

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

David Johnston

January 11, 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Variables and Probability Models</b>	<b>3</b>
2.1	Section outline . . . . .	3
2.1.1	Brief outline of probability gap models . . . . .	4
2.2	Semantics of observed and unobserved variables . . . . .	5
2.3	Events . . . . .	8
2.4	Standard probability theory . . . . .	8
2.5	Probabilistic models for causal inference . . . . .	10
2.6	Probability sets . . . . .	11
2.7	Probability sets defined by marginal and conditional probabilities	14
2.8	Probability gap models . . . . .	15
2.9	Example: invalidity . . . . .	16
2.9.1	Conditional independence . . . . .	16
2.10	Curried Markov kernels . . . . .	17
<b>3</b>	<b>Decision theoretic causal inference</b>	<b>19</b>
3.1	Decision problems . . . . .	19
3.2	Decisions as measurement procedures . . . . .	21
3.3	Causal models similar to see-do models . . . . .	21
3.4	See-do models and classical statistics . . . . .	22
<b>4</b>	<b>Conditional probabilities in sequential experiments</b>	<b>24</b>
4.1	Repeatable experiments . . . . .	24
4.2	When does a canonical “effect of a decision” exist? . . . . .	27
4.3	In decision problems, policies are extreme . . . . .	30
4.3.1	Extended conditional independence . . . . .	33
4.3.2	Graphical properties of conditional independence . . . . .	36
4.4	Results I use that don’t really fit into the flow of the text . . . . .	36
4.4.1	Repeated variables . . . . .	36
4.5	Validity . . . . .	39

4.6	Combs . . . . .	41
4.7	Comb conditional correspondence . . . . .	43
4.8	Representation of conditional probability models . . . . .	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).

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

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

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

---

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

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

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

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

## 2 Variables and Probability Models

### 2.1 Section outline

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  $\{\mathbb{P}^{X|W}, \mathbb{P}^{Z|WXY}\}$ , 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 3, 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  $\{\mathbb{P}^{X|W}, \mathbb{P}^{Z|WXY}\}$  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.

## 2.2 Semantics of observed and unobserved variables

We are interested in constructing *probabilistic models* which explain some part of the world. In a model, variables play the role of “pointing to the parts of the world the model is explaining”. Both observed and unobserved variables play important roles in causal modelling and we think it is worth clarifying what variables of either type refer to. Our approach is a standard one: a probabilistic model is associated with an experiment or measurement procedure that yields values in a well-defined set. Observable variables are obtained by applying well-defined functions to the result of this total measurement. We use a richer fundamental probability set that includes “unobserved variables” that are formally treated the same way as observed variables, but aren’t associated with any real-world counterparts.

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

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

No actual procedure can be guaranteed to return elements of a mathematical set known in advance – anything can fail – but we assume that we can study procedures reliable enough that we don’t lose much by making this assumption.

**Definition 2.1** (Measurement procedure). A *measurement procedure* is a procedure that involves interacting with the real world somehow and delivering an element of a mathematical set as a result. The set of possible values is known prior to the measurement taking place, but the value that it will yield is not known. A procedure is given the font  $\mathcal{B}$ , we say it takes values in  $X$  and  $\mathcal{B} \bowtie x$  is the proposition that the the procedure  $\mathcal{B}$  will yield the value  $x \in X$ .  $\mathcal{B} \bowtie A$  for  $A \subset X$  is the proposition  $\bigvee_{x \in A} \mathcal{B} \bowtie x$ . Two procedures  $\mathcal{B}$  and  $\mathcal{C}$  are the same if  $\mathcal{B} \bowtie x \iff \mathcal{C} \bowtie x$  for all  $x \in B$  (note that  $\mathcal{B}$  and  $\mathcal{C}$  could involve different actions in the real world).

Measurement procedures are like functions without well-defined domains. We can compose measurement procedures with functions to produce new measurement procedures.

**Definition 2.2** (Composition of functions with procedures). Given a procedure  $\mathcal{B}$  that takes values in some set  $B$ , and a function  $f : B \rightarrow C$ , define the “composition”  $f \circ \mathcal{B}$  to be any procedure  $\mathcal{C}$  that yields  $f(x)$  whenever  $\mathcal{B}$  yields  $x$ . We can construct such a procedure by describing the steps: first, do  $\mathcal{B}$  and secondly, apply  $f$  to the value yielded by  $\mathcal{B}$ .

For example,  $\mathcal{MA}$  is the composition of  $h : (x, y) \mapsto xy$  with the procedure  $(\mathcal{M}, \mathcal{A})$  that yields the mass and acceleration of the same object. Measurement procedure composition is associative:

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

$$\iff \mathcal{B} \text{ yields } f^{-1}(g^{-1}(x)) \quad (2)$$

$$\iff f \circ \mathcal{B} \text{ yields } g^{-1}(x) \quad (3)$$

$$\iff g \circ (f \circ \mathcal{B}) \text{ yields } x \quad (4)$$

One might wonder whether there is also some kind of “append” operation that takes a standalone  $\mathcal{M}$  and a standalone  $\mathcal{A}$  and returns a procedure  $(\mathcal{M}, \mathcal{A})$ . Unlike function composition, this would be an operation that acts on two procedures rather than a procedure and a function. Unlike composition, we can’t easily reason about such an operation mathematically, because of the fact that measurement procedures have a foot in the real world. Our approach here is to suppose that there is some master measurement procedure  $\mathcal{S}$  which takes values in  $\Psi$  that handles all of the “real world” interaction relevant to our problem. Specifically, we assume that any measurement procedure of interest to our problem can be written as the composition  $f \circ \mathcal{S}$  for some  $f$ .

For the model  $\mathcal{F} = \mathcal{MA}$ , for example, we could assume  $\mathcal{F} = f \circ \mathcal{S}$  for some  $f$  and  $(\mathcal{M}, \mathcal{A}) = g \circ \mathcal{S}$  for some  $g$ . In this case, we can get  $\mathcal{MA} = h \circ (\mathcal{M}, \mathcal{A}) =$

$(h \circ g) \circ \mathcal{S}$ . Note that each procedure is associated with a unique function with domain  $\Psi$ .

Given that measurement processes are in practice finite precision and with finite range,  $\Psi$  will generally be a finite set. We can therefore equip  $\Psi$  with the collection of measurable sets given by the power set  $\mathcal{E} := \mathcal{P}(\Psi)$ , and  $(\Psi, \mathcal{E})$  is a standard measurable space.  $\mathcal{E}$  stands for a complete collection of logical propositions we can generate that depend on the results yielded by the measurement procedure  $\mathcal{S}$ .

$(\Psi, \mathcal{E})$  defines is a “sample space” limited to observable variables. That is,  $(\Psi, \mathcal{E})$  is associated with a measurement procedure. Unobserved variables need not be associated with measurement procedures, and to accommodate these we use instead of  $(\Psi, \mathcal{E})$  a richer sample space  $(\Omega, \mathcal{F})$  which represents both observed and unobserved variables.

**Definition 2.3** (Sample space). The sample space  $(\Omega, \mathcal{F})$  is a set  $\Omega$  along with with a  $\sigma$ -algebra  $\mathcal{F}$  of subsets of  $\Omega$ .

Observables are represented by a function  $\mathcal{S} : \Omega \rightarrow \Psi$ , and values of  $\omega$  are related to propositions about measurement procedures via the criterion of *consistency with observation*.

**Definition 2.4** (Consistency with observation). An element  $\omega \in \Omega$  is *consistent with observation* if the result yielded by  $\mathcal{S} \bowtie \mathcal{S}(\omega)$

Thus the procedure  $\mathcal{S}$  restricts the observationally consistent elements of  $\Omega$ . If  $\mathcal{S}$  yield the result  $s$ , then the consistent values of  $\Omega$  will be  $\mathcal{S}^{-1}(s)$ . While two different sets of measurement outcomes  $\Psi$  and  $\Psi'$  entail a different measurement procedures  $\mathcal{S}$  and  $\mathcal{S}'$ , but different fundamental probability sets  $\Omega$  and  $\Omega'$  may be used to model a single procedure  $\mathcal{S}$ .

As far as we know, distinguishing variables from procedures is somewhat non-standard, but we feel it is useful to distinguish the formal elements of our theory (variables) from the semi-formal elements (measurement procedures). Both variables and procedures are often discussed in statistical texts. For example, Pearl (2009) offers the following two, purportedly equivalent, definitions of variables:

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

This is a minor generalization of the textbook definition, according to which a random variable is a mapping from the fundamental probability set (e.g., the set of elementary events) to the real line. In our definition, the mapping is from the fundamental probability set to any set of objects called “values,” which may or may not be ordered.

Our view is that the first definition is a definition of a procedure, while the second is a definition of a variable. Variables model procedures, but they are

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

### 2.3 Events

To recap, we have a procedure  $S$  yielding values in  $\Psi$  that measures everything we are interested in, a fundamental probability set  $\Omega$  and a function  $S$  that models  $S$  in the sense of Definition 2.4. We assume also that  $\Psi$  has a  $\sigma$ -algebra  $\mathcal{E}$  (this may be the power set of  $\Psi$ , as measurement procedures are typically limited to finite precision).  $\Omega$  is equipped with a  $\sigma$ -algebra  $\mathcal{F}$  such that  $\sigma(S) \subset \mathcal{F}$ . If a procedure  $\mathcal{X} = f \circ S$  then we define  $X : \Omega \rightarrow X$  by  $X := f \circ S$ .

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

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

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

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

### 2.4 Standard probability theory

**Definition 2.5** (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 2.6** (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 2.7** (Probability measures as Markov kernel). Given  $(X, \mathcal{X})$  and  $\mu \in \Delta(X)$ , the Markov kernel  $\mathbb{K} : \{\ast\} \rightarrow X$  given by  $\mathbb{K}(A|\ast) = \mu(A)$  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.

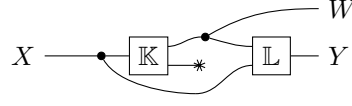


**Definition 2.8** (Delta measure). Given a measure space  $(X, \mathcal{X})$  and  $x \in X$ ,  $\delta_x \in \Delta(X)$  is the measure defined by  $\delta_x(A) = \mathbb{I}[x \in A]$ .

**Definition 2.9** (Probability space). A probability space is a triple  $(\mu, \Omega, \mathcal{F})$ , where  $\mu$  is a base measure on  $\mathcal{F}$ .

**Definition 2.10** (Marginal distribution with respect to a probability space). Given a probability space  $(\mu, \Omega, \mathcal{F})$  and a random variable  $\mathbf{X} : \Omega \rightarrow (X, \mathcal{X})$ , we can define the *marginal distribution* of  $\mathbf{X}$  with respect to  $\mu$ ,  $\mu^{\mathbf{X}} : \mathcal{X} \rightarrow [0, 1]$  by  $\mu^{\mathbf{X}}(A) := \mu(\mathbf{X} \boxtimes A)$  for any  $A \in \mathcal{X}$ .

**Definition 2.11** (Disintegration). Given a Markov kernel  $\mathbb{K} : W \rightarrow X \times Y$ , with  $W, X$  and  $Y$  standard measurable, any kernel  $\mathbb{L} : W \times X \rightarrow Y$  satisfying



$$\mathbb{K} = \quad (5)$$

is a  $W \times X \rightarrow Y$  *disintegration* of  $\mathbb{K}$ .

**Definition 2.12** (Conditional probability with respect to a probability space). Given a probability space  $(\mu, \Omega)$  and random variables  $\mathbf{X} : \Omega \rightarrow (X, \mathcal{X})$ ,  $\mathbf{Y} : \Omega \rightarrow (Y, \mathcal{Y})$ , the probability of  $\mathbf{Y}$  given  $\mathbf{X}$  is any  $X \rightarrow Y$  disintegration of  $\mu^{\mathbf{X}\mathbf{Y}}$ . That is,

$$\int_A \mu^{\mathbf{Y}|\mathbf{X}}(B|x) d\mu^{\mathbf{X}}(x) = \mu^{\mathbf{X}\mathbf{Y}}(x, y) \quad \forall A \in \mathcal{X}, B \in \mathcal{Y} \quad (6)$$

$$\iff \quad (7)$$

$$\quad (8)$$

**Lemma 2.13** (Marginal distribution as a kernel product). *Given a probability space  $(\mu, \Omega, \mathcal{F})$  and a random variable  $\mathbf{X} : \Omega \rightarrow (X, \mathcal{X})$ , define  $\mathbb{F}_{\mathbf{X}} : \Omega \rightarrow X$  by  $\mathbb{F}_{\mathbf{X}}(A|\omega) = \delta_{\mathbf{X}(\omega)}(A)$ , then*

$$\mu^{\mathbf{X}} = \mu \mathbb{F}_{\mathbf{X}} \quad (9)$$

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

$$\mu \mathbb{F}_{\mathbf{X}}(A) = \int_{\Omega} \delta_{\mathbf{X}(\omega)}(A) d\mu(\omega) \quad (10)$$

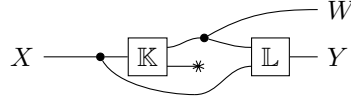
$$= \int_{\mathbf{X}^{-1}(A)} d\mu(\omega) \quad (11)$$

$$= \mu^{\mathbf{X}}(A) \quad (12)$$

□

Disintegration of arbitrary Markov kernels is possible in standard measurable spaces. We will assume that all spaces are standard measurable, such that whenever we have a Markov kernel we are able to disintegrate it.

**Lemma 2.14** (Disintegration existence in standard measurable Markov kernels). *For any Markov kernel  $\mathbb{K} : X \rightarrow W \times Y$  and  $X, W, Y$  are standard measurable, there exists  $\mathbb{L} : W \times X \rightarrow Y$  such that*



$$\mathbb{K} = \quad (13)$$

$\mathbb{L}$  is a disintegration of  $\mathbb{K}$ .

*Proof.* Cho and Jacobs (2019) Theorem 3.11 □

## 2.5 Probabilistic models for causal inference

The sample space  $(\Omega, \mathcal{F})$  along with our collection of variables is a “model skeleton” – it tells us what kind of data we might see. The process  $\mathcal{S}$  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 probability 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 correspond to a non-empty probability set.

## 2.6 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 2.15** (Probability set). A probability set  $\mathbb{P}_{\Omega}$  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  $\mathbb{P}_{\Omega}$ , we define marginal and conditional probabilities as probability measures and Markov kernels that satisfy Definitions 2.10 and 2.12 respectively for *all* base measures in  $\mathbb{P}_{\Omega}$ . 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 2.21) with respect to probability sets.

**Definition 2.16** (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  $\mathbb{P}_{\Omega}$ , the marginal distribution  $\mathbb{P}_{\Omega}^X = \mathbb{P}_{\alpha}^X$  for any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$  if a distribution satisfying this condition exists. Otherwise, it is undefined.

**Definition 2.17** (Conditional probability with respect to a probability set). Given a fundamental probability set  $\Omega$  variables  $\mathbf{X} : \Omega \rightarrow X$  and  $\mathbf{Y} : \Omega \rightarrow Y$  and a probability set  $\mathbb{P}_{\{\}}$ , a version of  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  is any Markov kernel  $X \rightarrow Y$  such that  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  is an  $X \rightarrow Y$  disintegration of  $\mathbb{P}_{\alpha}^{\mathbf{X}\mathbf{Y}}$  for all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$ . If no such Markov kernel exists,  $\mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}}$  is undefined.

Given a conditional probability with respect to a probability set, we can find other conditional probabilities by “pushing it forward”.

**Theorem 2.18** (Recursive pushforward). Suppose we have a sample space  $\Omega$  variables  $\mathbf{X} : \Omega \rightarrow X$  and  $\mathbf{Y} : \Omega \rightarrow Y$ ,  $\mathbf{Z} : \Omega \rightarrow Z$  and a probability set  $\mathbb{P}_{\{\}}$  such that  $\mathbb{P}_{\{\}}^{\mathbf{X}|\mathbf{Y}}$  is a  $\mathbf{Y}|\mathbf{X}$  conditional probability of  $\mathbb{P}_{\{\}}$  and  $\mathbf{Z} = f \circ \mathbf{Y}$  for some  $f : Y \rightarrow Z$ . Then there exists a  $\mathbf{Z}|\mathbf{X}$  conditional probability of  $\mathbb{P}_{\{\}}$  given by  $\mathbb{P}_{\{\}}^{\mathbf{Z}|\mathbf{X}} = \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \mathbb{F}_f$ .

*Proof.* For any  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$ ,  $x, z$

$$\mathbb{P}_{\alpha}^{\mathbf{X}}(x) \mathbb{P}_{\alpha}^{\mathbf{Z}|\mathbf{X}}(z|x) = \mathbb{P}_{\alpha}(\mathbf{X}^{-1}(x) \cap \mathbf{Z}^{-1}(z)) \quad (14)$$

$$= \mathbb{P}_{\alpha}(\mathbf{X}^{-1}(x) \cap \mathbf{Y}^{-1}(f^{-1}(z))) \quad (15)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}, \mathbf{Y}}(\{x\} \times f^{-1}(z)) \quad (16)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}}(x) \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{X}}(f^{-1}(z)|x) \quad (17)$$

□

We define the copy-product  $\odot$  as a shorthand for the operation in Equation ?? that combines a marginal with a disintegration to get the original Markov kernel back.

**Definition 2.19** (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

$$\mathbb{K} \odot \mathbb{L} := \text{copy}_X(\mathbb{K} \otimes \text{id}_X)(\text{copy}_Y \otimes \text{id}_X)(\text{id}_Y \otimes \mathbb{L}) \quad (18)$$

$$= \begin{array}{c} \text{---} Y \\ \text{---} X \text{---} \boxed{\mathbb{K}} \text{---} \bullet \text{---} \boxed{\mathbb{L}} \text{---} Z \\ \text{---} \end{array} \quad (19)$$

$$\iff \quad (20)$$

$$(\mathbb{K} \odot \mathbb{L})(A \times B|x) = \int_A \mathbb{L}(B|y, x) \mathbb{K}(dy|x) \quad A \in \mathcal{Y}, B \in \mathcal{Z} \quad (21)$$

**Lemma 2.20** (Semidirect product is associative). Given  $\mathbb{K} : X \rightarrow Y$ ,  $\mathbb{L} : Y \times X \rightarrow Z$  and  $\mathbb{M} : Z \times Y \times X \rightarrow W$

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (22)$$

$$(23)$$

*Proof.*

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \boxed{\text{od}(\mathbf{a4})} = \boxed{\text{od}(\mathbf{a5})} = \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (26)$$

□

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  $\mathbf{X}$  with each Markov kernel is identical.

**Definition 2.21** (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  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$ ,

$$\mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{K} = \mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{L} \quad (27)$$

**Lemma 2.22** (Conditional probabilities are almost surely equal). *If  $\mathbb{K} : X \rightarrow Y$  and  $\mathbb{L} : X \rightarrow Y$  are both versions of  $\mathbb{P}_{\Omega}^{Y|X}$ ,  $\mathbb{K} \stackrel{\mathbb{P}_{\Omega}}{\cong} \mathbb{L}$*

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

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

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

□

**Lemma 2.23** (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  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\Omega}$*

$$\mathbb{P}_{\alpha}^{Y|X} \odot \mathbb{K} \stackrel{a.s.}{\cong} \mathbb{P}_{\alpha}^{Y|X} \odot \mathbb{L} \quad (30)$$

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

$$\mathbb{P}_{\alpha}^{\mathbf{XY}} \odot \mathbb{K} = (\mathbb{P}_{\alpha}^{\mathbf{X}} \odot \mathbb{P}_{\Omega}^{Y|X}) \odot \mathbb{K} \quad (31)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\Omega}^{Y|X} \odot \mathbb{K}) \quad (32)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}} \odot (\mathbb{P}_{\Omega}^{Y|X} \odot \mathbb{L}) \quad (33)$$

□

**Lemma 2.24** (Semidirect product of conditionals is a joint conditional). *Given a probability set  $\mathbb{P}_{\Omega}$  on  $(\Omega, \mathcal{F})$  along with conditional probabilities  $\mathbb{P}_{\Omega}^{Y|X}$  and  $\mathbb{P}_{\Omega}^{Z|XY}$ ,  $\mathbb{P}_{\Omega}^{YZ|X}$  exists and is equal to*

$$\mathbb{P}_{\Omega}^{YZ|X} = \mathbb{P}_{\Omega}^{Y|X} \odot \mathbb{P}_{\Omega}^{Z|XY} \quad (34)$$

$$(35)$$

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

$$\mathbb{P}_\alpha^{\mathbf{XYZ}} = \mathbb{P}_\alpha^{\mathbf{X}} \odot \mathbb{P}_\alpha^{\mathbf{YZ|X}} \quad (36)$$

$$= \mathbb{P}_\alpha^{\mathbf{X}} \odot (\mathbb{P}_\alpha^{\mathbf{Y|X}} \odot \mathbb{P}_\alpha^{\mathbf{Z|YX}}) \quad (37)$$

$$= \mathbb{P}_\alpha^{\mathbf{X}} \odot (\mathbb{P}_{\{\}}^{\mathbf{Y|X}} \odot \mathbb{P}_{\{\}}^{\mathbf{Z|YX}}) \quad (38)$$

□

## 2.7 Probability sets defined by marginal and conditional probabilities

In the previous section we defined marginal and conditional probabilities for probability sets. Here we will go in the other direction: define probability sets by specifying key marginal or conditional probabilities. There is an issue to be careful of here: not all probability measures  $\mathbb{Q}^{\mathbf{X}}$  on  $X$  define nonempty sets of probability measures on  $\Omega$  with respect to the variable  $\mathbf{X}$ . Consider, for example,  $\mathbf{X} = (\mathbf{Y}, \mathbf{Y})$  for some  $\mathbf{Y} : \Omega \rightarrow \{0, 1\}$  and any measure  $\mathbb{Q}^{\mathbf{YY}}$  that assigns nonzero probability to the event  $(\mathbf{Y}, \mathbf{Y}) \bowtie (1, 0)$ . There is no base measure that pushes forward to such a  $\mathbb{P}^{\mathbf{YY}}$ , because two copies of the same variable must always be deterministically equal. A *valid distribution* is a distribution associated with a particular variable that defines a nonempty set of base measures on  $\Omega$  (Theorem 4.16).

**Definition 2.25** (Valid distribution). A valid  $\mathbf{X}$  probability distribution  $\mathbb{P}^{\mathbf{X}}$  is any probability measure on  $\Delta(X)$  such that  $\mathbf{X}^{-1}(A) = \emptyset \implies \mathbb{P}^{\mathbf{X}}(A) = 0$  for all  $A \in \mathcal{X}$ .

*Valid conditionals* not only define a nonempty set of base measures on  $\Omega$ , but the  $\odot$  product of two appropriately typed valid conditionals itself defines a nonempty set of base measures on  $\Omega$  (Lemma 4.19).

**Definition 2.26** (Valid conditional). Given  $(\Omega, \mathcal{F})$ ,  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbf{Y} : \Omega \rightarrow Y$  a *valid  $\mathbf{Y|X}$  conditional probability*  $\mathbb{P}^{\mathbf{Y|X}}$  is a Markov kernel  $X \rightarrow Y$  such that it assigns probability 0 to contradictions:

$$\forall B \in \mathcal{Y}, x \in \mathcal{X} : (\mathbf{X}, \mathbf{Y}) \bowtie \{x\} \times B = \emptyset \implies \left( \mathbb{P}^{\mathbf{Y|X}}(A|x) = 0 \right) \vee (\mathbf{X} \bowtie \{x\} = \emptyset) \quad (39)$$

Thus, given a collection of valid conditional probabilities  $\{\mathbb{P}_i^{\mathbf{X}_i|\mathbf{X}_{[i-1]}} | i \in [n]\}$  such that each adjacent pair can be combined with the  $\odot$  product, the sequential product of each conditional probability is a valid conditional probability and there is a non-empty set of probability measures on the sample space that with that conditional probability. 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). Lemma 4.19 establishes that recursive collections of conditional

probabilities define non-empty probability sets as long as all the conditional probabilities in the collection are valid.

**Definition 2.27** (Probability set defined by a valid conditional). If  $\mathbb{P}_{\{\}} is a probability set such that there is a valid conditional probability  $\mathbb{P}_{\{\}}^{Y|X} : X \rightarrow Y$  and for every  $\mu \in \Delta(\Omega)$  such that  $\mu^{Y|X} \stackrel{\mu}{\cong} \mathbb{P}_{\{\}}^{Y|X}$ , we say  $\mathbb{P}_{\{\}}^{\bar{Y}|\bar{X}} := \mathbb{P}_{\{\}}$  is the probability set defined by  $\mathbb{P}_{\{\}}^{Y|X}$ .$

## 2.8 Probability gap models

A probability set gives us a collection of different probability models that “could” model the consequences of an action. In addition to possibilities, we also want *choices* and a model that tells us, given a particular choice, we realise a particular possibility. This is what we call a *probability gap model*. The general form of a probability gap models is

- A fixed probability set  $\mathbb{P}_{\{\}} \subset \Delta(\Omega)$  which we call the *model*
- A collection of probability sets  $A \subset \Delta(\Omega)$  that we call *choices*
- A map  $\mathbb{P}_{\square} : A \rightarrow \mathcal{P}(\Delta(\Omega))$  defined by  $\mathbb{P}_{\alpha} := \mathbb{P}_{\square}(\alpha) = \mathbb{P}_{\{\}} \cap \alpha$

We require that the choices are compatible with the model in the sense that  $\mathbb{P}_{\{\}} \cap \alpha \neq \emptyset$  for all  $\alpha \in A$ . Here, we will limit our attention to a particular type of probability gap model, where we define the probability set  $\mathbb{P}_{\{\}}$  is defined by a conditional probability and each choice is defined by a marginal probability relative to the same variable.

**Definition 2.28** (Conditional probability model). A *conditional probability model*  $\mathbb{P}_{\square}$  is a probability gap model  $(\mathbb{P}_{\{\}}^{\bar{Y}|\bar{X}}, A)$  such that each  $\alpha \in A$  is some probability set defined by an  $X$ -valid marginal probability  $\alpha^{\bar{X}}$ .

We will compute the intersection  $\mathbb{P}_{\alpha}$  between the model  $\mathbb{P}_{\{\}}$  and a choice  $\alpha \in A$  as the probability set  $\mathbb{P}_{\alpha}^{\bar{X}|\bar{Y}}$  such that:

$$\mathbb{P}_{\alpha}^{\bar{X}|\bar{Y}} = \alpha^{\bar{X}} \odot \mathbb{P}_{\{\}}^{\bar{Y}|\bar{X}} \quad (40)$$

This is justified by Lemma 4.17, which says that the probability set defined by Equation 40 is equivalent to the intersection of  $\alpha$  and  $\mathbb{P}_{\{\}}$ .

If the conditional probability  $\mathbb{P}_{\{\}}^{\bar{Y}|\bar{X}}$  and all the marginal probabilities  $\alpha^{\bar{X}}$  are valid, then by Lemma 4.19  $\mathbb{P}_{\{\}} \cap \alpha \neq \emptyset$  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.

## 2.9 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

$$\nu^{\mathcal{B}|\mathcal{WH}} = \begin{array}{c} \mathcal{H} \xrightarrow{*} \\ \mathcal{W} \xrightarrow{*} \end{array} \triangleleft_{\delta_x} \text{---} \mathcal{B} \quad (41)$$

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

$$\begin{aligned} (\mathcal{B}, \mathcal{W}, \mathcal{H}) \bowtie \{(x, w, h)\} &= \{\omega | (\mathcal{W}, \mathcal{H})(\omega) = (w, h), \mathcal{B}(\omega) = \frac{w}{h^2}\} \\ &= \emptyset \end{aligned} \quad (42)$$

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

Validity rules out conditional probabilities like 41. 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 41 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  $x \mapsto P(\cdot | do(\mathcal{B} = x))$  as is the usual convention when discussing causal Bayesian networks.

### 2.9.1 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 4.3.1.

**Definition 2.29** (Conditional independence with respect to a probability set). For a probability set  $\mathbb{P}_{\{\}}$  and variables  $\mathcal{A}, \mathcal{B}, \mathcal{Z}$ , we say  $\mathcal{B}$  is conditionally inde-



pendent of  $A$  given  $C$ , written  $B \perp\!\!\!\perp_{\mathbb{P}_{\{\}} } A|C$ , if

$$\mathbb{P}_{\{\}}^{ABC} = \begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}^Z \text{ with a dot on its right side. From the dot, three lines emerge: one to a box labeled } \mathbb{P}^{A|C} \text{ leading to } A, \text{ one to a box labeled } \mathbb{P}^{B|C} \text{ leading to } B, \text{ and one to } C. \end{array} \quad (44)$$

Cho and Jacobs (2019) have shown that this definition coincides with the standard notion of conditional independence for a particular probability model. In particular, it satisfies the *semi-graphoid axioms*.

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  $A \perp\!\!\!\perp_{\mathbb{P}_{\square}} (B, C)|W$

**Theorem 2.30.** *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*

$$\mathbb{P}^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{P}_{\square}^{Y|W}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (45)$$

$$\iff \mathbb{P}^{Y|WX}(y|w, x) = \mathbb{P}^{Y|W}(y|w) \quad (46)$$

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

The semi-graphoid axioms hold for all probability measures  $\mathbb{P}$ , so in particular they hold for all  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\{\}}$ . Thus conditional independence with respect to a probability set also satisfies the semi-graphoid axioms.

**Definition 2.31** (Conditional independence with respect to a probability comb). Conditional independence  $A \perp\!\!\!\perp_{\mathbb{P}_{\square}} B|C$  holds for an arbitrary probability comb  $\mathbb{P}_{\square} : A \rightarrow \mathcal{P}(\Delta(\Omega))$  if  $A \perp\!\!\!\perp_{\mathbb{P}_{\alpha}} B|C$  holds for all probability models  $\mathbb{P}_{\alpha}$ ,  $\alpha \in A$ .

## 2.10 Curried Markov kernels

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  $\lambda \mathbb{K} : \{*\} \rightarrow Y^X$  by

$$\lambda \mathbb{K}((y_i)_{i \in X}) = \prod_{i \in X} \mathbb{K}(y_i|i) \quad (47)$$

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

$$\mathbb{K} = (\lambda \mathbb{K} \otimes \text{id}_X) \mathbb{F}_{\text{ev}} \quad (48)$$

Unlike the case of function currying,  $\lambda \mathbb{K}$  is not the unique Markov kernel for which 48 holds. In fact, we can substitute any  $\mathbb{L}$  such that, for any  $i \in X$

$$\sum_{y_{\{i\}^C} \in Y^{|X|-1}} \mathbb{L}((y_i)_{i \in X}) = \mathbb{K}(y_i | i) \quad (49)$$

Evaluation of a curried Markov kernel  $\lambda \mathbb{K}$  resembles the definition of *potential outcomes*; for outcomes  $Y : \Omega \rightarrow Y$  and treatments  $X : \Omega \rightarrow X$ , potential outcomes are described by a probability distribution  $\mathbb{P}^{Y^X}$  on  $Y^X$  and we have the relation

$$Y \stackrel{a.s.}{=} \text{ev}(Y^X, X) \quad (50)$$

Then

$$(\mathbb{P}^{Y^X} \otimes \text{id}_X) \mathbb{F}_{\text{ev}} \quad (51)$$

is some Markov kernel  $\mathbb{K} : X \rightarrow Y$ , which is equal to  $\mathbb{P}^{Y|X}$  if  $Y^X \perp\!\!\!\perp X$ . However, potential outcomes models typically do not explain what the kernel  $\mathbb{K}$  represents, and instead offer a definition of the variable  $Y^X$ . For  $x \in X$ , the component  $Y^x$  of  $Y^X$  is usually said to express “the outcomes that would have been observed, if  $X$  was  $x$ ”.

Our original motivating question was “when are potential outcomes well-defined?”. We’re not actually going to try to answer this question, because our aim is not to tell people using potential outcomes how to do it. Furthermore, that question invites controversy we are not particularly interested in joining; Dawid (2000) and Richardson and Robins (2013) have both argued that it is better to use equivalence classes of potential outcomes models induced by a criterion of distinguishability by experiment, while Pearl (2009) advocates for models that can make finer distinctions than this.

However, given a probability gap model  $\mathbb{P}_\square$ , we do have a natural notion of the well-definedness of a conditional probability  $\mathbb{P}_\square^{Y|X}$  – it is well-defined when  $\mathbb{P}_\alpha^{Y|X}$  is equal for all  $\alpha$  (Definition 2.17). Furthermore, the formal conditions that guarantee the existence of such a conditional probability very closely resemble the *stable unit treatment value assumption* (SUTVA), which is said to be necessary for the existence of potential outcomes 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.

The added emphasis is ours. In the next section, we offer formal criteria that correspond to these two statements.

start again here

### 3 Decision theoretic causal inference

People very often have to make decisions with some information they may consult to help them make the decision. We are going to examine how gappy probability models can formally represent problems of this type, which in turn allows us to make use of the theory of probability to help guide us to a good decision. Probabilistic models have a long history of being used to represent decision problems, and there exist a number of coherence theorems that show that preferences that satisfy certain kinds of constraints must admit representation by a probability model and a utility function of the appropriate type. Particularly noteworthy are the theorems of Ramsey (2016) and Savage (1954), which together yield a method for representing decision problems known as “Savage decision theory”, and the theorem of Bolker (1966); Jeffrey (1965) which yields a rather different method for representing decision problems known as “evidential decision theory”. Joyce (1999) extends Jeffrey and Bolker’s result to a representation theorem that subsumes both “causal decision theory” and “evidential decision theory”.

It is an open question whether the models induced by any of these theories are equivalent to probability gap models.

We do not have a comparable axiomatisation of preferences that yield a representation of decision problems in terms of utility and gappy probability. Such an undertaking could potentially clarify some choices that can be made in setting up a gappy probability model of decision making, but it is the subject of future work. Instead, we suppose that we are satisfied with a particular probabilistic model of a decision problem, based on convention rather than axiomatisation.

#### 3.1 Decision problems

Suppose we have an observation process  $\mathcal{X}$ , modelled by  $X$  taking values in  $X$  (we are *informed*). Given an observation  $x \in X$ , we suppose that we can choose a decision from a known set  $D$  (the set of decisions is *transparent*), and we suppose that choosing a decision results in some action being taken in the real world. As with processes of observation, we will mostly ignore the details of what “taking an action” involves. The process of choosing a decision that yields an element of  $D$  is a decision making process  $\mathcal{D}$  modelled by  $D$ . We might be able to introduce randomness to the choice, in which case the relation between  $X$  and  $D$  may be stochastic. We will assume that there is some  $\mathcal{Y}$  modelled by

$Y$  such that  $(X, D, Y)$  tell us everything we want to know for the purposes of deciding which outcomes are better than others.

We want a model that allows us to compare different stochastic *decision functions*  $Q_\alpha^{D|X} : X \rightarrow D$ , letting  $A$  be the set of all such functions available to be chosen. That is, we need a higher order function  $f$  that takes a decision function  $Q_\alpha^{X|D}$  and returns a probabilistic model of the consequences of selecting that decision function  $P_\alpha^{DXY}$ . An order 2 model  $(P_\square^{X, P_\square^Y | XD}, A)$  defines such a function, though there are many such functions that are not order 2 models. The key feature of probability gap models is that the map is by intersection of probability sets, so for example the conditional probability of  $X|D$  given a decision function  $Q_\alpha^{X|D}$  must actually be equal to  $Q_\alpha^{X|D}$ , and we can say the same for  $P_\square^X$  and  $P_\square^{Y|XD}$ . If we don't think all of these conditional probabilities are fixed, then we want something other than an order 2 model of the type discussed. We will define *ordinary decision problems* to be those for which the desired model  $P_\square$  is this type of order 2 probability gap model.

I think adding hypotheses at this point might make things unnecessarily confusing; on the other hand, they are useful for the connection to classical statistical decision theory. The "repeatable experiments" section shows how see-do models with certain assumptions induce an easier to understand class of hypotheses, and I could just save the idea of a hypothesis until I get there

We consider an additional kind of gap in our probability model. The nature of this gap is: we don't know exactly which order 2 model  $(P_\square^{X, P_\square^Y | XD}, A)$  we "ought" to use. To represent this gap we include an unobserved variable  $H$ , the *hypothesis*. We can interpret  $H$  as expressing the fact that, if we knew the value of  $H$  then we would know that our decision problem was represented by a unique order 2 model  $(P_{h, \square}^{X, P_{h, \square}^Y | XD}, A)$ . However,  $H$  is not known and in fact we do not know how to determine  $H$  (this is the nature of an *unobserved* variable – there is no process available to find the value it yields). Our model is thus given by

$$(P_\square^{X|H, P_\square^Y | HXD}, A)$$

**Definition 3.1** (Ordinary decision problem). An ordinary decision problem  $(\mathbb{P}, \Omega, H, (X, \mathcal{X}), (D, \mathcal{D}), (Y, \mathcal{Y}))$  consists of a fundamental probability set  $\Omega$ , hypotheses  $H : \Omega \rightarrow H$ , observations  $X : \Omega \rightarrow X$ , decisions  $D : \Omega \rightarrow D$  and consequences  $Y : \Omega \rightarrow Y$ , and the latter three random variables are associated with measurement processes. It is equipped with a probability gap model  $\mathbb{P} : \Delta(D)^X \rightarrow \Delta(\Omega)^H$  where  $\Delta(D)^X$  is the set of valid  $D|X$  Markov kernels  $X \rightarrow D$  and  $\Delta(\Omega)^H$  is the set of valid Markov kernels  $H \rightarrow \Omega$ . We require of  $\mathbb{P}$ :

1.  $\mathbb{P}_\alpha^{D|X} = Q_\alpha^{D|X}$  for all decision functions  $Q_\alpha^{D|X} \in \Delta(D)^X$
2.  $\mathbb{P}_\alpha^{X|H} = P_\alpha^{X|H}$  for all  $P_\alpha := \mathbb{P}(Q_\alpha^{D|X})$

3.  $\mathbb{P}^{Y|XDH} = \mathbb{P}_\alpha^{Y|XDH}$  for all  $\mathbb{P}_\alpha := \mathbb{P}(\mathbb{Q}_\alpha^{D|X})$

(1) reflects the assumption that the “probability of  $D$  given  $X$ ” based on the induced model is equal to the “probability of  $D$  given  $X$ ” based on the chosen decision function. (2) reflects the assumption that the observations should be modelled identically no matter which decision function is chosen. (3) reflects the assumption that given hypothesis, the observations and the decision, the model of  $Y$  does not depend any further on the decision function  $\alpha$ .

Under these assumptions  $\mathbb{P}_\square$  is an order 2 model  $(\mathbb{P}_\square^{X, \mathbb{P}_\square^Y | XD}, A)$  which we call a “see-do model”.

I need to update the proof for this claim

### 3.2 Decisions as measurement procedures

We have previously posited that observed variables are variables  $X$  – themselves purely mathematical objects – associated with a measurement process  $\mathcal{X}$  that has “one foot in the real world”. In the framework we have proposed here, decisions correspond to a special class of measurement procedure.

Suppose that we are only contemplating decision functions that map deterministically to  $D$ . Suppose furthermore that we will  $D$  according to a model  $\mathbb{P}_\square$ , a utility function on  $X \times D \times Y \rightarrow \mathbb{R}$  and a decision rule which is a function  $f$  from models, utility functions and decision rules to decisions. Note that models, utility functions and decision rules are all well-defined mathematical objects. If we are confident that our choice will in the end be an element of a well-defined set of objects of the appropriate type, then we are positing that we have a “measurement procedure”  $\mathcal{M}$  that yield models, utilities and decision rules. If so,  $f \circ \mathcal{M}$  – that is, the function that yields a decision – is itself a measurement procedure. This is what is unique about decisions: proposing a complete decision problem with models, utilities and decision rules, defines a measurement procedure for decisions. Other quantities of interest do not seem to have this property – we *require* a measurement process for observations in order to make the whole setup work, but we do not *define* it in the course of setting up a model for our decision problem.

I don’t know how important this observation is, but the fact that  $\mathcal{D}$  is an output of a formal decision making system makes it different from other things we might call decisions, and I wonder if I should call it something else in order to avoid ambiguity. The vague reason I think this matters is: whatever you might want to measure, you won’t learn more about  $\mathcal{D}$  from it than you already know once you have the model, the utility and the decision rule, this is not a property that other things we call “decisions” share and this distinction might be important regarding judgements of causal contractibility.

### 3.3 Causal models similar to see-do models

Lattimore and Rohde (2019a) and Lattimore and Rohde (2019b) consider an observational probability model and a collection of indexed interventional prob-

ability models, with the probability model tied to the interventional models by shared parameters. In these papers, they show how such a model can reproduce inferences made using Causal Bayesian Networks. This kind of model can be identified with a type of see-do model, where what we call hypotheses  $H$  are identified with the sequence of what Rohde and Lattimore call parameter variables.

The approach to decision theoretic causal inference described by Dawid (2020) is somewhat different:

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

The difference between the model described here and a see-do model is that a see-do model uses different variables  $X$  and  $Y$  to represent observations and consequences, while Dawid’s model uses the same variable  $(T, Y)$  to represent outcomes in interventional and observational regimes. In this work we associate one observed variable with each measurement process, while in Dawid’s approach  $(T, Y)$  seem to be doing double duty, representing measurement processes carried out during observations and after taking action. This can be thought of as the causal analogue of the difference between saying we have a sequence  $(X_1, X_2, X_3)$  of observations independent and identically distributed according to  $\mu \in \Delta(X)$  and saying that we have some observations distributed according to  $\mathbb{P}^X \in \Delta(X)$ . People usually understand what is meant by the latter, but if one is trying to be careful the former is a more precise statement of the model in question.

Heckerman and Shachter (1995) also explore a decision theoretic approach to causal inference. Our approach is quite close to their approach if we identify what we call hypotheses with what they call states and allow for probabilistic dependence between states, decisions and consequences. It is an open question whether their notion of limited unresponsiveness corresponds to any notion of conditional independence in our work.

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

### 3.4 See-do models and classical statistics

See-do models are capable of expressing the expected results of a particular choice of decision strategy, but they cannot by themselves tell us which strategies are more desirable than others. To do this, we need some measure of the desirability of our collection of results  $\{\mathbb{P}_\alpha | \alpha \in A\}$ . A common way to do this is

to employ the principle of expected utility. The classic result of Von Neumann and Morgenstern (1944) shows that all preferences over a collection of probability models that obey their axioms of completeness, transitivity, continuity and independence of irrelevant alternatives must be able to be expressed via the principle of expected utility. This does not imply that anyone knows what the appropriate utility function is.

A further property that may hold for some see-do models  $\mathbb{P}^{X|H \square Y|D}$  is  $Y \perp\!\!\!\perp_{\mathbb{P}}^2 X|(H, D)$ . This expresses the view that the consequences are independent of the observations, once the hypothesis and the decision are fixed. Such a situation could hold in our scenario above, where the observations are trial data, the decisions are recommendations to care providers and the consequences are future patient outcomes. In such a situation, we might suppose that the trial data are informative about the consequences only via some parameter such as effect size; if the effect size can be deduced from  $H$  then our assumption corresponds to the conditional independence above.

Given a see-do model  $\mathbb{P}^{X|H \square Y|D}$  along with the principle of expected utility to evaluate strategies, and the assumption  $Y \perp\!\!\!\perp_{\mathbb{P}}^2 X|(H, D)$  we obtain a statistical decision problem in the form introduced by Wald (1950).

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

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

$$\mathbb{E}_{\mathbb{P}_{\alpha}}[Y|h] = \sum_{x \in X} \sum_{d \in D} \sum_{y \in Y} \mathbb{P}^{X|H}(x|h) \mathbb{Q}_{\alpha}^{D|X}(d|x) \mathbb{P}^{Y|HD}(y|h, d) \quad (52)$$

$$= \sum_{x \in X} \sum_{d \in D} \mathbb{P}^{X|H}(x|h) \mathbb{Q}_{\alpha}^{D|X}(d|x) l(h, d) \quad (53)$$

$$= R(h, \alpha) \quad (54)$$

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

## 4 Conditional probabilities in sequential experiments

If we have a conditional probability model  $\mathbb{P}_{\square}^{\overline{Y|D}}$ , then by definition there is a conditional probability  $\mathbb{P}_{\square}^{Y|D}$  which has a curried representation. Our aim is to show when certain conditional probabilities exist with respect to a probability gap model, which in this case is a triviality.

The question becomes more interesting when we propose conditional probability model  $\mathbb{P}_{\square}^{\overline{Y|D}}$  of a sequential experiment. That is,  $Y := Y_M = (Y_i)_{i \in M}$  and  $D := D_M = (D_i)_{i \in M}$  and we say that  $Y_i$  is the consequence corresponding to the decision  $D_i$  for all  $i \in M$ . We will call a  $(D_i, Y_i)$  pair an experimental unit. In this case, the conditional probability  $\mathbb{P}_{\square}^{Y_i|D_i}$  does not generally exist. We might suppose, however, it would exist if the experiment was, in some sense, suitably regular or repeatable.

### 4.1 Repeatable experiments

We have a conditional probability model  $\mathbb{P}_{\square}^{\overline{Y_A|D_A}}$  with choices  $A$  that represents a sequential experiment. What might we mean when we say this experiment is repeatable? We’re going to propose two conditions. The first condition is *commutativity of exchange*, which is the assumption that swapping the choices that we apply at each step and then applying the corresponding inverse swap to consequences leaves the model unchanged. The second condition is *commutativity of marginalisation* – if we perform the whole experiment multiple times, making the same choice  $D_i$  at any point  $i$  gets the same results, regardless of what other choices are made.

Commutativity of exchange is similar to the condition of *post-treatment exchangeability* found in Dawid (2020), and commutativity of marginalisation is similar to the stable unit treatment distribution assumption (SUTDA) in the same, as well as the “no interference” part of the stable unit treatment value assumption (SUTVA) with which it shares a name. Commutativity of exchange is also very similar to the exchangeability assumption of GREENLAND and ROBINS (1986) for further discussions of exchangeability in the context of causal modelling, and note that both authors consider exchanging to be an operation that alters which person receives which treatment. The assumption of exchangeability found in Banerjee et al. (2017) can also be regarded as similar to commutativity of exchange.

I think the useful part is not that these ideas are conceptually new, but they have sharp definitions instead of

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

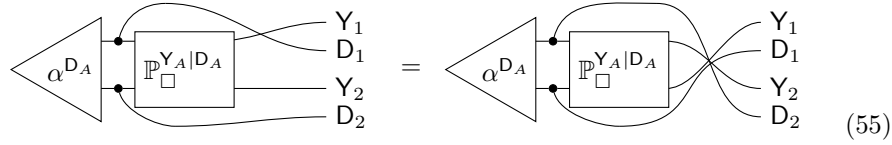
Commutativity of exchange is not equivalent to exchangeability in the sense of De Finetti’s well-known theorem de Finetti ([1937] 1992). The latter can be



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 experimental units receive 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):



Commutativity of exchange (swapping choices  $\sim$  swapping labels):

$$\boxed{\text{commutativity of exchange}} \quad (56)$$

Commutativity of exchange is a property of probability gap models, not a property of fixed probability model for which there is no analogue of “attaching a different choice” in that case.

—end not sure where to put—

More precisely, a conditional probability model “commutes with exchange” if applying any finite permutation to blind choices or separately applying the corresponding permutation to consequences each yields the same result. We can apply the exchange “before” multiplying by the conditional  $\mathbb{P}_{\square}^{Y|D}$  or after it and we get the same result.

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

**Definition 4.2** (Commutativity of exchange). Suppose we have a sample space  $(\Omega, \mathcal{F})$  and a conditional probability model  $(\mathbb{P}_{\square}^{Y|D}, A)$  with  $Y = Y_M$ ,  $D = D_M$ ,  $M \subseteq \mathbb{N}$ . If, for any two decision rules  $\alpha^{\bar{D}}, \beta^{\bar{D}} \in A$ ,

$$\alpha^{\bar{D}} \odot \text{swap}_{\rho(D)} \mathbb{P}_{\square}^{Y|D} = \alpha^{\bar{D}} \odot \mathbb{P}_{\square}^{Y|D} \text{swap}_{\rho(D \times Y)} \quad (57)$$

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

A do model is non interfering if it gives identical results for identical subsequences of different choices when we limit our attention to the corresponding subsequences of consequences. For example, if we have  $D = (D_1, D_2, D_3)$  and  $Y = (Y_1, Y_2, Y_3)$  and  $\alpha^{D_1 D_3} = \mathbb{P}_{\beta}^{D_1 D_3}$  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** (Commutativity of marginalisation). Suppose we have a sample space  $(\Omega, \mathcal{F})$  and a conditional probability model  $(\mathbb{P}_{\square}^{\mathbf{Y}|\mathbf{D}}, A)$  with  $\mathbf{Y} = \mathbf{Y}_M$ ,  $\mathbf{D} = \mathbf{D}_M$ ,  $M \subseteq \mathbb{N}$ . For any  $S = (s_i)_{i \in Q}$ ,  $Q \subset M$ , and  $i < j \implies p_i < p_j$  &  $q_i < q_j$ , let  $\mathbf{D}_S := (\mathbf{D}_i)_{i \in S}$  and  $\mathbf{D}_T := (\mathbf{D}_i)_{i \in T}$ . If for any  $\alpha, \beta \in R$

$$\mathbb{P}_{\alpha}^{\mathbf{D}^S} = \mathbb{P}_{\beta}^{\mathbf{D}^S} \quad (58)$$

$$\implies \mathbb{P}_{\alpha}^{(\mathbf{D}_i, \mathbf{Y}_i)_{i \in S}} = \mathbb{P}_{\beta}^{(\mathbf{D}_i, \mathbf{Y}_i)_{i \in S}} \quad (59)$$

then  $\mathbb{P}_{\square}$  commutes with marginalisation.

Neither condition implies the other.

**Lemma 4.4.** *Commutativity of exchange does not imply commutativity or vice versa.*

*Proof.* Suppose  $D = Y = \{0, 1\}$  and we have a conditional probability model  $(\mathbb{P}_{\square}^{\mathbf{Y}|\mathbf{D}}, A)$  where  $\mathbf{D} = (\mathbf{D}_1, \mathbf{D}_2)$ ,  $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2)$  and  $A$  contains all deterministic probability measures in  $\Delta(D^2)$ . If

$$\mathbb{P}_{\square}^{\mathbf{Y}_1 \mathbf{Y}_2 | \mathbf{D}_1 \mathbf{D}_2}(y_1, y_2 | d_1, d_2) = \mathbb{I}((y_1, y_2) = (d_1 + d_2, d_1 + d_2)) \quad (60)$$

Then  $\mathbb{P}_{\delta_{00}}^{\mathbf{Y}_1 \mathbf{D}_1}(y_1) = \mathbb{I}(y_1 = 0)$  while  $\mathbb{P}_{\delta_{01}}^{\mathbf{Y}_1} = \mathbb{I}(y_1 = 1)$ . However,  $\delta_0 0^{\mathbf{D}_1} = \delta_{01}^{\mathbf{D}_1} = \delta_0^{\mathbf{D}_1}$  so  $\mathbb{P}_{\square}$  does not commute with marginalisation. However, taking  $(d_i, d_j) := \delta_{d_i d_j} \in A$ ,

$$\mathbb{P}_{d_2, d_1}^{\mathbf{Y}_1 \mathbf{D}_1 \mathbf{Y}_2 \mathbf{D}_2}(y_1, d_1, y_2, d_2) = \mathbb{I}((y_1, y_2) = (d_2 + d_1, d_2 + d_1)) \quad (61)$$

$$= \mathbb{I}((y_2, y_1) = (d_1 + d_2, d_1 + d_2)) \quad (62)$$

$$= \mathbb{P}_{d_1, d_2}^{\mathbf{Y}_1 \mathbf{D}_1 \mathbf{Y}_2 \mathbf{D}_2}(y_2, d_2, y_1, d_1) \quad (63)$$

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

Alternatively, suppose the same setup, but define  $\mathbb{P}_{\square}$  instead by, for all  $\alpha \in A$

$$\mathbb{P}_{\square}^{\mathbf{Y}_1 \mathbf{Y}_2 | \mathbf{D}_1 \mathbf{D}_2}(y_1, y_2 | d_1, d_2) = \mathbb{I}((y_1, y_2) = (0, 1)) \quad (64)$$

Then  $\mathbb{P}_{\square}$  commutes with marginalisation. If  $\mathbb{P}_{\alpha}^{\mathbf{D}^S} = \mathbb{P}_{\beta}^{\mathbf{D}^S}$  for  $S \subset \{0, 1\}$  then

$$\mathbb{P}_{\alpha}^{\mathbf{Y}_S \mathbf{D}^S}(y_s, d_s) = \sum_{y'_2 \in \{0, 1\}^{S^C}} \mathbb{I}((y_1, y_2) = (0, 1)) \mathbb{P}_{\alpha}^{\mathbf{D}^S}(d_s) \quad (65)$$

$$= \mathbb{P}_{\beta}^{\mathbf{Y}_S \mathbf{D}^S}(y_s, d_s) \quad (66)$$

but not exchange. For all  $\alpha, \beta \in A$ :

$$\mathbb{P}_\alpha Y_1 Y_2(y_1, y_2) = \mathbb{I}[(y_1, y_2) = (0, 1)] \quad (67)$$

$$\neq \mathbb{P}_\beta Y_1 Y_2(y_2, y_1) \quad (68)$$

□

Although commutativity of marginalisation seems to be a bit like non-interference – the marginal distribution I get for  $Y_i$  depends only on the decision  $D_i$  – it still allows for some models in which we seem to have interference of a kind. For example: in the first experiment I flip a coin and decide either to pass the results to the second experiment ( $D_1 = 0$ ) or flip another coin and pass those results to the second experiment ( $D_1 = 1$ ). In the second I either copy the results I have been given ( $D_2 = 0$ ) or invert them ( $D_2 = 1$ ). Then

- The marginal distribution of both experiments is Bernoulli(0.5) no matter what choices I make, so it satisfies Definition 4.3
- Nevertheless, the choice for the first experiment seems to “affect” the result of the second experiment (affect in quotes because it is an intuitive judgement, not a formal property)

Here we are most interested in the conjunction of these assumptions, a condition we call *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.

## 4.2 When does a canonical “effect of a decision” exist?

The main result in this section is Theorem 4.6 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  $\mathbb{P}_\square^{Y_1|D_1H}$ . This can be interpreted as expressing the idea that all experimental units  $(Y_i, D_i)$  share a canonical but unknown “consequence function”  $D \rightarrow Y$ . As discussed already in Section 2.10, the existence of such a conditional probability implies the existence of a common unknown *curried* conditional probability for all experimental units, which resembles a potential outcomes model. In fact, we prove the existence of a curried representation first, in Lemma 4.2.

[Exchangeable curried representation] A conditional probability model  $(\mathbb{P}_\square^{\overline{Y|D}}, A)$  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

$$\mathbb{P}_{\square}^{Y|D} = \begin{array}{c} \text{triangle} \\ \text{P}^{Y^D} \\ \text{D} \end{array} \begin{array}{c} \text{box} \\ \text{L}^{D,Y^D} \end{array} \text{Y} \quad (69)$$

$$\iff \quad (70)$$

$$\mathbb{P}_{\square}^{Y|D}(y|d) = \mathbb{P}^{(Y_{d_i}^D)_{i \in \mathbb{N}}}(y) \quad (71)$$

Where  $\mathbb{P}^{Y^D}$  is an exchangeable probability measure on  $Y^{D \times \mathbb{N}}$ , for convenience we extend the sample space with the random variable  $Y^D := (Y_{ij}^D)_{i \in D, j \in \mathbb{N}}$  and  $\mathbb{L}^{D,Y^D}$  is the Markov kernel associated with the lookup function

$$l : D^{\mathbb{N}} \times Y^{D \times \mathbb{N}} \rightarrow Y \quad (72)$$

$$((d_i)_{i \in \mathbb{N}}, (y_{ij})_{i \in D, j \in \mathbb{N}}) \mapsto y_{d_i i} \quad (73)$$

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

Define

$$\mathbb{P}^{Y^D}((y_{ij})_{D \times \mathbb{N}}) := \mathbb{P}_e^Y((y_{|D|i+j})_{i \in D, j \in \mathbb{N}}) \quad (74)$$

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

$$\mathbb{Q} : D \rightarrow Y \quad (75)$$

$$\mathbb{Q} := \begin{array}{c} \text{triangle} \\ \text{P}^{Y^D} \\ \text{D} \end{array} \begin{array}{c} \text{box} \\ \text{L}^{D,Y^D} \end{array} \text{Y} \quad (76)$$

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

$$\sum_{y \in Y^{-1}(y_A)} \mathbb{Q}(y|d) = \sum_{y \in Y^{-1}(y_A)} \mathbb{P}^{(Y_{d_i}^D)_{i \in A}}(y) \quad (77)$$

$$= \sum_{y \in Y^{-1}(y_A)} \mathbb{P}_e^{(Y_{|D|d_i+i})_{i \in A}}(y) \quad (78)$$

$$= \mathbb{P}_e^{Y^B}(y_A) \quad (79)$$

$$= \mathbb{P}_d^{Y^A}(y_A) \quad \text{by causal contractibility} \quad (80)$$

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

$$\mathbb{Q}(y|d) = \mathbb{P}_d^{\mathbf{Y}}(y) \quad (81)$$

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

$$\mathbb{Q}(y|d) = \mathbb{P}_d^{\mathbf{Y}|\mathbf{D}}(y|d) \quad (82)$$

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

$$\mathbb{P}^{\mathbf{Y}|\mathbf{D}} = \mathbb{Q} \quad (83)$$

Next we will show  $\mathbb{P}^{\mathbf{Y}^D}$  is contractible. Consider any subsequences  $\mathbf{Y}_S^D$  and  $\mathbf{Y}_T^D$  of  $\mathbf{Y}^D$  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

$$\mathbb{P}^{\mathbf{Y}_S^D} = \mathbb{P}_e^{\mathbf{Y}_{\rho(S)}} \quad (84)$$

$$= \mathbb{P}_e^{\mathbf{Y}_{\rho(T)}} \quad \text{by contractibility of } \mathbb{P} \text{ and the equality } e_{\rho(S)} = e_{\rho(T)} \quad (85)$$

$$= \mathbb{P}^{\mathbf{Y}_T^D} \quad (86)$$

If: Suppose

$$\mathbb{P}^{\mathbf{Y}|\mathbf{D}} = \begin{array}{c} \triangle \\ \mathbb{P}^{\mathbf{Y}^D} \\ \text{D} \end{array} \begin{array}{c} \text{---} \\ \mathbb{L}^{\mathbf{D}, \mathbf{Y}^D} \end{array} \mathbf{Y} \quad (87)$$

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

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

By definition,

$$\mathbb{P}^{\mathbf{Y}_S|\mathbf{D}}(y_S|d) = \sum_{y_S^D \in Y^{|\mathbf{D}| \times |S|}} \mathbb{P}^{\mathbf{Y}_S^D}(y_S^D) \mathbb{L}^{\mathbf{D}_S, \mathbf{Y}^S}(y_S|d, y_S^D) \quad (88)$$

$$= \sum_{y_S^D \in Y^{|\mathbf{D}| \times |T|}} \mathbb{P}^{\mathbf{Y}_T^D}(y_S^D) \mathbb{L}^{\mathbf{D}_S, \mathbf{Y}^S}(y_S|d, y_S^D) \quad \text{by contractibility of } \mathbb{P}^{\mathbf{Y}_T^D} \quad (89)$$

$$= \mathbb{P}^{\mathbf{Y}_T|\mathbf{D}}(y_S|d) \quad (90)$$

□

The curried representation of Lemma 4.2 does not need to support an interpretation as a distribution of potential outcomes. For example, consider a series of bets on fair coinflips – in this case, the consequence  $Y_i$  is uniform on  $\{0, 1\}$  for any decision  $D_i$ . The  $D = Y = \{0, 1\}$  and  $\mathbb{P}_\alpha^{Y^n}(y) = \prod_{i \in [n]} 0.5$  for all  $n, y \in Y^n$ ,  $\alpha \in R$ . Then the construction in Lemma 4.2 yields  $\mathbb{P}^{Y^D}(y_i^D) = \prod_{j \in D} 0.5$  for all  $y_i^D \in Y^D$ . That is,  $Y_i^0$  and  $Y_i^1$  are independent and uniformly distributed. However, if we wanted  $Y_i^0$  to represent “what would happen if I bet on outcome 0 on turn  $i$ ” and  $Y_i^1$  to represent “what would happen if I bet on outcome 1 on turn  $i$ ”, then it seems that we ought to have  $Y_i^0 = 1 - Y_i^1$ .

We could suppose that Lemma 4.2 provides necessary but not sufficient conditions for the existence of a potential outcomes representation of a conditional probability model. However, it doesn’t seem to succeed at that either. We note, for example, that Rubin (2005) does not assume that the distribution of potential outcomes is exchangeable. A non-exchangeable  $\mathbb{P}^{Y^D}$  does not induce a causally contractible conditional probability model, and at the same time commutativity with marginalisation is not sufficient for a conditional probability model to support a curried representation in the sense of Lemma 4.2. What seems to be missing is an additional assumption that consequences are mutually independent of one another given the associated decision.

We can also represent contractible conditional probability models repeated copies of an unknown “consequence function”, a Markov kernel that maps from decisions to probability distributions over consequences, coupled by a common hypothesis  $H$ .

**Theorem 4.6.** *Suppose we have a fundamental probability set  $\Omega$  and a do model  $(\mathbb{P}, D, Y, R)$  such that  $D := (D_i)_{i \in \mathbb{N}}$  and  $Y := (Y_i)_{i \in \mathbb{N}}$ .  $\mathbb{P}$  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*

$$\mathbb{P}^{Y | HD} = \begin{array}{c} H \\ D \end{array} \begin{array}{c} \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\Pi_i} \boxed{\mathbb{P}^{Y_0 | HD_0}} Y_i \\ i \in \mathbb{N} \end{array} \quad (91)$$

$$\iff \quad (92)$$

$$Y_i \perp\!\!\!\perp Y_{\mathbb{N} \setminus i}, D_{\mathbb{N} \setminus i} | HD_i \quad \forall i \in \mathbb{N} \quad (93)$$

$$\wedge \mathbb{P}^{Y_i | HD_i} = \mathbb{P}^{Y_0 | HD_0} \quad \forall i \in \mathbb{N} \quad (94)$$

*Proof.* We make use of Lemma 4.2 to show that we can represent the conditional probability as an exchangeable tabular probability distribution. We then use the property of exchangeability of the columns of that distribution in conjunction with De Finetti’s theorem to derive the result.  $\square$

### 4.3 In decision problems, policies are extreme

Theorem 4.6 shows that a formal property of conditional probability models – causal contractibility – implies the existence of a map  $D \times H \rightarrow Y$  that,

when coupled with a distribution over  $H$ , yields the original model. We have also suggested that causal contractibility might be a property of a model of an experiment where we expect that experimental outcomes will be the same if we shuffle the planned experimental actions and apply the inverse shuffle to results, and also if we believe that outcomes for units of interest will be unchanged if we substitute the experimental plan for any other plan that prescribes the same actions for the same units.

One could be inclined to leave things there and say that the question of whether causal contractibility holds is up to the judgement of the practitioner. However, we believe that making this judgement is a subtle issue that deserves further commentary.

The subtlety is illustrated by comparing the following three situations:

1. I am considering different policies for treating 100 patients for lower back pain, and want to predict the outcomes of each policy
2. I am treating 50 patients for lower back pain, not necessarily according to a fixed policy, and subsequently want to predict the outcomes of 50 additional patients treated according to a fixed policy

The two situations seem similar - in all cases, I treat 100 patients for lower back pain and I want to predict some outcomes. The only difference is the presence or absence of a fixed policy.

If we accept for the sake of argument that we are going to use probability models to predict these outcomes, it seems reasonable that a conditional probability model for the first situation should be causally contractible. I have no means of distinguishing one ordering of patients from the other, and so it seems that I should use the same probability to predict outcomes for any ordering. Furthermore, the treatment of each patient seems like it is a separate matter - there's no reason to expect patient  $i$ 's outcome to differ if I give different treatments to the rest of the patients.

Causal contractibility implies that, given enough patients and any policy with full support, I will converge to a single Markov kernel representing the response of any patient to a treatment, which will be equal to the conditional probability of recovery given treatment. On the other hand, it is well understood that I should not always assume that the conditional probability of recovery given treatment is a good guide to my selection of policy. In particular, it would be unwise to estimate the probability of recovery given treatment for the first 50 patients in the second example and assume this is the same as the conditional probability of recovery given treatment for the second 50 patients.

The question is: what is being accomplished by the use of a policy that distinguishes situation 1 from 2? The answer, we think, is that the word "policy" implies more than just "a probability distribution associated with a variable something that we happen to call a decision". In particular, suppose policy  $\alpha_1$  is to always choose decision 1 and policy  $\alpha_2$  is to always choose decision 2. Choosing a mixed policy  $0.5\alpha_1 + 0.5\alpha_2$  does not mean "I choose to be uncertain by degree  $a$  over which decision I will choose", it means "I will consult a random

number generator that I know to yield 1 half of the time and 2 the other half, then choose  $\alpha_i$  according to the result”.

We can use probability gap models for many things. Maybe we want  $\mathbb{P}_\alpha$  on Monday we do some work to figure out  $\mathbb{P}_\square^{Y|X}$  and on Tuesday we were going to do a bit more to work out  $\mathbb{P}_\alpha^X$ , but actually we found the answer to our question would be the same in any case. Whether or not we explicitly set it up, we’ve made use of a conditional probability model to come to this conclusion.

However, we actually want to model decision problems. For a decision problem we compare different policies  $\alpha$  and choose the best one according to some criterion. We choose decisions according to the policy we arrive at. Decisions can be represented by variables, but they have a very important property: they must be deterministic for any policy choice  $\alpha$ .

**Definition 4.7** (Decision variable). Given a probability gap model  $(\mathbb{P}_\square, A)$  on  $\Omega$ , a variable  $D$  is a *decision variable* if there is some variable  $R$  which we call *purely random* such that, for  $\alpha \in A$ ,  $\mathbb{P}_\alpha^{D|R}$  is deterministic.

We don’t attempt to define what a “purely random” variable is here.

Because  $D$  is deterministic given  $R$ , we have  $Y \perp\!\!\!\perp_{\mathbb{P}_\square} R|D$ , and so we can marginalise over  $R$  and talk about policies as if they address marginal probabilities over  $D$ :  $\mathbb{P}_\alpha^D$ .

The reason we include this definition is that it constrains the measurement procedures that are allowed to be associated with decision variables. In particular, all such measurement procedures must, with probability 1, yield a function of the policy choice  $\alpha$  and a purely random measurement procedure  $\mathcal{R}$ . In particular, this means that if I talk about a “mixed policy”  $a\alpha + (1 - a)\beta$  I mean: consult a purely random process and, if the result is consistent with  $a$  act according to  $\alpha$ , while if the result is consistent with  $(1 - a)$  act according to  $\beta$ . This is what people usually mean when they talk about mixing policies, and we think this is important enough to bake into the theory rather than leave it as an implicit side constraint.

This is quite important, because it helps us to distinguish between these two apparently similar situations:

1. Some doctor is going to treat 49 (otherwise unknown) patients for lower back pain, then I will treat one patient for the same
2. I am preparing to treat 50 (otherwise unknown) patients for lower back pain

In the second case, because decisions are deterministic, I don’t learn anything about the patient when I treat them or not. In the first case, because the other doctor’s decisions are uncertain from my point of view, their treatment decisions can inform me about the patient *and* determine whether the patient takes the medicine or not.

If we go ahead and set this up formally, we could use probability gap models for both;



In the first case, we know it is unwise to assume that the result of giving patient 50 treatment  $x$  will be probabilistically the same as the results experienced by all previous patients who received  $x$ . On the other hand, in the second case, it seems reasonable to expect the result of giving treatment  $x$  to patient 50 is, a priori, the same as giving treatment  $x$  to any other patient. That is, causal contractibility seems plausible if we imagine ourselves in a situation of planning a sequence of treatments, but a rather similar assumption seems unwise if we imagine the treatments having already taken place.

There is a subtle asymmetry between these situations: while “patient  $i$  receives treatment  $x$ ” is a sensible proposition for both situations, *it corresponds to a different measurement procedure*.

#### 4.3.1 Extended conditional independence

Needs a support condition

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 Constantinou and Dawid (2017) *extended conditional independence*, 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:

**Theorem 4.8.** *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  $\mathbb{P}_{\alpha}^{VW} := \mathbb{P}_{\alpha}^W \odot \mathbb{P}^{V|W}$*

$$A \perp\!\!\!\perp_{\mathbb{P}, ext} (B, D)|(C, E) \iff A^* \perp\!\!\!\perp_{\mathbb{P}_{\alpha}} (B^*, D^*)|(C^*, E^*) \quad (95)$$

Where  $\perp\!\!\!\perp_{\mathbb{P}_{\alpha}}$  is order 0 conditional independence.

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

**Theorem 4.9.** *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  $\mathbb{P}_{\alpha}^{VW} := \mathbb{P}_{\alpha}^W \odot \mathbb{P}^{V|W}$*

$$A \perp\!\!\!\perp_{\mathbb{P}, ext}^1 (B, D)|(C, E) \iff A^* \perp\!\!\!\perp_{\mathbb{P}} (B^*, D^*)|(C^*, E^*) \quad (96)$$

*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 4.8,  $A \perp\!\!\!\perp_{\mathbb{P}, ext} (B, D)|(C, E)$ .

Only if:

For any  $\beta$ ,  $\mathbb{P}_{\beta}^{ABC|DE} = \mathbb{P}_{\beta}^{DE} \odot \mathbb{P}^{ABC|DE}$ . By Lemma 2.14, we have  $\mathbb{P}^{A|BCDE}$  such that

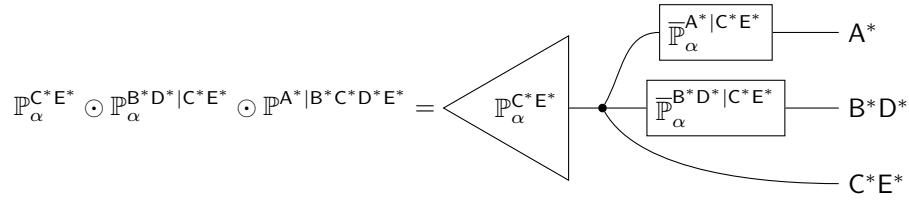
$$\mathbb{P}_\beta^{A^*B^*C^*D^*E^*} = \mathbb{P}_\beta^{D^*E^*} \odot \mathbb{P}^{B^*C^*|D^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (97)$$

$$= \mathbb{P}_\beta^{B^*C^*D^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (98)$$

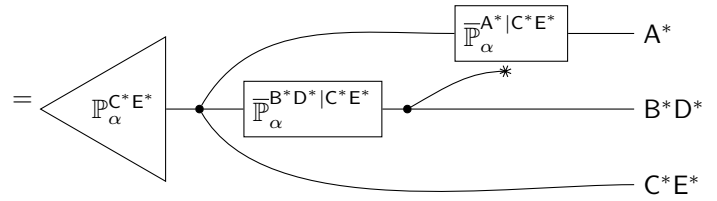
$$= \mathbb{P}_\beta^{C^*E^*} \odot \mathbb{P}_\beta^{B^*D^*|C^*E^*} \odot \mathbb{P}^{A^*|B^*C^*D^*E^*} \quad (99)$$

By Theorem 4.8, we have some  $\alpha$  such that  $\mathbb{P}_\alpha^{D^*E^*}$  is strictly positive on the range of  $(D^*, E^*)$  and  $A^* \perp_{\mathbb{P}_\alpha} (B^*, D^*)|(C^*, E^*)$ .

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



(100)



(101)

$$= \mathbb{P}_\alpha^{C^*E^*} \odot \mathbb{P}_\alpha^{B^*D^*|C^*E^*} \odot (\mathbb{P}_\alpha^{A^*|C^*E^*} \otimes \text{erase}_{BD}) \quad (102)$$

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

(103)

$$= \begin{array}{c} \text{---} \triangle \text{---} \mathbb{P}_\beta^{C^*E^*} \text{---} \bullet \begin{array}{l} \boxed{\mathbb{P}_\alpha^{A^*|C^*E^*}} \text{---} A^* \\ \boxed{\mathbb{P}_\beta^{B^*D^*|C^*E^*}} \text{---} B^*D^* \\ \text{---} C^*E^* \end{array} \end{array} \quad (104)$$

☐

We can deduce conditional independences in probability combs when conditional probabilities exist and they are *unresponsive* to some input variables.

**Definition 4.10** (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

$$\mathbb{P}_{\square}^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{P}_{\square}^{Y|W}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (105)$$

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

**Definition 4.11** (Domination). Given a probability gap model  $\mathbb{P}_{\square} : A \rightarrow \Delta(\Omega)$ ,  $\alpha \in A$  dominates  $A$  if  $\mathbb{P}_{\beta}(B) > 0 \implies \mathbb{P}_{\alpha}(B) > 0$  for all  $\beta$  in  $A$ ,  $B \in \mathcal{F}$ .

**Theorem 4.12** (Conditional independence from kernel unresponsiveness). *Given discrete  $\Omega$ , variables  $W : \Omega \rightarrow W$ ,  $X : \Omega \rightarrow X$ ,  $Y : \Omega \rightarrow Y$  and a probability gap model  $\mathbb{P}_{\square} : A \rightarrow \Delta(\Omega)$  with conditional probability  $\mathbb{P}_{\square}^{Y|WX}$  and such that there is  $\alpha \in A$  dominating  $A$ ,  $Y \perp\!\!\!\perp_{\mathbb{P}_{\square}} X|W$  if and only if  $\mathbb{P}_{\square}^{Y|WX}$  is unresponsive to  $W$ .*

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

$$\mathbb{P}_\alpha^{Y|WX} = W \text{ --- } \boxed{\mathbb{P}_\alpha^{Y|W}} \text{ --- } Y \quad (106)$$

$X \text{ --- } *$

And so, by Theorem 2.30,  $Y \perp\!\!\!\perp_{\mathbb{P}_\alpha} X|W$  for all  $\alpha \in A$ , and so  $Y \perp\!\!\!\perp_{\mathbb{P}_\square} X|W$ . Only if: For  $\alpha$  dominating  $A$ , by Theorem 2.30, there exists a version of  $\mathbb{P}_\alpha^{Y|WX}$  unresponsive to  $W$ . Because  $\alpha$  dominates  $A$ , every version of  $\mathbb{P}_\alpha^{Y|WX}$  is a version of  $\mathbb{P}_\beta^{Y|WX}$  for all  $\beta \text{ in } A$ , thus it is a version of  $\mathbb{P}_\square^{Y|WX}$  also.  $\square$

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

We also need 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  $X$ .

### 4.3.2 Graphical properties of conditional independence

It is well-known that directed acyclic graphs are able to represent some conditional independence properties of probability models via the graphical property of *d-separation*. String diagrams are similar to directed acyclic graphs, and string diagrams can be translated into directed acyclic graphs and vice-versa (Fong, 2013). Thus we expect that a property analogous to d-separation can be defined for string diagrams.

We can reason from graphical properties of model disintegrations to graphical properties of models as Theorem 4.12. A general theory akin to d-separation for string diagrams may facilitate a more general understanding of how conditional independence properties of a model relate to conditional independence properties of its components.

## 4.4 Results I use that don't really fit into the flow of the text

### 4.4.1 Repeated variables

Lemmas 4.13 and 4.14 establish that models of repeated variables must connect the repetitions with a copy map.

**Lemma 4.13** (Output copies of the same variable are identical). *For any  $\Omega$ ,  $X, Y, Z$  random variables on  $\Omega$  and conditional probability  $\mathbb{K}^{YZ|X}$ , there is a conditional probability  $\mathbb{K}^{YYZ|X}$  unique up to impossible values of  $X$  such that*

$$X \text{ --- } \boxed{\mathbb{K}^{YYZ|X}} \begin{matrix} \text{---}^* \\ \text{---} \\ \text{---} \end{matrix} \begin{matrix} Y \\ Z \end{matrix} = \mathbb{K}^{YZ|X} \quad (107)$$

and it is given by

$$\mathbb{K}^{YYZ|X} = X \text{ --- } \boxed{\mathbb{K}^{YZ|X}} \begin{matrix} \text{---} Y \\ \text{---} Y \\ \text{---} Z \end{matrix} \quad (108)$$

$$\iff \quad (109)$$

$$\mathbb{K}^{YYZ|X}(y, y', z|x) = \llbracket y = y' \rrbracket \mathbb{K}^{YZ|X}(y, z|x) \quad (110)$$

$$(111)$$

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

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

$$\mathbb{K}^{YYZ|X}(y, y', z|x) = \llbracket y = y' \rrbracket \mathbb{L}(y, z|x) \quad (112)$$

$$(113)$$

But then

$$\mathbb{K}^{YZ|X}(y, z|x) = \sum_{y' \in Y} \mathbb{K}^{YYZ|X}(y, y', z|x) \quad (114)$$

$$= \mathbb{L}(y, z|x) \quad (115)$$

$$(116)$$

□

**Lemma 4.14** (Copies shared between input and output are identical).

*This got mixed up at some point and needs ot be unmixed-up*

For any  $\mathbb{K} : (X, Y) \rightarrow (X, Z)$ ,  $\mathbb{K}$  is a model iff there exists some  $\mathbb{L} : (X, Y) \rightarrow Z$  such that

$$\mathbb{K} = \begin{array}{c} X \\ \text{---} \bullet \\ Y \end{array} \begin{array}{c} \text{---} \curvearrowright X \\ \boxed{\mathbb{K}^{Z|XY}} \\ \text{---} Z \end{array} \quad (117)$$

$$\iff \quad (118)$$

$$\mathbb{K}_{x,y}^{Ix',z} = \llbracket x = x' \rrbracket \mathbb{L}_{x,y}^z \quad (119)$$

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

$$\begin{array}{c} X \\ \text{---} \end{array} \begin{array}{c} \boxed{\mathbb{K}^{XZ|XY}} \\ \text{---} Y \end{array} \begin{array}{c} \text{---} * \\ \text{---} Z \end{array} = \mathbb{K}^{XZ|XY} \quad (120)$$

and it is given by

$$\mathbb{K}^{XZ|XY} = X \text{---} \begin{array}{c} \boxed{\mathbb{K}^{YZ|X}} \\ \text{---} Y \\ \text{---} Z \end{array} \begin{array}{c} \text{---} Y \\ \text{---} Y \\ \text{---} Z \end{array} \quad (121)$$

$$\iff \quad (122)$$

$$\mathbb{K}^{XZ|XY}(x, z|x', y) = \llbracket x = x' \rrbracket \mathbb{K}^{Z|XY}(z|x', y) \quad (123)$$

$$(124)$$

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

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

$$\mathbb{K}^{XZ|XY}(x', z|x, y) = 0 = \llbracket x = x' \rrbracket \mathbb{L}(z|x, y) \quad (125)$$

$$(126)$$

But then

$$\mathbb{K}^{Z|XY}(y, z|x) = \sum_{x' \in X} \mathbb{K}^{XZ|XY}(x', z|x, y) \quad (127)$$

$$= \mathbb{L}(z|x, y) \quad (128)$$

$$(129)$$

□

**Theorem 4.15** (Existence of valid conditional probabilities). *Given a probability gap model  $\mathbb{P}_\square : A \rightarrow \Delta(\Omega)$  along with a valid conditional probability  $\mathbb{P}_\square^{XY|W}$ , there exists a valid conditional probability  $\mathbb{P}_\square^{Y|WX}$ .*

*Proof.* From Lemma 2.14, we have the existence of some Markov kernel  $\mathbb{P}_\square^{Y|WX} : W \times X \rightarrow Y$  such that

$$\mathbb{P}_\square^{XY|W} = \mathbb{P}_\square^{X|W} \odot \mathbb{P}_\square^{Y|WX} \quad (130)$$

By definition of conditional probability, for any insert  $\alpha \in A$  there exists  $\mathbb{P}_\alpha^W \in \Delta(W)$  such that

$$\mathbb{P}_\alpha^{WXY} = \mathbb{P}_\alpha^W \odot \mathbb{P}_\square^{XY|W} \quad (131)$$

Thus

$$\mathbb{P}_\alpha^{WXY} = \mathbb{P}_\alpha^W \odot (\mathbb{P}_\square^{X|W} \odot \mathbb{P}_\square^{Y|WX}) \quad (132)$$

$$= (\mathbb{P}_\alpha^W \odot \mathbb{P}_\square^{X|W}) \odot \mathbb{P}_\square^{Y|WX} \quad (133)$$

Let  $\text{erase}_Y : Y \rightarrow \{*\}$  be the erase function on  $Y$  (as opposed to the erase kernel) and  $\text{idf}_{W \times X}$  be the identity function on  $W \times X$ . Noting that

$$(W, X) = (\text{idf}_{W \times X} \otimes \text{erase}_Y) \circ (W, X, Y) \quad (134)$$

By Lemma ?? together with Theorem 2.18 we have for all  $\alpha$ :

$$\mathbb{P}_\alpha^{XW} = \mathbb{P}_\alpha^{WXY}(\text{id}_{W \times X} \otimes \text{erase}_Y) \quad (135)$$

$$= \mathbb{P}_\alpha^W \odot (\mathbb{P}_\square^{X|W} \odot \mathbb{P}_\square^{Y|WX})(\text{id}_{W \times X} \otimes \text{erase}_Y) \quad (136)$$

$$= \mathbb{P}_\alpha^W \odot \mathbb{P}_\square^{X|W} \quad (137)$$

Then

$$\mathbb{P}_\alpha^{XWY} = (\mathbb{P}_\alpha^{XW}) \odot \mathbb{P}_\square^{Y|WX} \quad (138)$$

And so  $\mathbb{P}_\square^{Y|WX}$  is a  $Y|WX$  conditional probability. We also want it to be valid, so we will verify that it can be chosen as such.

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

## 4.5 Validity

**Theorem 4.16** (Validity). *Given  $(\Omega, \mathcal{F})$ ,  $X : \Omega \rightarrow X$ ,  $\mathbb{J} \in \Delta(X)$  with  $\Omega$  and  $X$  standard measurable, there exists some  $\mu \in \Delta(\Omega)$  such that  $\mu^X = \mathbb{J}$  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  $X \bowtie x = \emptyset$  but  $\mathbb{J}(x) > 0$ . Then

$$\mu^X(x) = \mu(X \bowtie x) \quad (139)$$

$$= \sum_{x' \in X} \mathbb{J}(x') \mathbb{K}(X \bowtie x|x') \quad (140)$$

$$= 0 \quad (141)$$

$$\neq \mathbb{J}(x) \quad (142)$$

$\square$

**Lemma 4.17** (Copy-product is an intersection of probability sets). *Given  $(\Omega, \mathcal{F})$ ,  $X : \Omega \rightarrow (X, \mathcal{X})$ ,  $Y : \Omega \rightarrow (Y, \mathcal{Y})$ ,  $Z : \Omega \rightarrow (Z, \mathcal{Z})$  all standard measurable and valid candidate conditionals  $\mathbb{P}_\square^{Y|X}$  and  $\mathbb{Q}_\square^{Z|YX}$  defining probability sets  $\mathbb{P}_\square$  and  $\mathbb{Q}_\square$ , then the probability set  $\mathbb{R}_\square$  defined by  $\mathbb{R}_\square^{YZ|X} := \mathbb{P}_\square^{Y|X} \odot \mathbb{Q}_\square^{Z|YX}$  is equal to  $\mathbb{P}_\square \cap \mathbb{Q}_\square$ .*

*Proof.* By assumption

$$\mathbb{R}_{\{\}}^{\mathbf{YZ}|\mathbf{X}} := \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (143)$$

Therefore for any  $\mathbb{R}_a \in \mathbb{R}_{\{\}}$

$$\mathbb{R}_a^{\mathbf{XYZ}} = \mathbb{R}_a^{\mathbf{X}} \odot \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (144)$$

$$\implies \mathbb{R}_a^{\mathbf{XY}} = \mathbb{R}_a^{\mathbf{X}} \odot \mathbb{P}_{\{\}}^{\mathbf{Y}|\mathbf{X}} \quad (145)$$

$$\wedge \mathbb{R}_a^{\mathbf{XYZ}} = \mathbb{R}_a^{\mathbf{XY}} \odot \mathbb{Q}_{\{\}}^{\mathbf{Z}|\mathbf{YX}} \quad (146)$$

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 4.18** (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

$$\forall A \in \mathcal{X} : \mathbf{X} \bowtie A = \emptyset \implies \mathbb{P}^{\mathbf{X}|\ast}(A|\ast) = 0 \quad (147)$$

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

$$\forall A \in \mathcal{X} : \mathbf{X} \bowtie A = \emptyset \implies \mathbb{P}^{\mathbf{X}}(A) = 0 \quad (148)$$

$\square$

**Lemma 4.19** (Copy-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  $x \in X(\Omega)$ , 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  $C \in \mathcal{Z}$

$$\mathbb{R}^{\mathbf{YZ}|\mathbf{X}}(B \times C|x) = \int_B \mathbb{Q}^{\mathbf{Z}|\mathbf{YX}}(C|y, x) \mathbb{P}^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (149)$$

$$\leq \mathbb{P}^{\mathbf{Y}|\mathbf{X}}(B|x) \quad (150)$$

$$= 0 \quad (151)$$

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:

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

$$= 0 \quad (153)$$

□

**Corollary 4.20** (Valid conditionals are validly extendable to valid distributions). *Given  $\Omega$ ,  $\mathbf{U} : \Omega \rightarrow U$ ,  $\mathbf{W} : \Omega \rightarrow W$  and a valid candidate conditional  $\mathbb{T}^{\mathbf{W}|\mathbf{U}}$ , then for any valid candidate conditional  $\mathbb{V}^{\mathbf{U}}$ ,  $\mathbb{V}^{\mathbf{U}} \odot \mathbb{T}^{\mathbf{W}|\mathbf{U}}$  is a valid candidate probability.*

*Proof.* Applying Lemma 4.19 choosing  $\mathbf{X} = *$ ,  $\mathbf{Y} = \mathbf{U}$ ,  $\mathbf{Z} = \mathbf{W}$  and  $\mathbb{P}^{\mathbf{Y}|\mathbf{X}} = \mathbb{V}^{\mathbf{U}|*}$  and  $\mathbb{Q}^{\mathbf{Z}|\mathbf{YX}} = \mathbb{T}^{\mathbf{W}|\mathbf{U}*}$  we have  $\mathbb{R}^{\mathbf{WU}|*} := \mathbb{V}^{\mathbf{U}|*} \odot \mathbb{T}^{\mathbf{W}|\mathbf{U}*}$  is a valid conditional probability. Then  $\mathbb{R}^{\mathbf{WU}} \cong \mathbb{R}^{\mathbf{WU}|*}$  is valid by Theorem 4.18. □

**Theorem 4.21** (Validity of conditional probabilities). *Suppose we have  $\Omega$ ,  $\mathbf{X} : \Omega \rightarrow X$ ,  $\mathbf{Y} : \Omega \rightarrow Y$ , with  $\Omega$ ,  $X$ ,  $Y$  discrete. A conditional  $\mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is valid if and only if for all valid candidate distributions  $\mathbb{V}^{\mathbf{X}}$ ,  $\mathbb{V}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}}$  is also a valid candidate distribution.*

*Proof.* If: this follows directly from Corollary 4.20.

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

$$(\mathbb{V}^{\mathbf{X}} \odot \mathbb{T}^{\mathbf{Y}|\mathbf{X}})(x, y) = \mathbb{T}^{\mathbf{Y}|\mathbf{X}}(y|x) \mathbb{V}^{\mathbf{X}}(x) \quad (154)$$

$$= \mathbb{T}^{\mathbf{Y}|\mathbf{X}}(y|x) \quad (155)$$

$$> 0 \quad (156)$$

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

## 4.6 Combs

**Theorem ??** (Equivalence of comb representations). *Given sample space  $(\Omega, \mathcal{F})$ , a finite collection of variables  $\mathbf{X}_i : \Omega \rightarrow (X_i, \mathcal{X}_i)$  for  $i \in [n]$ ,  $X_i$  discrete, and a*

disassembled probability comb  $\{\mathbb{P}_{\square}^{\mathbf{X}_i|\mathbf{X}_{[i-1]}}|i \in \mathbf{X}_{[n]_{\text{odd}}}\}$ , for any  $l \in [n]_{\text{odd}}$  and any  $\mathbb{K} : X_{[l-1]} \rightarrow X_l$

$$\left(\bigodot_{j \in [l-1]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}\right) \odot \mathbb{K} \stackrel{\mathbb{P}_{\square}}{\cong} \left(\bigodot_{j \in [l]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}\right) \quad (157)$$

$$\implies \mathbb{K} \stackrel{\mathbb{P}_{\square}}{\cong} \mathbb{P}_{\square}^{\mathbf{X}_l|\mathbf{X}_{[l-1]}} \quad (158)$$

*Proof.* Equality is trivial for  $l = 1$ .

For a sequence  $x_{[l-2]} \in X_{[l-2]}$ , let  $e_{[l-2]}$  be the even indices of  $x_{[l-2]}$  and  $o_{[l-2]}$  be the odd indices.

For any  $e_{l-1} \in X_{l-1}$ ,  $A \in X_l$  let  $C_{A,e_{l-1}}^> \in \mathcal{X}_{[l-2]}$  be the set of points  $C_{A,e_{l-1}}^> := \{x_{[l-2]} | \mathbb{K}(A|e_{l-1}, x_{[l-2]}) > \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(A|e_{l-1}, x_{[l-2]})\}$ , and  $C_{A,e_{l-1}}^<$  the obvious analog. Then, defining  $C_{A,e_{l-1}} = C_{A,e_{l-1}}^> \cup C_{A,e_{l-1}}^<$ ,

$$\left(\bigodot_{j \in [l-1]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}\right)(A \times C_{A,e_{l-1}}^> | e_{[l-3]}, e_{l-1}) = 0 \quad (159)$$

$$\left(\bigodot_{j \in [l-1]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}\right)(A \times C_{A,e_{l-1}}^< | e_{[l-3]}, e_{l-1}) = 0 \quad (160)$$

$$\implies \left(\bigodot_{j \in [l-1]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}\right)(A \times C_{A,e_{l-1}} | e_{[l-3]}, e_{l-1}) = 0 \quad (161)$$

$$= \sum_{o_{[l-2]} \in C_{A,e_{l-1}}, \text{odd}} \mathbb{K}(A|e_{l-1}, x_{[l-2]}) \prod_{j \in [l-2]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(o_j) \quad (162)$$

for all  $e_{[l-3]} \in X_{[l-3]_{\text{even}}}$ .

Consider arbitrary  $\mathbb{P}_{\alpha} \in \mathbb{P}_{\square}$ ,  $A \subset X_l$ ,  $C \subset X_{[l-1]}$ :

$$\mathbb{P}_{\alpha}^{\mathbf{X}_{[l]}}(A \times C) = \sum_{x_{[l-2]} \in C} \mathbb{P}_{\square}(A|e_{l-1}, x_{[l-2]}) \prod_{j \in [l-2]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(x_j | o_{[j-2]}, e_{[j-1]}) \mathbb{P}_{\alpha}^{\mathbf{X}_{j+1}|\mathbf{X}_{[j]}}(x_{j+1} | o_{[j]}, e_{[j-1]}) \quad (163)$$

$$= \sum_{x_{[l-2]} \in C_{A,e_{l-1}}} \mathbb{P}_{\square}(A|e_{l-1}, x_{[l-2]}) \prod_{j \in [l-2]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(x_j | o_{[j-2]}, e_{[j-1]}) \mathbb{P}_{\alpha}^{\mathbf{X}_{j+1}|\mathbf{X}_{[j]}}(x_{j+1} | o_{[j]}, e_{[j-1]}) \quad (164)$$

$$+ \sum_{x_{[l-2]} \in C_{A,e_{l-1}}^C} \mathbb{P}_{\square}(A|e_{l-1}, x_{[l-2]}) \prod_{j \in [l-2]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(x_j | o_{[j-2]}, e_{[j-1]}) \mathbb{P}_{\alpha}^{\mathbf{X}_{j+1}|\mathbf{X}_{[j]}}(x_{j+1} | o_{[j]}, e_{[j-1]}) \quad (165)$$

$$= 0 + \sum_{x_{[l-2]} \in C_{A,e_{l-1}}^C} \mathbb{K}(A|e_{l-1}, x_{[l-2]}) \prod_{j \in [l-2]_{\text{odd}}} \mathbb{P}_{\square}^{\mathbf{X}_j|\mathbf{X}_{[j-1]}}(x_j | o_{[j-2]}, e_{[j-1]}) \mathbb{P}_{\alpha}^{\mathbf{X}_{j+1}|\mathbf{X}_{[j]}}(x_{j+1} | o_{[j]}, e_{[j-1]}) \quad (166)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{X}_{[l-1]}} \odot \mathbb{K}(A \times C) \quad (167)$$

$$\implies \mathbb{K} \stackrel{\mathbb{P}_{\square}}{\cong} \mathbb{P}_{\square}^{\mathbf{X}_l|\mathbf{X}_{[l-1]}} \quad (168)$$

□

## 4.7 Comb conditional correspondence

**Theorem ??** (Comb-conditional correspondence). *Given a probability comb  $\{\mathbb{P}_{\square}^{\mathbf{X}_i|\mathbf{X}_{[i-1]}}|i \in \mathbf{X}_{D_{odd}}\}$  and a blind choice  $\alpha$*

$$\mathbb{P}_{\square}^{\mathbf{X}_{D_{odd}}|\mathbf{X}_{D_{even}}} \cong \mathbb{P}_{\alpha} = \mathbb{P}_{\alpha}^{\mathbf{X}_{D_{odd}}|\mathbf{X}_{D_{even}}} \quad (169)$$

*Proof.* Consider  $n \in D$ . The correspondence is immediate for  $n = 1$ :

$$\mathbb{P}_{\square}^{\mathbf{X}_1|\mathbf{X}_0} \stackrel{\mathbb{P}_{\alpha}}{\cong} \mathbb{P}_{\alpha}^{\mathbf{X}_1|\mathbf{X}_0} \quad (170)$$

Suppose for induction the correspondence holds for odd  $n - 2$ . For any blind

$\alpha$

$$\mathbb{P}_\alpha^{X_{[n]}|X_0} = \mathbb{P}_\alpha^{X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_{n-1}|X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_n|X_{[n-1]}} \quad (171)$$

$$= \mathbb{P}_\alpha^{X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_{n-1}|X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_n|X_{[n-1]}} \quad (172)$$

$$= \mathbb{P}_\alpha^{X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_{n-1}|X_{[n-2]}} \cdot \mathbb{P}_\alpha^{X_n|X_{[n-1] \cup \{0\}}} \quad (173)$$

$$= \mathbb{P}_\alpha^{X_{[n-2]_e}} \cdot \mathbb{P}_\alpha^{X_{[n-2]_o}|X_{[n-1]_e}} \cdot \mathbb{P}_\alpha^{X_n|X_{[n-1]}} \quad (174)$$

$$= \mathbb{P}_\alpha^{X_{[n-1]_e}} \cdot \mathbb{P}_\alpha^{X_{[n-2]_o}|X_{[n-1]_e}} \cdot \mathbb{P}_\alpha^{X_n|X_{[n-1]}} \quad (175)$$

and we also have

$$(\mathbb{P}_\alpha^{X_{[n]_{\text{even}}}|X_0} \odot \mathbb{P}_\square^{X_{D_{\text{odd}}}|X_{D_{\text{even}}}})(A \times B|x) = \quad (176)$$

□

## 4.8 Representation of conditional probability models

**Theorem 4.6.** Suppose we have a fundamental probability set  $\Omega$  and a do model  $(\mathbb{P}, D, Y, R)$  such that  $D := (D_i)_{i \in \mathbb{N}}$  and  $Y := (Y_i)_{i \in \mathbb{N}}$ .  $\mathbb{P}$  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

$$\mathbb{P}^{Y|HD} = \begin{array}{c} \text{H} \\ \text{D} \end{array} \begin{array}{c} \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\Pi_i} \end{array} \begin{array}{c} \boxed{\mathbb{P}^{Y_0|HD_0}} \end{array} \text{---} Y_i \quad i \in \mathbb{N} \quad (177)$$

$$\iff \quad (178)$$

$$Y_i \perp\!\!\!\perp_{\mathbb{P}} Y_{\mathbb{N} \setminus i}, D_{\mathbb{N} \setminus i} | HD_i \quad \forall i \in \mathbb{N} \quad (179)$$

$$\wedge \mathbb{P}^{Y_i|HD_i} = \mathbb{P}^{Y_0|HD_0} \quad \forall i \in \mathbb{N} \quad (180)$$

*Proof.* If: By the assumptions of independence and identical conditionals, for any deterministic decision functions  $d, d' \in D$  with equal subsequences  $d_S = d'_T$

$$\mathbb{P}_d^{Y_S|HD}(y|d) = \int_H \prod_{i \in S} \mathbb{P}^{Y_0|HD_0}(y_i|h, d_i) d\mathbb{P}^H(h) \quad (181)$$

$$= \int_H \prod_{i \in T} \mathbb{P}^{Y_0|HD_0}(y_i|h, d'_i) d\mathbb{P}^H(h) \quad \text{by equality of subsequences} \quad (182)$$

$$= \mathbb{P}_{d'}^{Y_T|HD}(y|d) \quad (183)$$

Only if: We have

$$\mathbb{P}^{Y|D} = \begin{array}{c} \text{D} \end{array} \begin{array}{c} \bullet \end{array} \begin{array}{c} \boxed{\mathbb{L}^{D, Y^D}} \end{array} \text{---} Y \quad (184)$$

Also, by contractibility of  $\mathbb{P}^{Y^D}$  and De Finetti's theorem, there is some  $H$  such that

$$\begin{array}{c} \text{---} \end{array} \begin{array}{c} \triangleleft \mathbb{P}^H \end{array} \begin{array}{c} \bullet \end{array} \begin{array}{c} \boxed{\mathbb{P}^{Y_0^D|H}} \end{array} \text{---} Y_i^D \quad i \in \mathbb{N}$$

$$\mathbb{P}^{Y^D} = \quad (185)$$

In particular, let  $Y_{\cdot i}^D := (Y_{ji}^D)_{j \in D}$  and  $Y_{\{i\}^c}^D = (Y_{jk}^D)_{j \in D, k \in \mathbb{N} \setminus \{i\}}$ , and

$$Y_{.i}^D \perp\!\!\!\perp_{\mathbb{P}} Y_{\{i\}^c}^D | H \quad \text{representation theorem} \quad (186)$$

$$Y^D H \perp\!\!\!\perp_{\mathbb{P}} D \quad \text{by Theorem 4.12 and existence of } \mathbb{P}^{Y^D H} \quad (187)$$

$$Y_{.i}^D \perp\!\!\!\perp_{\mathbb{P}} D | Y_{\{i\}^c}^D H \quad \text{weak union on Eq. 187} \quad (188)$$

$$Y_{.i}^D \perp\!\!\!\perp_{\mathbb{P}} D Y_{\{i\}^c}^D | H \quad \text{contraction on Eqs. 186 and 187} \quad (189)$$

$$Y_{.i}^D \perp\!\!\!\perp_{\mathbb{P}} D_{\{i\}^c} Y_{\{i\}^c}^D | H D_i \quad \text{weak union on Eq. 189} \quad (190)$$

$$D_i \perp\!\!\!\perp_{\mathbb{P}} Y_{\{i\}^c}^D D_{\{i\}^c} | H D_i Y_{.i}^D \quad \text{due to conditioning on } D_i \quad (191)$$

$$Y_i^D D_i \perp\!\!\!\perp_{\mathbb{P}} D_{\{i\}^c} Y_{\{i\}^c}^D | H D_i \quad \text{contraction on Eqs. 190 and 191} \quad (192)$$

$$(193)$$

Now, note that  $(Y_i, D_i)$  is a deterministic function of  $(Y_i^D, D_i)$  and  $(Y_{\{i\}^c}^D, D_{\{i\}^c})$  is a deterministic function of  $(Y_{\{i\}^c}^D, D_{\{i\}^c})$ . Therefore

$$Y_i \perp\!\!\!\perp_{\mathbb{P}} D_{\{i\}^c} Y_{\{i\}^c} | H D_i \quad (194)$$

So, by Theorem 4.12,  $\mathbb{P}^{Y_i | H D_i}$  exists and by contractibility of  $\mathbb{P}^{Y^D}$ , for any  $i, j \in \mathbb{N}$

$$\mathbb{P}^{Y_i | H D_i}(y_i | h, d_i) = \mathbb{P}^{Y_{d_i i}^D | H}(y_i | h) \quad (195)$$

$$= \mathbb{P}^{Y_{d_i j}^D | H}(y_i | h) \quad (196)$$

$$= \mathbb{P}^{Y_j | H D_j}(y_i | h, d_i) \quad (197)$$

□

## References

- 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>.
- Ethan D. Bolker. Functions Resembling Quotients of Measures. *Transactions of the American Mathematical Society*, 124(2):292–312, 1966. ISSN 0002-9947. doi: 10.2307/1994401. URL <https://www.jstor.org/stable/1994401>. Publisher: American Mathematical Society.

- 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.
- R.P. Feynman. *The Feynman lectures on physics*. Le cours de physique de Feynman. Intereditions, France, 1979. ISBN 978-2-7296-0030-3. INIS Reference Number: 13648743.
- Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv:1301.6201 [math]*, January 2013. URL <http://arxiv.org/abs/1301.6201>. arXiv: 1301.6201.
- 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>.
- Alan Hájek. What Conditional Probability Could Not Be. *Synthese*, 137(3): 273–323, December 2003. ISSN 1573-0964. doi: 10.1023/B:SYNT.0000004904.91112.16. URL <https://doi.org/10.1023/B:SYNT.0000004904.91112.16>.
- Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In Mikoaj Bojaczek and Alex Simpson, editors, *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing, 2019. ISBN 978-3-030-17127-8.
- Richard C. Jeffrey. *The Logic of Decision*. University of Chicago Press, July 1965. ISBN 978-0-226-39582-1.
- James M. Joyce. *The Foundations of Causal Decision Theory*. Cambridge University Press, Cambridge ; New York, April 1999. ISBN 978-0-521-64164-7.
- Alfred Korzybski. *Science and sanity; an introduction to Non-Aristotelian systems and general semantics*. Lancaster, Pa., New York City, The International Non-Aristotelian Library Publishing Company, The Science Press Printing Company, distributors, 1933. URL <http://archive.org/details/sciencesanityint00korz>.
- Finnian Lattimore and David Rohde. Causal inference with Bayes rule. *arXiv:1910.01510 [cs, stat]*, October 2019a. URL <http://arxiv.org/abs/1910.01510>. arXiv: 1910.01510.
- Finnian Lattimore and David Rohde. Replacing the do-calculus with Bayes rule. *arXiv:1906.07125 [cs, stat]*, December 2019b. URL <http://arxiv.org/abs/1906.07125>. arXiv: 1906.07125.
- Karl Menger. Random Variables from the Point of View of a General Theory of Variables. In Karl Menger, Bert Schweizer, Abe Sklar, Karl Sigmund, Peter Gruber, Edmund Hlawka, Ludwig Reich, and Leopold Schmetterer, editors, *Selecta Mathematica: Volume 2*, pages 367–381. Springer, Vienna, 2003. ISBN 978-3-7091-6045-9. doi: 10.1007/978-3-7091-6045-9\_31. URL [https://doi.org/10.1007/978-3-7091-6045-9\\_31](https://doi.org/10.1007/978-3-7091-6045-9_31).
- 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>.
- Frank P. Ramsey. Truth and Probability. In Horacio Arló-Costa, Vincent F. Hendricks, and Johan van Benthem, editors, *Readings in Formal Epistemology: Sourcebook*, Springer Graduate Texts in Philosophy, pages 21–45. Springer International Publishing, Cham, 2016. ISBN 978-3-319-20451-2. doi: 10.1007/978-3-319-20451-2\_3. URL [https://doi.org/10.1007/978-3-319-20451-2\\_3](https://doi.org/10.1007/978-3-319-20451-2_3).
- 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.
- Jin Tian and Judea Pearl. A general identification condition for causal effects. In *Aaai/iaai*, pages 567–573, July 2002.
- J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Theory of games and economic behavior. Princeton University Press, Princeton, NJ, US, 1944.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England, 1950.

## Appendix: