

# Causal questions are questions that are answered by a function

David Johnston

September 16, 2021

The typical way to construct a probability model is to define a sample space and an ambient probability measure, which together form a “probability space”. This is unweildly for causal models where we are often interested in pulling probability models apart and reassembling them in different ways. We introduce a different approach to building probability models which we call “modular probability” that uses a *modelling context* instead of a sample space, in which random variables are akin to types in a computer program and probability measures and Markov kernels are akin to functions. We illustrate the use of modular probability with examples of decision theoretic causal models, Causal Bayesian Networks and Potential Outcomes models.

## Contents

<b>1</b>	<b>Technical prerequisites</b>	<b>2</b>
1.1	Markov kernels . . . . .	2
1.2	Cartesian and tensor products . . . . .	3
1.3	Delta measures, erase maps, copy maps . . . . .	3
1.4	Products . . . . .	4
1.5	Labeled Markov kernels, conditional probabilities . . . . .	4
1.6	Modelling context . . . . .	5
1.7	Conditional independence . . . . .	7
1.8	Uniqueness of disintegrations . . . . .	7
1.9	Existence of modelling context . . . . .	7
1.10	Standard probability models . . . . .	8
<b>2</b>	<b>See-do models</b>	<b>8</b>
2.1	See-do models and classical statistics . . . . .	9
2.2	Combs . . . . .	9
<b>3</b>	<b>Causal Bayesian Networks</b>	<b>10</b>

<b>4</b>	<b>Potential outcomes with and without counterfactuals</b>	<b>11</b>
4.1	Potential outcomes in see-do models . . . . .	13
4.2	Parallel potential outcomes representation theorem . . . . .	14
<b>5</b>	<b>Appendix:see-do model representation</b>	<b>17</b>

# 1 Technical prerequisites

Our theory makes heavy use of *Markov kernels* or *stochastic functions*, which are taken from probability theory. However, the manner in which we use them is non-standard. The usual way to apply probability theory to model building is to assume we have a probability space  $(\mathbb{P}, \Omega, \mathcal{F})$  with random variables defined as functions with domain  $\Omega$ , and all aspects of the model of interest are supposed to be captured by this. Under our approach, we instead consider components, represented by Markov kernels  $\mathbf{K} : E \rightarrow \Delta(F)$  along with labeled inputs and outputs. The labels do the same job that random variables do in the usual formulation. These components can be composed or broken apart, but we do not assume that there is an overarching probability space from which all components can be derived.

In addition, we introduce a graphical notation for Markov kernels that is the subject of a coherence theorem: two Markov kernels represented by pictures that differ only by planar deformations are identical (Selinger, 2010).

## 1.1 Markov kernels

Markov kernels can be thought of as measurable functions that map to probability distributions. A conditional probability  $\mathbb{P}(Y|X)$ , which maps from values of  $X$  to probability distributions over  $Y$ , and an interventional map  $x \mapsto \mathbb{P}(Y|do(X = x))$  that likewise maps values of  $X$  to probability distributions on  $Y$ , are both Markov kernels.

Our theory is substantially simplified by restricting our attention to discrete sets – that is, sets  $X$  with at most a countable number of elements endowed with the  $\sigma$ -algebra made up of every subset of  $X$ , also called the discrete  $\sigma$ -algebra.

In the discrete setting, we can represent probability distributions as covectors, Markov kernels as matrices and measurable functions as vectors.

Given a set  $X$ , a probability distribution  $\mathbb{P}$  on  $X$  is a covector in  $\mathbb{R}^{|X|}$ , which we will write  $\mathbb{P} := (\mathbb{P}^i)_{i \in X}$ . To be a probability distribution we require

$$0 \leq P_i \leq 1 \quad \forall i \in X \quad (1)$$

$$\sum_i P_i = 1 \quad (2)$$

Given discrete sets  $X$  and  $Y$ , a Markov kernel  $\mathbf{K} : X \rightarrow \Delta(Y)$  is a matrix

in  $\mathbb{R}^{|X| \times |Y|}$ ;  $\mathbf{K} = (K_i^j)_{i \in X, j \in Y}$  where

$$0 \leq K_i^j \leq 1 \quad \forall i, j \quad (3)$$

$$\sum_{i \in X} K_i^j = 1 \quad \forall j \quad (4)$$

Rows of Markov kernel are probability distributions:  $\mathbf{K}_x := (K_x^j)_{j \in Y}$ . Alternatively, we can consider probability distributions to be Markov kernels with one row.

Graphically, we represent a Markov kernel as a box and a probability distribution as a triangle:

$$\mathbf{K} := \boxed{\mathbf{K}} \quad (5)$$

$$\mathbb{P} := \triangleleft \mathbf{K} \quad (6)$$

## 1.2 Cartesian and tensor products

The Cartesian product  $X \times Y := \{(x, y) | x \in X, y \in Y\}$ .

Given kernels  $\mathbf{K} : W \rightarrow Y$  and  $\mathbf{L} : X \rightarrow Z$ , the tensor product  $\mathbf{K} \otimes \mathbf{L} : W \times X \rightarrow \Delta(Y \times Z)$  is defined by  $(\mathbf{K} \otimes \mathbf{L})_{(w, x)}^{(y, z)} := K_w^y L_x^z$ .

Graphically, the tensor product is represented by parallel juxtaposition:

$$\mathbf{K} \otimes \mathbf{L} := \boxed{\mathbf{K}} \quad \boxed{\mathbf{L}} \quad (7)$$

## 1.3 Delta measures, erase maps, copy maps

The Iverson bracket  $\llbracket \cdot \rrbracket$  evaluates to 1 if  $\cdot$  is true and 0 otherwise.

For any  $X$  and any  $x \in X$ ,  $\delta[x]$  is the probability measure defined by  $\delta[x]^i = \llbracket x = i \rrbracket$ . The identity map  $\text{Id}[X] : X \rightarrow \Delta(X)$  is given by  $x \mapsto \delta[x]$ .

Graphically, the identity map is a bare line:

$$\text{Id}[X] := \text{---} \quad (8)$$

The erase map  $*[A] : A \rightarrow \{1\}$  is the map  $*[A]^i = 1$ . It is the unique Markov kernel with domain  $A$  and only one column.

Graphically, the stopper is a fuse:

$$\text{Id}[X] := \text{---} * \quad (9)$$

The copy map  $\Upsilon[X] : X \rightarrow \Delta(X \times X)$  is the Markov kernel defined by  $\Upsilon_x := \delta_x \otimes \delta_x$ . Graphically it is a fork with a dot at the point where it splits:

$$\Upsilon[X] := \text{---} \curvearrowright \quad (10)$$

## 1.4 Products

Two Markov kernels  $\mathbf{L} : X \rightarrow \Delta(Y)$  and  $\mathbf{M} : Y \rightarrow \Delta(Z)$  have a product  $\mathbf{LM} : X \rightarrow \Delta(Z)$  given by the usual matrix-matrix product:  $\mathbf{LM}_x^z = \sum_y \mathbf{L}_x^y \mathbf{M}_y^z$ . Graphically, we write represent products by joining kernel wires together:

$$\mathbf{LM} := \boxed{\mathbf{K}} \text{---} \boxed{\mathbf{M}} \quad (11)$$

## 1.5 Labeled Markov kernels, conditional probabilities

A labeled Markov kernel  $(\mathbf{K}, \mathbf{A}_C, \mathbf{B}_D)$  is a Markov kernel  $\mathbf{K} : X \rightarrow \Delta(Y)$  along with a sequence of *domain labels*  $\mathbf{A}_C := (\mathbf{A}_i)_{i \in C}$  and *codomain labels*  $\mathbf{B}_D := (\mathbf{B}_i)_{i \in D}$  such that

- Each label  $\mathbf{A}_i$  has an associated discrete set  $A_i$  (and similarly  $\mathbf{B}_i$  is associated with  $B_i$ )
- $X = \times_{i \in C} A_i$  and  $Y = \times_{i \in D} B_i$

Repeated labels are okay only if there's a valid diagrammatic representation of  $\mathbf{K}$  such that the repeated labels are connected by a wire with no boxes in between (copy map dots are OK)

A labeled probability distribution  $\mathbb{P} \in \Delta(Y)$  comes with a sequence of codomain labels  $(\mathbf{B}_i)_{i \in D}$  only, satisfying  $Y = \times_{i \in D} B_i$ .

A conditional probability  $\mathbb{L}[\mathbf{A}_C | \mathbf{B}_D]$  is a labeled kernel  $(\mathbf{K}, \mathbf{A}_C, \mathbf{B}_D)$  along with an *ambient conditional probability* (Definition 1.4)  $\mathbb{L}$ .

Graphically, we place the labels on the wires of a conditional probability and the name of the ambient conditional probability in the centre of the box:

$$\mathbb{L}[\mathbf{B}_1 \mathbf{B}_2 | \mathbf{A}_1 \mathbf{A}_2] := \begin{matrix} \mathbf{A}_1 & & \mathbf{B}_1 \\ & \boxed{\mathbb{L}} & \\ \mathbf{A}_2 & & \mathbf{B}_2 \end{matrix} \quad (12)$$

A sequence of labels is itself a label, so we can also bundle wires and their corresponding labels together:

$$\mathbb{L}[\mathbf{B}_1 \mathbf{B}_2 | \mathbf{A}_1 \mathbf{A}_2] = (\mathbf{A}_1, \mathbf{A}_2) \text{---} \boxed{\mathbb{L}} \text{---} (\mathbf{B}_1, \mathbf{B}_2) \quad (13)$$

Because it saves a lot of space, we will generally hold to the convention that a label  $\mathbf{X}$  is associated with the set  $X$ . However, this convention sometimes fails, for example when we have two labels  $\mathbf{X}_1$  and  $\mathbf{X}_2$  that are associated with the same set – in such cases, we will explicitly define the relationship.

The trivial label  $*$  always corresponds to the 1-element set  $\{*\}$ . Because  $\{*\} \times A$  is isomorphic to  $A$  for any  $A$ , we can consider any label sequence to be isomorphic to the same label sequence with any number of copies of the trivial

label appended. A trivial sequence of labels is simply a sequence of labels that consists entirely of trivial labels.

If two conditional probabilities  $\mathbb{L}[A_C|B_D]$  and  $\mathbb{K}[A_C|B_D]$  share the same kernel, we will say  $\mathbb{L}[A_C|B_D] \stackrel{krn}{=} \mathbb{K}[A_C|B_D]$

## 1.6 Modelling context

A *modelling context*  $\mathcal{M}$  is a collection of conditional probabilities. It can be thought of as a namespace. We place the following requirements on elements of  $\mathcal{M}$ :

- If  $\mathbb{K}[XZ|Y] \in \mathcal{M}$  and  $\mathbb{L}[XW|V] \in \mathcal{M}$ , the label  $X$  is associated with the same set  $X$  both conditional probabilities
- If  $\mathbb{K}[X|Y] \in \mathcal{M}$  and  $\mathbb{K}[ZX|Y] \in \mathcal{M}$  then  $\mathbb{K}[X|Y]$  is a marginal (Definition 1.2) of  $\mathbb{K}[ZX|Y]$
- If  $\mathbb{K}[XZ|Y] \in \mathcal{M}$  and  $\mathbb{K}[X|YZ] \in \mathcal{M}$  then  $\mathbb{K}[X|YZ]$  is a disintegration (Definition 1.3) of  $\mathbb{K}[XZ|Y]$
- If  $\mathbb{K}[X|Y] \in \mathcal{M}$  then all marginals of  $\mathbb{K}[X|Y]$  are also in  $\mathcal{M}$
- If  $\mathbb{K}[X|Y] \in \mathcal{M}$  then all disintegrations of  $\mathbb{K}[X|Y]$  are also in  $\mathcal{M}$
- If  $\mathbb{K}[ZX|YQ] \in \mathcal{M}$  is extendable by  $\mathbb{L}[W|XQR] \in \mathcal{M}$  then  $\mathbb{K}[ZX|YQ] \Rightarrow \mathbb{L}[W|XQR] \in \mathcal{M}$
- If  $\mathbb{K}[X|Y] \in \mathcal{M}$  is an ambient conditional probability (Definition 1.4), then  $\mathbb{K}[X|Y] = \mathbb{K}$

Given two conditional probabilities from a modelling context, we can extend conditional probabilities by matching labels on inputs of one with the labels on the inputs and outputs of the other. Roughly speaking, we can extend a conditional probability with a second if the second comes “after” the first.

Given two conditional probabilities  $\mathbb{K}[B_C|A_D]$  and  $\mathbb{L}[F_H|E_G]$ , we say  $\mathbb{K}[B_C|A_D]$  is extendable by  $\mathbb{L}[F_H|E_G]$  if there is no label in  $A_D$  that matches a label in  $F_H$  (so no outputs from the “after” conditional probability match inputs from the “before” conditional probability) and there is no label in  $F_H$  that matches a label in  $B_C$  (so no shared outputs).

**Definition 1.1** (extension). Consider two arbitrary conditional probabilities  $F$  and  $S$  in  $\mathcal{M}$  where  $F$  is before  $S$ . Let  $Z$  as the (possibly trivial) sequence of all labels that appear only in the output of  $F$ ,  $X$  the sequence of all labels that appear in the output of  $F$  and the input of  $S$ ,  $Y$  the sequence of all labels that appear only in the input of  $F$ ,  $Q$  the sequence of all labels shared by the inputs of  $F$  and  $S$ ,  $R$  the sequence of all labels that appear only in the input of  $S$  and  $W$  the sequence of all labels that appear only in the output of  $S$ . Given the assumption that  $F$  is before  $S$ , we can in general write  $F = \mathbb{K}[ZX|YQ]$  and

$S = \mathbb{L}[W|XQR]$  for some model names  $\mathbb{K}$  and  $\mathbb{L}$  (the necessary condition for this is that the label sequences are exhaustive).

Then Equations 15 and 16 are equivalent definitions of extension:

$$\mathbb{K}[ZX|YQ] \Rightarrow \mathbb{L}[W|XQR] := \mathbb{M}[ZXW|YQR] \quad (14)$$

$$:= \begin{array}{c} \text{Y} \text{---} \boxed{\mathbb{K}} \text{---} \text{Z} \\ \text{Q} \text{---} \bullet \text{---} \text{X} \\ \text{R} \text{---} \text{---} \boxed{\mathbb{L}} \text{---} \text{W} \end{array} \quad (15)$$

$$\mathbb{M}[ZXW|YQR]_{yqr}^{zxw} = \mathbb{K}[ZX|YQ]_{yq}^{zx} \mathbb{L}[W|XQR]_{xqr}^w \quad (16)$$

Here we assume that we have some way of choosing the ambient conditional probability  $\mathbb{M}$ . If the two kernels are appropriate marginals and disintegrations of the same ambient conditional probability  $\mathbb{K}$  (see Definition 1.3), then the extension has ambient conditional probability  $\mathbb{K}$  also. Otherwise, we can always consistently choose  $\mathbb{M} := \mathbb{M}[ZXW|YQR]$ .

I don't know if there are other cases where it would be sensible to make a different choice. The intuition I have is that if two ambient conditional probabilities are equal in kernel and label set then they should be equal, but not other conditional probabilities.

Equation 15 can be broken down to the product of four Markov kernels, each of which is itself a tensor product of a number of other Markov kernels:

$$\mathbb{M}[ZXW|YQR] = \left[ \begin{array}{c} \text{Y} \text{---} \\ \text{Q} \text{---} \bullet \text{---} \\ \text{R} \text{---} \end{array} \right] \left[ \begin{array}{c} \boxed{\mathbb{K}} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \bullet \text{---} \\ \text{---} \end{array} \right] \left[ \begin{array}{c} \text{---} \text{Z} \\ \text{---} \text{X} \\ \boxed{\mathbb{L}} \text{---} \text{W} \end{array} \right] \quad (17)$$

$$(18)$$

$\mathbb{M}[ZXW|YQR]$  is itself a conditional probability with some ambient conditional probability  $\mathbb{M}$  and is an element of  $\mathcal{M}$ .

Prove that  $\Rightarrow$  is associative

**Definition 1.2** (marginal). Given a conditional probability  $\mathbb{K}[XY|W]$ , the marginal  $\mathbb{K}[X|W]$  is defined as

$$\mathbb{K}[X|W] := \text{W} \text{---} \boxed{\mathbb{K}} \text{---} \text{X} \quad (19)$$

$$\mathbb{K}[X|W]_w^x = \sum_{y \in Y} \mathbb{K}[XY|W]_w^{xy} \quad (20)$$

**Definition 1.3** (disintegration).  $\mathbb{K}[Y|XW]$  is a disintegration of  $\mathbb{K}[XY|W]$  if

$$\mathbb{K}[X|W] \Rightarrow \mathbb{K}[Y|XW] = \mathbb{K}[XY|W] \quad (21)$$

Any Markov kernel  $\mathbf{L}$  with the property

$$\mathbf{L}_{xw}^y = \frac{\mathbb{K}[\mathbf{XY}|\mathbf{W}]_w^{xy}}{\sum_{x \in X} \mathbb{K}[\mathbf{XY}|\mathbf{W}]_w^{xy}} \quad \forall w, y : \text{the denominator is positive} \quad (22)$$

is a version of  $\mathbb{K}[\mathbf{Y}|\mathbf{XW}]$ .

**Definition 1.4** (ambient conditional probability). A conditional probability  $\mathbb{K}[\mathbf{Y}|\mathbf{X}]$  is an *ambient conditional probability* relative to  $\mathcal{M}$  if there is no other conditional probability in  $\mathcal{M}$  such that  $\mathbb{K}[\mathbf{Y}|\mathbf{X}]$  is either a marginal or a disintegration of this conditional probability.

Recall that, if  $\mathbb{K}[\mathbf{Y}|\mathbf{X}]$  is an ambient conditional probability, then  $\mathbf{K}[\mathbf{Y}|\mathbf{X}] = \mathbf{K}$ .

## 1.7 Conditional independence

Given  $\mathbb{K}[\mathbf{X}|\mathbf{WZ}]$  in general we have no definition of  $\mathbb{K}[\mathbf{X}|\mathbf{Z}]$ . However, we can define such a “conditional probability” if we have the additional fact that  $\mathbf{X}$  is independent of  $\mathbf{W}$  given  $\mathbf{Z}$  relative to  $\mathbb{K}$ .

Given  $\mathbb{K}[\mathbf{X}|\mathbf{WZ}]$  we say  $\mathbf{X}$  is independent of  $\mathbf{W}$  given  $\mathbf{Z}$  relative to  $\mathbb{K}$ , notated  $\mathbf{X} \perp_{\mathbb{K}} \mathbf{W}|\mathbf{Z}$  iff  $\mathbb{K}[\mathbf{X}|\mathbf{WZ}]_{wz}^x = \mathbb{K}[\mathbf{X}|\mathbf{WZ}]_{w'z}^x$  for all  $w, w' \in W$ ,  $x \in X$  and  $z \in Z$ .

Given  $\mathbb{K}[\mathbf{X}|\mathbf{WZ}]$  such that  $\mathbf{X} \perp_{\mathbb{K}} \mathbf{W}|\mathbf{Z}$ , we define  $\mathbb{K}[\mathbf{X}|\mathbf{Z}]$  to be any kernel satisfying  $\mathbb{K}[\mathbf{X}|\mathbf{Z}]_z^x = \mathbb{K}[\mathbf{X}|\mathbf{DZ}]_{dz}^x$  for all  $x, z, d$ .

## 1.8 Uniqueness of disintegrations

Every conditional probability  $\mathbf{K}[\mathbf{X}|\mathbf{Y}]$  is unique up to an equivalence class defined with respect to the ambient conditional probability  $\mathbb{K}$ .

Proof sketch: if it is an ambient conditional probability, then it is unique. If not, it is obtained from an ambient conditional probability by a sequence of marginalisations and disintegrations. Defining the equivalence class to be “equal up to measure 0 sets”, marginalisations and disintegrations are both unique.

## 1.9 Existence of modelling context

Take a collection of Markov kernels and give them label sets consistent with their type signatures and respecting the rule that identical labels require identical spaces. Add all the recursive disintegrations + marginals. Add all valid extensions assigning a new model name for any result not already in the modelling context. Add all recursive disintegrations and marginals of valid extensions, etc.

Then: disintegration, marginalisation, extension operations all preserve label consistency rules. By construction, marginals, disintegrations and extensions are included. Marginalisation + disintegration preserves uniqueness of ambient conditional probability. Extension + assigning a new model name also preserves uniqueness of ambient conditional probability.

## 1.10 Standard probability models

The operation of combining two conditional probabilities which do not share a model name and obtaining a conditional probability relative with a new model name is unique to our approach. The standard approach to probability modelling features an ambient probability distribution defined on a sample space, along with “labels” that each correspond to a measurable functions on the sample space. With this setup, we can define marginals and disintegrations with respect to any sequences of labels. It is an open question whether there is a way to construct a modelling context with a single model that is equivalent to a standard probability model.

## 2 See-do models

Modular probability is useful when we want to combine different Markov kernels in such a way that “variables” refer to something consistent even though they don’t necessarily have a unique distribution. The first example we will present is using modular probability to model decision problems.

Suppose we will be given an observation  $x \in X$  and in response to this we can select any decision or stochastic mixture of decisions from a set  $D$ ; that is we can choose a “strategy” as any Markov kernel  $\mathbf{S}_\alpha : X \rightarrow \Delta(D)$ . We are interested in forecasting some consequences that take values in some set  $Y$ , and comparing the forecasts for different strategy choices so as to choose a best strategy.

How can we model this? One way to proceed is as follows: Define a model context  $\mathcal{M}$  to which we add the conditional probabilities mentioned hereafter. For each strategy  $\mathbf{S}_\alpha[D|X]$ , our forecast will be represented by some joint probability in  $\mathbb{P}_\alpha[XDY|H]$  where  $H$  is associated with a set of hypotheses  $H$  representing different choices that we think might be reasonable to make that may lead to different forecasts. Because observations come before we execute our strategy, we assume that  $\mathbb{P}_\alpha[X|H] = P_\beta[X|H]$  for all  $\alpha, \beta$ . Our chosen strategy is the probability of  $D$  given  $X$ :  $\mathbb{P}_\alpha[D|X] \stackrel{krn}{=} \mathbf{S}_\alpha[D|X]$ . Finally, our forecast of  $Y$  is the same for all strategies holding the observations, the decision and the hypothesis fixed:  $\mathbb{P}_\alpha[Y|HD] = P_\beta[Y|HD]$  for all  $\alpha, \beta$ .

Under these assumptions, there exists  $\mathbb{T}[XY|HD] \in \mathcal{M}$  with  $X \perp\!\!\!\perp_{\mathbb{T}} D|H$  such that for all  $\alpha$ ,

$$\mathbb{P}_\alpha[XDY|H] \stackrel{krn}{=} \mathbb{T}[X|H] \Rightarrow \mathbf{S}_\alpha[D|X] \Rightarrow \mathbb{T}[Y|XHD] \quad (23)$$

The proof is given in Appendix 5. Note that  $\mathbb{T}[X|H]$  exists by virtue of the fact  $X \perp\!\!\!\perp_{\mathbb{T}} D|H$ . While this independence is what enables Equation 23, in general  $X \not\perp\!\!\!\perp_{\mathbb{P}_\alpha} D|H$ , so  $\mathbb{T}$  cannot be a disintegration of  $\mathbb{P}_\alpha$ . Modular probability allows us to specify  $\mathbb{T}$ , which we call a *see-do model*, as a partial forecast to be completed with a strategy  $\mathbf{S}_\alpha$  while also being able to use consistent names for variables that represent the same things (observations, decisions, consequences,



hypotheses) whether their distributions are given by  $\mathbb{P}_\alpha, \mathbb{T}$ , which are mutually incompatible conditional probabilities.

## 2.1 See-do models and classical statistics

A *statistical model* (or *statistical experiment*) is a collection of probability distributions indexed by some set  $\Theta$ . We can observe that  $\{\mathbb{T}[X|H]_h\}_{h \in H}$  is a collection of probability distributions indexed by  $H$ .

In statistical decision theory, as introduced by Wald (1950), we are given a statistical experiment  $\{\mathbb{P}_\theta \in \Delta(X)\}_\Theta$ , a decision set  $D$  and a loss  $l : \Theta \times D \rightarrow \mathbb{R}$ . A strategy  $\mathbb{S}_\alpha : X \rightarrow \Delta(D)$  is evaluated according to the risk functional  $R(\theta, \mathbb{S}_\alpha) = \sum_{x \in X} \sum_{d \in D} \mathbb{P}_\theta^x(\mathbb{S}_\alpha)_x^d l(h, d)$ .

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

$$\mathbb{E}_{\mathbb{P}_\alpha[XDY|H]}[Y] = \sum_{x \in X} \sum_{d \in D} \sum_{y \in Y} y (\mathbb{T}[X|H] \Rightarrow \mathbb{S}_\alpha[D|X] \Rightarrow \mathbb{T}[Y|XHD])_h^{xdy} \quad (24)$$

$$= \sum_{x \in X} \sum_{d \in D} \sum_{y \in Y} \mathbb{T}[X|H]_h^x \mathbb{S}_\alpha[D|X]_x^d \mathbb{T}[Y|HD]_{h,d}^y \quad (25)$$

$$= \sum_{x \in X} \sum_{d \in D} \mathbb{T}[X|H]_h^x (\mathbb{S}_\alpha)_x^d l(h, d) \quad (26)$$

$$= R(h, \mathbb{S}_\alpha) \quad (27)$$

That is, if we are given a see-do model where we interpret  $\mathbb{T}[X|H]$  as a statistical experiment and  $Y$  as a reversed utility, the expectation of the utility under the strategy forecast given in equation 23 is the risk of that strategy under hypothesis  $h$ .

## 2.2 Combs

The see-do model  $\mathbb{T}[XY|HD]$  is known as a *comb*. This structure was introduced by Chiribella et al. (2008) in the context of quantum circuit architecture, and Jacobs et al. (2019) adapted the concept to causal modelling.

A comb is a Markov kernel with a “hole” in it. We combine the see-do model with a strategy by putting the strategy “in the middle” of the see-do model (Equation 23), rather than attaching it to one end. While it is not a well-formed diagram in the language described in this paper, we can visualise combs as Markov kernels with holes:

$$\mathbb{T}[XY|HD] = \begin{array}{c} \text{H} \text{---} \boxed{\text{T}} \text{---} \text{X} \text{---} \text{D} \text{---} \boxed{\text{T}} \text{---} \text{Y} \\ \text{---} \text{---} \text{---} \end{array} \quad (28)$$

$$= \begin{array}{c} \text{H} \text{---} \boxed{\text{T}} \text{---} \text{X} \text{---} \text{D} \text{---} \boxed{\phantom{\text{T}}} \text{---} \text{Y} \\ \text{---} \text{---} \text{---} \end{array} \quad (29)$$

We can take any strategy  $\mathbb{S}_\alpha[D|X]$  and drop it into the “hole” in 29 (as described in Equation 23) to get a forecast of the outcome of that strategy.

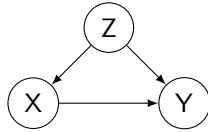
### 3 Causal Bayesian Networks

Causal Bayesian Networks are also see-do models. A Causal Bayesian Network posits a set of observational probability distributions, for example  $\{\mathbb{P}_h(X, Y)\}_H$ , and a set of interventional distributions, for example  $\{\mathbb{P}_h(X, Y|do(X = x))\}_{x \in X, h \in H}$ . Here we use notation similar to typical notation used for Causal Bayesian Networks and don’t intend for these to necessarily be elements of any modelling context. For simplicity, we will consider a Causal Bayesian Network with only hard interventions on  $X$ .

We can consider this an instance of a see-do model. To do so, we need to distinguish observation and intervention variables - let the former retain the labels  $X, Y$  and call the latter  $X', Y'$ . Let  $D = \{do(X = x)\}_{x \in X}$ . Then a Causal Bayesian Network can be considered a see-do model  $\mathbb{T}[XX'Y'Y|HD]$  by identifying  $\mathbb{T}[XY|H]_h := \mathbb{P}_h(X, Y)$  and  $\mathbb{T}[X'Y'|HD]_{h, do(X=x)} := P_h(X, Y|do(X = x))$ .

We need to rename the consequence variables because otherwise we would have  $\mathbb{T}[XXYY|HD]$  and the two  $X$ ’s and the two  $Y$ ’s would be deterministically equal by the “identical labels” rule

We can say a bit more about Causal Bayesian Networks. Suppose we have the network



Then, letting  $\mathbb{T}[XYZ|H]$  be the observational “see” model and  $\mathbb{T}[X'Y'Z'|HD]$  be the interventional “do” model with  $D$  the set of interventions  $\{do(X = x)\}_{x \in X}$  where we write  $x := do(X = x)$  for short, then we know by the backdoor adjustment rule that  $\mathbb{T}[X'Y'Z'|HD]_{hx}^{x'yz} \stackrel{krn}{=} \mathbb{T}[Z|H]_h^z \delta[x]^{x'} \mathbb{T}[Y|XZH]_{hx'z}^y$ .

Let  $\mathbb{U}[ZY|XH] = \mathbb{T}[Z|H] \Rightarrow \mathbb{T}[Y|XZH]$ , call  $\mathbb{T}[X|H]$  the “observational strategy” and  $\mathbb{D}_x[X|D]_x^{x'} \stackrel{krn}{=} \delta[x]^{x'}$  the interventional strategies for all  $x \in X$ . Then we have

$$\mathbb{T}[XYZ|H] = \mathbb{U}[Z|H] \Rightarrow \mathbb{T}[X|H] \Rightarrow \mathbb{U}[Y|XHZ] \quad (30)$$

$$\mathbb{T}[X'Y'Z'|HD] \stackrel{krn}{=} \mathbb{U}[Z|H] \Rightarrow \mathbb{D}[X|D] \Rightarrow \mathbb{U}[Y|XHZ] \quad (31)$$

So this simple example of a Causal Bayesian network is a “nested comb” where the outer comb  $\mathbb{T}[XYZX'Y'Z'|HD]$  is the “see” and “do” models, which are themselves generated by the inner comb  $\mathbb{U}[ZY|XH]$  with different choices  $\mathbb{T}[X|H]$  and  $\mathbb{D}[X|D]$  for the insert.

This is a simple example, but Jacobs et al. (2019) has used an “inner comb” representation of a general class of Causal Bayesian Networks to prove a sufficient identification condition which is itself slightly more general than the identification condition given by Tian and Pearl (2002).

## 4 Potential outcomes with and without counterfactuals

I still need to convert this section to “modular probability”, which I think will simplify it. The short story is that potential outcomes models are nested combs just like CBNs, and we can define potential outcome combs that accept “counterfactual” inserts as well as potential outcome combs that accept ordinary decision/observation inserts, the latter is all anyone actually uses and can be given a perfectly satisfactory non-counterfactual interpretation.

Potential outcomes is a widely used approach to causal modelling characterised by its use of “potential outcome” random variables. Potential outcome random variables are typically noted for being given counterfactual interpretations. For example, suppose have something we want to model, call it TYT (“The Y Thing”), which we represent with a variable  $Y$ . Suppose we want to know how TYT behaves under different regimes 0 and 1 under which we want to know about TYT, and we use a variable  $W$  to indicate which regime holds at a given point in time. A potential outcomes model will introduce the two additional “potential outcome” variables ( $Y(0), Y(1)$ ). What these variables represent can be given a counterfactual interpretation like “ $Y(0)$  represents what TYT would be under regime 0, whether or not regime 0 is the actual regime” and similarly “ $Y(1)$  represents what TYT would be under regime 1, whether or not regime 1 is the actual regime”. Note that we say “what TYT would be” rather than “what  $Y$  would be” as “what would  $Y$  be if  $W$  was 0 if  $W$  was actually 1” is not a question we can ask of random variables, but it is one that might make sense for the things we use random variables to model.

This is a key point, so it is worth restating: the assumption that potential outcome variables agree with “the value TYT would take” under fixed regimes regardless of the “actual” value of the regime seems to be a critical assumption that distinguishes potential outcome variables from arbitrary random variables

that happen to take values in the same space as  $Y$ . However, this assumption can only be stated by making reference to the informally defined “TYT” and the informal distinction between the supposed and the actual value of the regime.

The potential outcomes framework features other critical assumptions that relate potential outcome variables to things that are only informally defined. For example, Rubin (2005) defines the *Stable Unit Treatment Value Assumption* (SUTVA) as:

SUTVA (stable unit treatment value assumption) [...] 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

“Versions of treatments” do not appear within typical potential outcomes models, so this is also an assumption about how “the thing we are trying to model” behaves rather than an assumption stated within the model.

Given informal assumptions like this, one may be motivated to “formalize” them. More specifically, one might be motivated to ask whether there is some larger class of models that, under conditions corresponding to the informal conditions above yield regular potential outcome models?

I have a vague intuition here that you always need some kind of assumption like “my model is faithful to the real thing”, but if you are stating fairly specific conditions in English you should also be able to state them mathematically. Among other reasons, this is useful because it’s easier for other people to know what you mean when you state them.

The approach we have introduced here, motivated by decision problems, has in the past been considered a means of avoiding counterfactual statements, which has been considered a positive by some (Dawid, 2000) and a negative by others:

[...] Dawid, in our opinion, incorrectly concludes that an approach to causal inference based on “decision analysis” and free of counterfactuals is completely satisfactory for addressing the problem of inference about the effects of causes. (Robins and Greenland, 2000)

It may be surprising to some, then, that we can use see-do models to formally state these key assumptions associated with potential outcomes models. Furthermore, we will argue that potential outcomes are typically a strategy to motivate inductive assumptions in see-do models, and we will show that the counterfactual interpretation is unnecessary for this purpose.

## 4.1 Potential outcomes in see-do models

A basic property of potential outcomes models is the relation between variables representing actual outcomes and variables representing potential outcomes, which was stated informally in the opening paragraph of this section.

In the following definition,  $Y(W) = (Y(w))_{w \in W}$ .

**Definition 4.1** (Potential outcomes). Given a Markov kernel space  $(\mathbf{K}, E, F)$ , a collection of variables  $\{Y, Y(W), W\}$  where  $Y$  and  $Y(W)$  are random variables and  $W$  could be either a state or a random variable is a *potential outcome submodel* if  $\mathbf{K}[Y|WY(W)]$  exists and  $\mathbf{K}[Y|WY(W)]_{ij_1j_2 \dots j_{|W|}} = \delta[j_i]$ .

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

We allow  $X$  to be a state or a random variable to cover the cases where potential outcomes models feature as submodels of observation models (in which case  $X$  is a random variable) or as submodels of consequence models (in which case  $X$  may be a state variable).

As an aside that we could define stochastic potential outcomes if we allow the variables  $Y(x)$  to take values in  $\Delta(Y)$  rather than in  $Y$ , and then require  $\mathbf{K}[Y|XY(X)]_{ij_1j_2 \dots j_{|X|}} = j_i$  (where  $j_i$  is an element of  $\Delta(Y)$ ). This is more complex to work with and rarely seen in practice, but it is worth noting that Definition 4.1 can be generalised to cover models where  $Y(x)$  describes the value  $Y$  would take if  $X$  were  $x$  *with uncertainty*.

An arbitrary see-do model featuring potential outcome submodels does not necessarily allow for the formal statement of the counterfactual interpretation of potential outcomes. Here we use TYT (“the actual thing”) and “regime” to refer to the things we are actually trying to model. We require that  $Y \stackrel{a.s.}{=} Y(w)$  conditioned on  $W = w$ . If we add an interpretation to this model saying  $Y$  represents TYT and  $W$  represents the regime, then we have “for all  $w$ ,  $Y(w)$  is equal to  $Y$  which represents TYT under the regime  $w$ ”. However, this does not guarantee that our model has anything that reasonably represents “what TYT would be equal to under supposed regime  $w$  if the regime is actually  $w'$ ”.

We propose *parallel potential outcome submodels* as a means of formalising statements about what how TYT behaves under “supposed” and “actual” regimes:

**Definition 4.2** (Parallel potential outcomes). Given a Markov kernel space  $(\mathbf{K}, E, F)$ , a collection of variables  $\{Y_i, Y(W), W_i\}$ ,  $i \in [n]$ , where  $Y_i$  and  $Y(W)$  are random variables and  $W_i$  could be either a state or random variables is a *parallel potential outcome submodel* if  $\mathbf{K}[Y_i|W_iY(W)]$  exists and  $\mathbf{K}[Y_i|W_iY(W)]_{kj_1j_2 \dots j_{|W|}} = \delta[j_k]$ .

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

A parallel potential outcomes model features a sequence of  $n$  “parallel” outcome variables  $Y_i$  and  $n$  “regime proposals”  $W_i$ , with the property that if the

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

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

## 4.2 Parallel potential outcomes representation theorem

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

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

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

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

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

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

**Definition 4.3** (Exchangeability). Given a Markov kernel space  $(\mathbf{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $Y_i$  random variables is *exchangeable* conditional on  $Z$  if, defining  $Y_{[n]} = (Y_i)_{i \in [n]}$  and  $D_{[n]} = (D_i)_{i \in [n]}$ ,  $\mathbf{K}[Y_{[n]}|D_{[n]}Z]$  exists and for any bijection  $\pi : [n] \rightarrow [n]$   $\mathbf{K}[Y_{\pi([n])}|D_{\pi([n])}Z] = \mathbf{K}[Y_{[n]}|D_{[n]}Z]$ .

**Definition 4.4** (Extension). Given a Markov kernel space  $(\mathbf{K}, E, F)$ ,  $(\mathbf{K}', E', F')$  is an *extension* of  $(\mathbf{K}, E, F)$  if there is some random variable  $X$  and some state variable  $U$  such that  $\mathbf{K}'[X|U]$  exists and  $\mathbf{K}'[X|U] = \mathbf{K}$ .

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

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

**Definition 4.5** (Extendably exchangeable). Given a Markov kernel space  $(\mathbf{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  and a state variable  $Z$  with  $Y_i$  random variables is *extendably exchangeable* if there exists an extension of  $\mathbf{K}$  with respect to which  $((D_i, Y_i))_{i \in \mathbb{N}}$  is exchangeable conditional on  $Z$ .

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

Deterministically reproducible sequences have the property that repeating the same decision gets the same response with probability 1. This could be a model of an experiment that exhibits no variation in results (e.g. every time I put green paint on the page, the page appears green), or an assumption about collections of “what-ifs” (e.g. if I went for a walk an hour ago, just as I actually did, then I definitely would have stubbed my toe, just like I actually did). Incidentally, many consider that this assumption is false concerning what-if questions about things that exhibit quantum behaviour.

**Definition 4.6** (Deterministically reproducible). Given a Markov kernel space  $(\mathbf{K}, E, F)$ , a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $Y_i$  random variables is *deterministically reproducible* conditional on  $Z$  if  $n \geq 2$ ,  $\mathbf{K}[Y_{[n]}|D_{[n]}Z]$  exists and  $\mathbf{K}[Y_{\{i,j\}}|D_{\{i,j\}}Z]_{kk}^{lm} = \mathbb{I}[l = m]\mathbf{K}[Y_i|D_iZ]_k^l$  for all  $i, j, k, l, m$ .

**Theorem 4.7** (Potential outcomes representation). *Given a Markov kernel space  $(\mathbf{K}, E, F)$  along with a sequence of variables  $((D_i, Y_i))_{i \in [n]}$  with  $n \geq 2$  and a conditioning variable  $Z$ ,  $(\mathbf{K}, E, F)$  can be extended with a set of variables  $Y(D) := (Y(i))_{i \in D}$  such that  $\{Y_i, Y(D), D_i\}$  is a parallel potential outcome submodel if and only if  $((D_i, Y_i))_{i \in [n]}$  is extendably exchangeable and deterministically reproducible conditional on  $Z$ .*

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

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

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

$$\mathbf{K}[X_n|D_nZ]_{dz}^x = \sum_{y \in Y^{|D|}} \mathbf{K}[X_{[n]}|D_nZY(D)]_{dzy}^x \mathbf{K}[Y(D)|D_{[n]}Z]_{dz}^y \quad (32)$$

$$= \sum_{y \in Y^{|D|}} \prod_{i \in [n]} \delta[y_{d_i}]^{x_i} \mathbf{K}[Y(D)|D_nZ]_{dz}^y \quad (33)$$

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



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

$$\mathbf{K}[X_{[n]}|D_{[n]}Z]_{dz}^x = \sum_{y \in Y^{[D]}} \prod_{i \in Q_d} \delta[y_{d_i}]^{x_i} \mathbf{K}[Y(D)|D_{[n]}Z]_{dz}^y \quad (34)$$

$$= \sum_{y' \in Y^{[R_d^C]}} \mathbf{K}[Y(R_d)Y(R_d^C)|D_{Q_d}D_{Q_d^C}Z]_{d_{Q_d}d_{Q_d^C}z}^{x_{Q_d}y'} \quad (35)$$

$$= \sum_{y' \in Y^{[R_d^C]}} \mathbf{K}[Y_{R_d}Y_{R_d^C}|D_{Q_d}D_{Q_d^C}Z]_{dz}^{x_{Q_d}y'} \quad (36)$$

$$= \sum_{y' \in Y^{[R_d^C]}} \mathbf{K}[Y_{[n]}|D_{[n]}Z]_{dz}^{x_{Q_d}y'} \quad (\text{using exchangeability}) \quad (37)$$

Note that

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

## 5 Appendix:see-do model representation

### Modularise the treatment of probability

**Theorem 5.1** (See-do model representation). *Suppose we have a decision problem that provides us with an observation  $x \in X$ , and in response to this we can select any decision or stochastic mixture of decisions from a set  $D$ ; that is we can choose a “strategy” as any Markov kernel  $\mathbf{S} : X \rightarrow \Delta(D)$ . We have a utility function  $u : Y \rightarrow \mathbb{R}$  that models preferences over the potential consequences of our choice. Furthermore, suppose that we maintain a denumerable set of hypotheses  $H$ , and under each hypothesis  $h \in H$  we model the result of choosing some strategy  $\mathbf{S}$  as a joint probability over observations, decisions and consequences  $\mathbb{P}_{h,\mathbf{S}} \in \Delta(X \times D \times Y)$ .*

*Define  $X, Y$  and  $D$  such that  $X_{xdy} = x$ ,  $Y_{xdy} = y$  and  $D_{xdy} = d$ . Then making the following additional assumptions:*

1. *Holding the hypothesis  $h$  fixed the observations as have the same distribution under any strategy:  $\mathbb{P}_{h,\mathbf{S}}[X] = \mathbb{P}_{h,\mathbf{S}'}[X]$  for all  $h, \mathbf{S}, \mathbf{S}'$  (observations are given “before” our strategy has any effect)*
2. *The chosen strategy is a version of the conditional probability of decisions given observations:  $\mathbf{S} = \mathbb{P}_{h,\mathbf{S}}[D|X]$*
3. *There exists some strategy  $\mathbf{S}$  that is strictly positive*
4. *For any  $h \in H$  and any two strategies  $\mathbf{Q}$  and  $\mathbf{S}$ , we can find versions of each disintegration such that  $\mathbb{P}_{h,\mathbf{Q}}[Y|DX] = \mathbb{P}_{h,\mathbf{S}}[Y|DX]$  (our strategy tells*

us nothing about the consequences that we don't already know from the observations and decisions)

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

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

$$\mathbb{P}[\mathbf{XDY}]^{ijk} = \mathbb{P}[\mathbf{X}]^i \mathbb{P}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{P}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (38)$$

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

$$\mathbf{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^k \mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{XD}]_{kj}^l \quad (39)$$

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

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

We then have

$$\mathbf{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj}^k = \mathbf{T}\mathbf{X}'_{hj}^k \quad (40)$$

$$= \sum_l \mathbf{T}_{hj}^{kl} \quad (41)$$

$$= \mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^k \quad (42)$$

$$= \mathbf{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj'}^k \quad (43)$$

Thus  $\mathbf{X}' \perp\!\!\!\perp_{\mathbf{T}} \mathbf{D}'|\mathbf{H}'$  and so  $\mathbf{T}[\mathbf{X}'|\mathbf{H}']$  exists (section 1.7) and  $(\mathbf{T}, \mathbf{H}', \mathbf{D}', \mathbf{X}', \mathbf{Y}')$  is a see-do model.

Applying Equation 38 to  $\mathbb{P}_{h,\mathbf{S}}$ :

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{XDY}]^{ijk} = \mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbf{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (44)$$

$$= \mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (45)$$

$$= \mathbb{P}_{h,\mathbf{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbf{T}[\mathbf{X}'\mathbf{Y}'|\mathbf{H}'\mathbf{D}']_{hj}^{ik} \quad (46)$$

$$= \mathbf{S}_i^j \mathbf{T}[\mathbf{X}'\mathbf{Y}'|\mathbf{H}'\mathbf{D}']_{hj}^{ik} \quad (47)$$

$$= \mathbf{S}_i^j \mathbf{T}[\mathbf{X}'|\mathbf{H}'\mathbf{D}']_{hj}^i \mathbf{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (48)$$

$$= \mathbf{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbf{S}_i^j \mathbf{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (49)$$

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

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^i = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{X}]^i \quad (50)$$

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \text{ for some version of } \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}] \quad (51)$$

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

$$\mathbf{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{X}]^k \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}]_{kj}^l \quad (52)$$

Then by substitution of  $\mathbf{Q}$  for  $\mathbf{S}$  in Equation 44 and working through the same steps

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{XDY}]^{ijk} = \mathbf{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbf{Q}_i^j \mathbf{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (53)$$

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

## References

- G. Chiribella, Giacomo D’Ariano, and P. Perinotti. Quantum Circuit Architecture. *Physical review letters*, 101:060401, September 2008. doi: 10.1103/PhysRevLett.101.060401.
- 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>.
- 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.
- James M. Robins and Sander Greenland. Causal Inference Without Counterfactuals: Comment. *Journal of the American Statistical Association*, 95(450):431–435, 2000. ISSN 0162-1459. doi: 10.2307/2669381. URL <http://www.jstor.org/stable/2669381>. Publisher: [American Statistical Association, Taylor & Francis, Ltd.].
- Donald B. Rubin. Causal Inference Using Potential Outcomes. *Journal of the American Statistical Association*, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880. URL <https://doi.org/10.1198/016214504000001880>.
- Peter Selinger. A survey of graphical languages for monoidal categories. *arXiv:0908.3347 [math]*, 813:289–355, 2010. doi: 10.1007/978-3-642-12821-9\_4. URL <http://arxiv.org/abs/0908.3347>. arXiv: 0908.3347.

- Jin Tian and Judea Pearl. A general identification condition for causal effects.  
In *Aaai/iaai*, pages 567–573, July 2002.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions.  
Wiley, Oxford, England, 1950.

## Appendix: