\!\perp_{\mathbb{P}, ext} (B, D)|(C, E) \iff A^* \perp\!\!\!\perp_{\mathbb{P}_{\alpha}} (B^*, D^*)|(C^*, E^*) \quad (162)$$

This result implies a close relationship between order 1 conditional independence and extended conditional independence.

**Theorem 5.18.** *Let  $A^* = A \circ V$ ,  $B^* = B \circ V$ ,  $C^* = C \circ V$  ( $(A, B, C)$  are  $\mathcal{V}$ -measurable) and  $D^* = D \circ W$ ,  $E^* = E \circ W$  where  $V, W$  are discrete and  $W = (D^*, E^*)$ . Then letting  $\mathbb{P}_{\alpha}^{VW} := \mathbb{P}_{\alpha}^W \odot \mathbb{P}^{V|W}$*

$$A \perp\!\!\!\perp_{\mathbb{P}, ext}^1 (B, D)|(C, E) \iff A^* \perp\!\!\!\perp_{\mathbb{P}} (B^*, D^*)|(C^*, E^*) \quad (163)$$

*Proof.* If:

By assumption,  $A^* \perp\!\!\!\perp_{\mathbb{P}_{\alpha}} (B^*, D^*)|(C^*, E^*)$  for all  $\mathbb{P}_{\alpha}^{D^*E^*}$ . In particular, this holds for some  $\mathbb{P}_{\alpha}^{D^*E^*}$  such that  $(d, e) \in (D^*, E^*)(\Omega) \implies \mathbb{P}_{\alpha}^{D^*E^*}(d, e) > 0$ . Then by Theorem 5.17,  $A \perp\!\!\!\perp_{\mathbb{P}, ext} (B, D)|(C, E)$ .

Only if:

For any  $\beta$ ,  $\mathbb{P}_{\beta}^{ABC|DE} = \mathbb{P}_{\beta}^{DE} \odot \mathbb{P}^{ABC|DE}$ . By Lemma 5.4, we have  $\mathbb{P}^{A|BCDE}$  such that

$$\mathbb{P}_{\beta}^{A^*B^*C^*D^*E^*} = \mathbb{P}_{\beta}^{D^*E^*} \odot \mathbb{P}^{B^*C^*|D^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (164)$$

$$= \mathbb{P}_{\beta}^{B^*C^*D^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (165)$$

$$= \mathbb{P}_{\beta}^{C^*E^*} \odot \mathbb{P}_{\beta}^{B^*D^*|C^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (166)$$

By Theorem 5.17, we have some  $\alpha$  such that  $\mathbb{P}_{\alpha}^{D^*E^*}$  is strictly positive on the range of  $(D^*, E^*)$  and  $A^* \perp\!\!\!\perp_{\mathbb{P}_{\alpha}} (B^*, D^*)|(C^*, E^*)$ .

By independence, for some version of  $\mathbb{P}^{A|BCDE}$ :

$$\mathbb{P}_\alpha^{C^*E^*} \odot \mathbb{P}_\alpha^{B^*D^*|C^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} = \begin{array}{c} \triangleleft \mathbb{P}_\alpha^{C^*E^*} \begin{array}{l} \boxed{\overline{\mathbb{P}}_\alpha^{A^*|C^*E^*}} \text{---} A^* \\ \boxed{\overline{\mathbb{P}}_\alpha^{B^*D^*|C^*E^*}} \text{---} B^*D^* \\ \text{---} C^*E^* \end{array} \end{array}$$

(167)

$$= \begin{array}{c} \triangleleft \mathbb{P}_\alpha^{C^*E^*} \begin{array}{l} \boxed{\overline{\mathbb{P}}_\alpha^{A^*|C^*E^*}} \text{---} A^* \\ \boxed{\overline{\mathbb{P}}_\alpha^{B^*D^*|C^*E^*}} \text{---} B^*D^* \\ \text{---} C^*E^* \end{array} \end{array}$$

(168)

$$= \mathbb{P}_\alpha^{C^*E^*} \odot \mathbb{P}_\alpha^{B^*D^*|C^*E^*} \odot (\mathbb{P}_\alpha^{A^*|C^*E^*} \otimes \text{erase}_{BD})$$

(169)

Thus for any  $(a, b, c, d, e) \in A \times B \times C \times D \times E$  such that  $\mathbb{P}_\alpha^{B^*C^*D^*E^*}(b, c, d, e) > 0$ ,  $\mathbb{P}^{A^*|B^*C^*D^*E^*}(a|b, c, d, e) = \mathbb{P}_\alpha^{A^*|C^*E^*}(a|c, e)$ . However, by assumption,  $\mathbb{P}_\alpha^{B^*C^*D^*E^*}(b, c, d, e) = 0 \implies \mathbb{P}_\beta^{B^*C^*D^*E^*}(b, c, d, e) = 0$ , and so  $\mathbb{P}_\beta^{A^*|B^*C^*D^*E^*} = \mathbb{P}_\alpha^{A^*|C^*E^*}(a|c, e)$  everywhere except a set of  $\mathbb{P}_\beta$ -measure 0. Thus

$$\mathbb{P}_\beta^{A^*B^*C^*D^*E^*} = \begin{array}{c} \triangleleft \mathbb{P}_\beta^{C^*E^*} \begin{array}{l} \boxed{\overline{\mathbb{P}}_\alpha^{A^*|C^*E^*}} \text{---} A^* \\ \boxed{\overline{\mathbb{P}}_\beta^{B^*D^*|C^*E^*}} \text{---} B^*D^* \\ \text{---} C^*E^* \end{array} \end{array}$$

(170)

$$= \begin{array}{c} \triangleleft \mathbb{P}_\beta^{C^*E^*} \begin{array}{l} \boxed{\overline{\mathbb{P}}_\alpha^{A^*|C^*E^*}} \text{---} A^* \\ \boxed{\overline{\mathbb{P}}_\beta^{B^*D^*|C^*E^*}} \text{---} B^*D^* \\ \text{---} C^*E^* \end{array} \end{array}$$

(171)

□

Conditional independence is a property of variables, we define “unresponsiveness” as a property of Markov kernels.

**Definition 5.19** (Unresponsiveness). Given discrete  $\Omega$ , a probability gap model  $\mathbb{P}_\square : A \rightarrow \Delta(\Omega)$ , variables  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$ , if there is some version of the conditional probability  $\mathbb{P}^{Y|WX}$  and  $\mathbb{P}_\square^{Y|W}$  such that

$$\mathbb{P}_\square^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{P}_\square^{Y|W}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (172)$$

then  $\mathbb{P}_\square^{Y|WX}$  is *unresponsive* to  $X$ .

**Definition 5.20** (Domination). Given a probability set  $\mathbb{P}_\{\}$   $\subset \Delta(\Omega)$ ,  $\mathbb{P}_\alpha$  dominates  $\mathbb{P}_\{\}$  if  $\mathbb{P}_\beta(B) > 0 \implies \mathbb{P}_\alpha(B) > 0$  for all  $\mathbb{P}_\beta \in \mathbb{P}_\{\}$ ,  $B \in \mathcal{F}$ .

**Theorem 5.21** (Conditional independence from kernel unresponsiveness). *Given standard measurable  $\Omega$ , variables  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and a probability set  $\mathbb{P}_\{\} : A \rightarrow \Delta(\Omega)$  with conditional probability  $\mathbb{P}_\{\}^{Y|WX}$  such that there is  $\mathbb{P}_\alpha \in \mathbb{P}_\{\}$  dominating  $\mathbb{P}_\{\}$ ,  $Y \perp_{\mathbb{P}_\{\}} X|W$  if and only if there is a version of  $\mathbb{P}_\{\}^{Y|WX}$  unresponsive to  $W$ .*

*Proof.* If: For every  $\alpha \in A$  we can write

$$\mathbb{P}_\alpha^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{P}_\alpha^{Y|W}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (173)$$

And so, by Theorem 5.15,  $Y \perp_{\mathbb{P}_\alpha} X|W$  for all  $\alpha \in A$ , and so  $Y \perp_{\mathbb{P}_\{\}} X|W$ . Only if: For  $\mathbb{P}_\alpha$  dominating  $\mathbb{P}_\{\}$ , by Theorem 5.15, there exists a version of  $\mathbb{P}_\alpha^{Y|WX}$  unresponsive to  $W$ . Because  $\mathbb{P}_\alpha$  dominates  $\mathbb{P}_\{\}$ ,  $\mathbb{P}_\alpha^{Y|WX}$  differs from  $\mathbb{P}_\beta^{Y|WX}$  on a set of measure 0 for any  $\mathbb{P}_\beta \in \mathbb{P}_\{\}$ , thus  $\mathbb{P}_\alpha^{Y|WX}$  is a version of  $\mathbb{P}_\{\}^{Y|WX}$  also.  $\square$

**Corollary 5.22.** *Given standard measurable  $\Omega$ , variables  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and a probability set  $\mathbb{P}_\{\} : A \rightarrow \Delta(\Omega)$  with conditional probability  $\mathbb{P}_\{\}^{Y|WX}$ ,  $\mathbb{P}_\{\}^{Y|W}$  exists if  $Y \perp_{\mathbb{P}_\{\}} X|W$ .*

*Proof.* By Theorem 5.21, there is  $\mathbb{K} : W \rightarrow Y$  such that for all  $\alpha$

$$\mathbb{P}_\alpha^{WY} = \begin{array}{c} \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_\alpha^{WX} \text{ has two outputs. One output goes to a box labeled } \mathbb{P}_{\{\}}^{Y|WX}. \text{ The other output goes to a box labeled } \mathbb{K}. \text{ The box } \mathbb{P}_{\{\}}^{Y|WX} \text{ has an input from } W \text{ and an output to } Y. \text{ The box } \mathbb{K} \text{ has an input from } W \text{ and an output to } Y. \end{array} \\ (174) \end{array}$$

$$= \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_\alpha^{WX} \text{ has two outputs. One output goes to a box labeled } \mathbb{K}. \text{ The other output goes to a box labeled } \mathbb{K}. \text{ The box } \mathbb{K} \text{ has an input from } W \text{ and an output to } Y. \end{array} \quad (175)$$

$$= \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_\alpha^{W} \text{ has one output going to a box labeled } \mathbb{K}. \text{ The box } \mathbb{K} \text{ has an input from } W \text{ and an output to } Y. \end{array} \quad (176)$$

Thus  $\mathbb{K}$  is a version of  $\mathbb{P}_{\{\}}^{Y|W}$ .  $\square$

This result can fail to hold in the absence of the domination condition. Consider  $A$  a collection of inserts that all deterministically set a variable  $X$ ; then for any variable  $Y$   $Y \perp\!\!\!\perp_{\mathbb{P}_\square} X$  because  $X$  is deterministic for any  $\alpha \in A$ . But  $\mathbb{P}_{\square}^{Y|X}$  is not necessarily unresponsive to  $X$ .

Note that in the absence of the assumption of the existence of  $\mathbb{P}_{\square}^{Y|WX}$ ,  $Y \perp\!\!\!\perp_{\mathbb{P}_\square} X|W$  does *not* imply the existence of  $\mathbb{P}_{\square}^{Y|W}$ . If we have, for example,  $A = \{\alpha, \beta\}$  and  $\mathbb{P}_\alpha^{XY}$  is two flips of a fair coin while  $\mathbb{P}_\beta^{XY}$  is two flips of a biased coin, then  $Y \perp\!\!\!\perp_{\mathbb{P}} X$  but  $\mathbb{P}^Y$  does not exist.

## References

- The Basic Symmetries. In Olav Kallenberg, editor, *Probabilistic Symmetries and Invariance Principles*, Probability and Its Applications, pages 24–68. Springer, New York, NY, 2005. ISBN 978-0-387-28861-1. doi: 10.1007/0-387-28861-9\_2. URL [https://doi.org/10.1007/0-387-28861-9\\_2](https://doi.org/10.1007/0-387-28861-9_2).
- A. V. Banerjee, S. Chassang, and E. Snowberg. Chapter 4 - Decision Theoretic Approaches to Experiment Design and External Validity. We thank Esther Duflo for her leadership on the handbook and for extensive comments on earlier drafts. Chassang and Snowberg gratefully acknowledge the support of NSF grant SES-1156154. In Abhijit Vinayak Banerjee and Esther Duflo, editors, *Handbook of Economic Field Experiments*, volume 1 of *Handbook of Field Experiments*, pages 141–174. North-Holland, January 2017. doi: 10.1016/bs.hefe.2016.08.005. URL <https://www.sciencedirect.com/science/article/pii/S2214658X16300071>.



- Abhijit V. Banerjee, James Berry, Esther Duflo, Harini Kannan, and Shobhini Mukerji. Mainstreaming an Effective Intervention: Evidence from Randomized Evaluations of 'Teaching at the Right Level' in India. SSRN Scholarly Paper ID 2843569, Social Science Research Network, Rochester, NY, September 2016. URL <https://papers.ssrn.com/abstract=2843569>.
- Vladimir Bogachev and Ilya Malofeev. Kantorovich problems and conditional measures depending on a parameter. *Journal of Mathematical Analysis and Applications*, 486:123883, June 2020. doi: 10.1016/j.jmaa.2020.123883.
- Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- 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.
- Panayiota Constantinou and A. Philip Dawid. EXTENDED CONDITIONAL INDEPENDENCE AND APPLICATIONS IN CAUSAL INFERENCE. *The Annals of Statistics*, 45(6):2618–2653, 2017. ISSN 0090-5364. URL <http://www.jstor.org/stable/26362953>. Publisher: Institute of Mathematical Statistics.
- A. P. Dawid. Causal Inference without Counterfactuals. *Journal of the American Statistical Association*, 95(450):407–424, June 2000. ISSN 0162-1459. doi: 10.1080/01621459.2000.10474210. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.2000.10474210>. Publisher: Taylor & Francis \_eprint: <https://www.tandfonline.com/doi/pdf/10.1080/01621459.2000.10474210>.
- 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.
- Bruno de Finetti. Foresight: Its Logical Laws, Its Subjective Sources. In Samuel Kotz and Norman L. Johnson, editors, *Breakthroughs in Statistics: Foundations and Basic Theory*, Springer Series in Statistics, pages 134–174. Springer, New York, NY, [1937] 1992. ISBN 978-1-4612-0919-5. doi: 10.1007/978-1-4612-0919-5\_10. URL [https://doi.org/10.1007/978-1-4612-0919-5\\_10](https://doi.org/10.1007/978-1-4612-0919-5_10).
- 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.

- 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>.
- SANDER GREENLAND and JAMES M ROBINS. Identifiability, Exchangeability, and Epidemiological Confounding. *International Journal of Epidemiology*, 15(3):413–419, September 1986. ISSN 0300-5771. doi: 10.1093/ije/15.3.413. URL <https://doi.org/10.1093/ije/15.3.413>.
- 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>.
- KyleW (<https://math.stackexchange.com/users/252356/kylew>). Distribution of universal quantifiers over implication. Mathematics Stack Exchange. URL <https://math.stackexchange.com/q/1377555>. URL:<https://math.stackexchange.com/q/1377555> (version: 2015-07-29).
- 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 Bojaczuk 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.
- 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.
- 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>.

- Thomas S Richardson and James M Robins. Single world intervention graphs (swigs): A unification of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social Sciences, University of Washington Series. Working Paper*, 128(30):2013, 2013.
- 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>.
- Leonard J. Savage. *The Foundations of Statistics*. Wiley Publications in Statistics, 1954.
- Eyal Shahar. 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.
- Peter Walley. *Statistical Reasoning with Imprecise Probabilities*. Chapman & Hall, 1991.

## Appendix: