

Understanding Causal Primitives Using Modular Probability

David Johnston

September 22, 2021

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 index¹ 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.

¹the 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.

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, 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.

To do this, we use a novel approach to probability modelling that we find is well suited to building causal models. A typical approach to probability modelling is to construct a probability space $(\mathbb{P}, \Omega, \mathcal{F})$ that serves as a top level model, along with a collection of random variables defined by measurable functions on this space, such that the particular quantities of interest can be obtained from conditional and marginal distributions on this space. Instead we consider a modelling context \mathcal{M} that contains a collection of *probability components*, which are Markov kernels with named inputs and outputs. The names correspond to variables in the standard setting. Probability components with the right input and output types can be *connected*, an operation that yields a new probability component. We relate this back to the standard approach by equipping each probability component with a probability space and requiring that all components are the conditional probability distributions on their assigned spaces corresponding to their input and output labels.

Equipped with this foundation, we apply it to a variety of approaches to causal modelling, showing how it can enable understanding of different approaches in a common framework, and how it can represent assertions that were previously made “outside the model”. First, we consider causal decision problems and derive *see-do models*, which reduce to statistical decision problems when augmented with the principle of expected utility. See-do models are a particular kind of probability component that we call a *comb*, which can be thought of as a probability model that needs something to be inserted into the middle. We consider causal graphical models, and show how under a very slight modification to the standard notation they induce see-do models, which allows us to formally connect *interventions* to *actions*. Finally, we consider potential out-

comes models and show how we can formalise the typical assertion (which again, lives “outside the model”) that potential outcomes represent counterfactual values. Potential outcomes models as typically used do not contain counterfactual assertions and in fact feature comb and insert components almost but not quite identical to combs and inserts found in causal graphical models.

I’m probably going to have to cut some of the above

Contents

1	Introduction	1
2	Technical prerequisites	3
2.1	Markov kernels	4
2.2	Cartesian and tensor products	5
2.3	Delta measures, erase maps, copy maps	5
2.4	Products	5
2.5	Label dictionaries, equality and inclusion	5
2.6	Labeled Markov kernels, conditional probabilities	6
2.7	Ambient model	7
2.8	Modelling context axioms	7
2.9	Connection	8
2.10	Conditional independence	12
2.11	Uniqueness of disintegrations	13
2.12	Existence of modelling context	13
2.13	Standard probability models	13
3	See-do models	13
3.1	See-do models and classical statistics	14
3.2	Combs	15
4	Causal Bayesian Networks	15
5	Potential outcomes with and without counterfactuals	16
5.1	Potential outcomes in see-do models	18
5.2	Parallel potential outcomes representation theorem	19
6	Appendix:see-do model representation	22

2 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 Ω , 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).

2.1 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 σ -algebra made up of every subset of X , also called the discrete σ -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:

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

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

2.2 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:

$$\mathbf{K} \otimes \mathbf{L} := \begin{array}{c} \boxed{\mathbf{K}} \\ \boxed{\mathbf{L}} \end{array} \quad (7)$$

2.3 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:

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

2.4 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:

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

2.5 Label dictionaries, equality and inclusion

A *modelling context* \mathcal{M} is a collection of labels, a label dictionary (defined in this section) and a collection of conditional probabilities (defined in section 2.6).

It can be thought of as a namespace that ensures that the properties defined in section 2.8 hold.

, if we respect the naming rules (section 2.8), any two conditional probabilities with the same name will be the same and any model containing two instances of the label name will assign probability 0 to any point at which the two labels differ.

A label collection is a collection of labels, each of which is assigned to a measurable set. For example, $\{X : X, Y : Y, X_1 : X_1, X_2 : X_2\}$. A sequence of labels is associated with the cartesian product of the sets associated with the individual labels. For example, $(X, Y) : X \times Y$.

A label dictionary is a set of statements of equality between labels and between labels and sequences of labels. For example, we might have the set $\{X = (X_1, X_2), W = X\}$. The sets associated with the labels on each side of all equality statements must be the same.

If an equality $W = X$ appears in the label dictionary then any true expression involving conditional probabilities or labels remains true when all instances of X are replaced by W or vice-versa. A label X_1 is contained in X iff it is possible to derive from the label dictionary a statement of equality between X and a sequence of labels containing X_1 .

Because it saves a lot of space, we will generally hold to the convention that a label X is associated with the set X . However, this convention sometimes fails, for example when we have two labels X_1 and 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 sequence of labels that consists entirely of trivial labels is equivalent to the trivial label. An empty sequence is equivalent to the trivial label.

2.6 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* $\mathbf{A}_C := (A_i)_{i \in C}$ and *codomain labels* $\mathbf{B}_D := (B_i)_{i \in D}$. A label assignment is valid if the Markov kernel's domain matches the sets associated with its domain labels and its codomain is the same as the sets associated with its codomain labels.

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

A conditional probability $\mathbb{L}[A_C | B_D; \mathbf{K}]$ is a labeled kernel $(\mathbf{K}, \mathbf{A}_C, \mathbf{B}_D)$ along with an *ambient model* (Definition 2.6) \mathbb{L} .

Graphically, we place the labels on the wires of a conditional probability and include the ambient model in the name of the Markov kernel:

$$\mathbb{L}[B_1 B_2 | A_1 A_2; \mathbf{K}] := \begin{matrix} A_1 \\ A_2 \end{matrix} \boxed{\mathbf{K}[\mathbb{L}]} \begin{matrix} B_1 \\ B_2 \end{matrix} \quad (12)$$

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

$$\mathbb{L}[B_1 B_2 | A_1 A_2; \mathbf{K}] = (A_1, A_2) \boxed{\mathbf{K}[\mathbb{L}]} (B_1, B_2) \quad (13)$$

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

2.7 Ambient model

An ambient model is a tuple $\mathbb{L} := (\mathbf{L}, D, E, \mathcal{C}, \mathcal{R})$ where $\mathbf{L} : D \rightarrow \Delta(E)$ is a Markov kernel, D and E are the domain and codomain, \mathcal{C} is a set of *choice variable definitions* and \mathcal{R} is a set of *random variable definitions*. \mathcal{C} and \mathcal{R} each assign labels to measurable functions on D and E respectively. Given a label X associated with a set X , if the label assignment $X : f \in \mathcal{R}$ then it must be the case that $f : D \rightarrow X$ or $f : E \rightarrow X$.

If the label dictionary contains the assignment $X = (Y, Z)$ and the variable assignments $\{X : f_X, Y : f_Y, Z : f_Z\}$ are a subset of either \mathcal{R} or \mathcal{C} , then it must be the case that $f_X(h) = (f_Y(h), f_Z(h))$.

Let D be the sequence of all choice variables $D = (A)_{A \in \mathcal{D}}$. Then f_D must be an invertible function of D .

The same label can appear in both \mathcal{C} and \mathcal{R} and can appear multiple times in \mathcal{R} . However, each label can appear at most once in \mathcal{C} . In addition, the assignment of repeated labels to functions must be compatible with the first axiom of modelling contexts (2.8); namely, the random variables that are given the same name must be almost surely equal.

Given $X : f_X \in \mathcal{R}$, $\mathbb{L}[X|D] := (d, A) \mapsto \sum_{i \in f_X^{-1}(A)} \mathbf{L}_d^i$ is the *conditional probability of X given D* under \mathbb{L} .

Given $\{X : f_X, Y : f_Y\} \in \mathcal{R}$, any Markov kernel \mathbf{M} such that $\mathbb{L}[XY|D]_d^{A \times B} = \sum_{i \in f_X^{-1}(A)} \mathbb{L}[X|D]_d^i \mathbf{M}(d, i)^B$ is a version of the *conditional probability of Y given (X, D)* under \mathbb{L} .

2.8 Modelling context axioms

We place the following requirements on elements of \mathcal{M} :

1. Any model containing two instances of the same label will assign probability 0 to any point at which the two labels differ
2. The kernel $\mathbf{K} : A^C \rightarrow \Delta(B^D)$ of the conditional probability $\mathbb{L}[A_C | B_D; \mathbf{K}]$ is a version of the *probability of A_C given B_D* under the ambient model \mathbb{L}

The first axiom can be broken down into two cases.

Case 1 Given $\mathbb{M}[XY|XZ; \mathbf{K}]$, we require that there exist some $\mathbf{H} : X \times Z \rightarrow \Delta(Y)$ such that

$$\mathbf{K} = \begin{array}{c} \text{X} \quad \text{X} \\ \quad \diagdown \quad \diagup \\ \quad \text{---} \bullet \text{---} \\ \quad \quad \quad \text{H} \\ \quad \quad \quad \text{---} \text{Y} \\ \text{Z} \text{---} \end{array} \quad (14)$$

$$\iff \quad (15)$$

$$\mathbf{K}_{xz}^{x'y} = \delta[x]^{x'} \mathbf{L}_{xz}^y \quad (16)$$

Case 2 Given $\mathbb{M}[XXY|Z; \mathbf{K}]$, we require that there exist some $\mathbf{H} : Z \rightarrow \Delta(X \times Y)$ such that

$$\mathbf{K} = \begin{array}{c} \text{X} \\ \diagup \quad \diagdown \\ \text{Z} \text{---} \text{H} \text{---} \bullet \text{---} \begin{array}{c} \text{X} \\ \text{X} \\ \text{Y} \end{array} \end{array} \quad (17)$$

$$\iff \quad (18)$$

$$\mathbf{K}_z^{xx'y} = \delta[x]^{x'} \mathbf{L}_z^{xy} \quad (19)$$

2.9 Connection

Connection is an operation \Rightarrow that “joins” two labeled Markov kernels where the labels can be matched and preserves unmatched inputs and outputs. A key property of extension is that, if both input Markov kernels satisfy Axiom 1, then the output also satisfies axiom 1. Depending on the labels of the inputs, the extension operation can reduce to:

- The tensor product
- The matrix product
- The operation that combines a marginal and a conditional probability to yield a joint probability

We overload \Rightarrow to also refer to an conditional probabilities. In general, the result of connecting of two conditional probabilities sharing an ambient model does not satisfy Axiom 2. We assume the result of connecting two conditional probabilities is equipped with an arbitrary ambient model satisfying Axiom 2. We show that it is always possible to make some choice for this and, in particular cases, an ambient model matching the inputs can be chosen.

Given two labeled Markov kernels $(\mathbf{K}, \mathbf{I}_F, \mathbf{O}_F)$ and $(\mathbf{L}, \mathbf{I}_S, \mathbf{O}_S)$, make the

following label identifications:

$$\mathbf{O}_{F\cdot} := \mathbf{O}_F \setminus \mathbf{l}_S \quad (20)$$

$$\mathbf{O}_{FS} := \mathbf{O}_F \cap \mathbf{l}_S \quad (21)$$

$$\mathbf{l}_{F\cdot} := \mathbf{l}_F \setminus \mathbf{l}_S \quad (22)$$

$$\mathbf{l}_{FS} := \mathbf{l}_F \cap \mathbf{l}_S \quad (23)$$

$$\mathbf{l}_{\cdot S} := \mathbf{l}_S \setminus \mathbf{l}_F \quad (24)$$

$$\mathbf{O}_{I_F O_S^*} := \mathbf{O}_S \cap \mathbf{l}_F \quad (25)$$

$$\mathbf{O}_{O_F O_S^*} := \mathbf{O}_F \cap \mathbf{O}_S \setminus \mathbf{l}_S \quad (26)$$

$$(27)$$

The output labels may contain duplicates, but not the input labels. We use the convention that at most one of each label in \mathbf{O}_F belongs to \mathbf{O}_{FS} and the multiplicity of each label in $(\mathbf{O}_{F\cdot}, \mathbf{O}_{FS}, \mathbf{O}_{\cdot S})$ is equal to the multiplicity of each label in $(\mathbf{O}_F, \mathbf{O}_S)$. Thus if a label is shared between \mathbf{O}_F and \mathbf{l}_S and appears 3 times in \mathbf{O}_F , one copy is assigned to \mathbf{O}_{FS} and two to $\mathbf{O}_{F\cdot}$.

$(\mathbf{K}, \mathbf{l}_F, \mathbf{O}_F)$ can be connected to $(\mathbf{L}, \mathbf{l}_S, \mathbf{O}_S)$ iff $\mathbf{O}_{I_F O_S^*}$ is trivial (the output of the “second” kernel cannot be connected to the inputs of the “first”) and $\mathbf{O}_{O_F O_S^*}$ is also trivial (the two kernels do not propose conflicting models of the same label).

Definition 2.1 (extension). Consider two labeled Markov kernels $(\mathbf{K}, \mathbf{l}_F, \mathbf{O}_F)$ which can be connected to $(\mathbf{L}, \mathbf{l}_S, \mathbf{O}_S)$. Because they can be conected, we can write $(\mathbf{K}, (\mathbf{l}_{F\cdot}, \mathbf{l}_{FS}), (\mathbf{O}_{F\cdot}, \mathbf{O}_{FS}))$ and $(\mathbf{L}, (\mathbf{l}_{FS}, \mathbf{l}_{\cdot S}), \mathbf{O}_S)$.

Then Equations 29 and 30 are equivalent definitions of extension:

$$(\mathbf{K}, (\mathbf{l}_{F\cdot}, \mathbf{l}_{FS}), (\mathbf{O}_{F\cdot}, \mathbf{O}_{FS})) \Rightarrow (\mathbf{L}, (\mathbf{l}_{FS}, \mathbf{l}_{\cdot S}), \mathbf{O}_S) := (\mathbf{J}, (\mathbf{l}_{F\cdot}, \mathbf{l}_{FS}, \mathbf{l}_{\cdot S}), (\mathbf{O}_{F\cdot}, \mathbf{O}_{FS}, \mathbf{O}_S)) \quad (28)$$

$$:= \begin{array}{c} \mathbf{l}_{F\cdot} \text{---} \boxed{\mathbf{K}} \text{---} \mathbf{O}_{F\cdot} \\ \mathbf{l}_{FS} \text{---} \bullet \text{---} \mathbf{O}_{FS} \\ \mathbf{l}_{\cdot S} \text{---} \bullet \text{---} \mathbf{O}_S \end{array} \quad (29)$$

$$\mathbf{J}_{yqr}^{zxw} = \mathbf{K}_{yq}^{zx} \mathbf{L}_{xqr}^w \quad (30)$$

Equation 29 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:

$$(\mathbf{J}, (\mathbf{l}_{F\cdot}, \mathbf{l}_{FS}, \mathbf{l}_{\cdot S}), (\mathbf{O}_{F\cdot}, \mathbf{O}_{FS}, \mathbf{O}_S)) = \left[\begin{array}{c} \mathbf{l}_{F\cdot} \text{---} \\ \mathbf{l}_{FS} \text{---} \bullet \text{---} \\ \mathbf{l}_{\cdot S} \text{---} \end{array} \right] \left[\begin{array}{c} \boxed{\mathbf{K}} \\ \text{---} \end{array} \right] \left[\begin{array}{c} \text{---} \bullet \text{---} \\ \text{---} \end{array} \right] \left[\begin{array}{c} \text{---} \mathbf{O}_{F\cdot} \\ \text{---} \mathbf{O}_{FS} \\ \text{---} \mathbf{O}_S \end{array} \right] \quad (31)$$

$$(32)$$

Lemma 2.2 (Extension is associative up to permutation of labels). *Given labeled Markov kernels $(\mathbf{K}, \mathbf{l}_K, \mathbf{O}_K)$, $(\mathbf{L}, \mathbf{l}_L, \mathbf{O}_L)$ and $(\mathbf{J}, \mathbf{l}_J, \mathbf{O}_J)$,*

$$(\mathbf{K} \rightrightarrows \mathbf{L}) \rightrightarrows \mathbf{J} \stackrel{perm}{=} \mathbf{K} \rightrightarrows (\mathbf{L} \rightrightarrows \mathbf{J}) \quad (33)$$

Where $\stackrel{perm}{=}$ indicates equality up to permutation of labels and corresponding swaps applied to the Markov kernel.

Proof. This will be proven with string diagrams, and consequently generalises to the operation defined by Equation 17 in other Markov kernel categories.

Define

$$\mathbf{l}_{K..} := \mathbf{l}_K \setminus \mathbf{l}_L \setminus \mathbf{l}_J \quad (34)$$

$$\mathbf{l}_{KL.} := \mathbf{l}_K \cap \mathbf{l}_L \setminus \mathbf{l}_J \quad (35)$$

$$\mathbf{l}_{K.J} := \mathbf{l}_K \cap \mathbf{l}_J \setminus \mathbf{l}_L \quad (36)$$

$$\mathbf{l}_{KLJ} := \mathbf{l}_K \cap \mathbf{l}_L \cap \mathbf{l}_J \quad (37)$$

$$\mathbf{l}_{.L.} := \mathbf{l}_L \setminus \mathbf{l}_K \setminus \mathbf{l}_J \quad (38)$$

$$\mathbf{l}_{.LJ} := \mathbf{l}_L \cap \mathbf{l}_J \setminus \mathbf{l}_K \quad (39)$$

$$\mathbf{l}_{..J} := \mathbf{l}_J \setminus \mathbf{l}_K \setminus \mathbf{l}_L \quad (40)$$

$$\mathbf{O}_{K..} := \mathbf{O}_K \setminus \mathbf{l}_N \setminus \mathbf{l}_J \quad (41)$$

$$\mathbf{O}_{KL.} := \mathbf{O}_K \cap \mathbf{l}_L \setminus \mathbf{l}_J \quad (42)$$

$$\mathbf{O}_{K.J} := \mathbf{O}_K \cap \mathbf{l}_J \setminus \mathbf{l}_L \quad (43)$$

$$\mathbf{O}_{KLJ} := \mathbf{O}_K \cap \mathbf{l}_L \cap \mathbf{l}_J \quad (44)$$

$$\mathbf{O}_{L.} := \mathbf{O}_L \setminus \mathbf{l}_J \quad (45)$$

$$\mathbf{O}_{LJ} := \mathbf{O}_L \cap \mathbf{l}_J \quad (46)$$

Also define

$$(\mathbf{P}, \mathbf{l}_P, \mathbf{O}_P) := \mathbf{K} \rightrightarrows \mathbf{L} \quad (47)$$

$$(\mathbf{Q}, \mathbf{l}_Q, \mathbf{O}_Q) := \mathbf{L} \rightrightarrows \mathbf{J} \quad (48)$$

Then

$$(\mathbf{K} \Rightarrow \mathbf{L}) \Rightarrow \mathbf{J} = \mathbf{P} \Rightarrow \mathbf{J} \quad (49)$$

$$= \begin{array}{c} \begin{array}{c} l_{P.} \\ l_{P.J} \end{array} \begin{array}{c} \boxed{\mathbf{P}} \\ \bullet \end{array} \begin{array}{c} o_{P.} \\ o_{P.J} \end{array} \\ \begin{array}{c} l_{J.} \end{array} \begin{array}{c} \boxed{\mathbf{J}} \\ \bullet \end{array} o_J \end{array} \quad (50)$$

$$= \begin{array}{c} \begin{array}{c} l_{K..} \\ l_{KL.} \\ l_{L.} \\ l_{K.J} \\ l_{KLJ} \\ l_{LJ} \\ l_{..J} \end{array} \begin{array}{c} \boxed{\mathbf{K}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \\ \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\mathbf{L}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \\ \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\mathbf{J}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \end{array} \quad (51)$$

$$\stackrel{perm}{=} \begin{array}{c} \begin{array}{c} l_{K..} \\ l_{KL.} \\ l_{K.J} \\ l_{KLJ} \\ l_{L.} \\ l_{LJ} \\ l_{..J} \end{array} \begin{array}{c} \boxed{\mathbf{K}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \\ \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\mathbf{L}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \\ \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} \boxed{\mathbf{J}} \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \begin{array}{c} o_{K..} \\ o_{KL.} \\ o_{K.J} \\ o_{KLJ} \\ o_{L.} \\ o_{LJ} \\ o_J \end{array} \end{array} \quad (52)$$

$$= \begin{array}{c} \begin{array}{c} l_{K.} \\ l_{KQ} \end{array} \begin{array}{c} \boxed{\mathbf{K}} \\ \bullet \end{array} \begin{array}{c} o_{K.} \\ o_{KQ} \end{array} \\ \begin{array}{c} l_{Q.} \end{array} \begin{array}{c} \boxed{\mathbf{Q}} \\ \bullet \end{array} o_Q \end{array} \quad (53)$$

□

Theorem 2.3 (Extension is compatible with a modelling context). *Given $\mathbb{M}[\mathbf{ZX}|\mathbf{YQ}; \mathbf{K}]$ and $\mathbb{N}[\mathbf{W}|\mathbf{XQR}; \mathbf{L}]$ in a modelling context \mathcal{M} , let $\mathbf{J} = \mathbb{M}[\mathbf{ZX}|\mathbf{YQ}; \mathbf{K}] \Rightarrow \mathbb{N}[\mathbf{W}|\mathbf{XQR}; \mathbf{L}]$. Then there exists an ambient model \mathbb{O} such that \mathcal{M} along with $\mathbb{O}[\mathbf{ZXW}|\mathbf{YQR}; \mathbf{J}]$ is a valid modelling context, where we can choose $\mathbb{O} = \mathbb{M}$ if $\mathbb{M} = \mathbb{N}$ and the non-shared inputs \mathbf{Y} and \mathbf{R} are trivial.*

Proof. By inspecting the label set in $\mathbb{O}[\mathbf{ZXW}|\mathbf{YQR}; \mathbf{J}]$, we can see that no labels from either of the inputs are duplicated. We need to verify that if either of the inputs had a duplicated label, then the result of the extension still satisfies axiom 1.

We have a number of cases to deal with Introduce $\mathbb{I}[\mathbf{X}|\mathbf{X}; \text{Id}_X]$ and note that in the case of Equation 14 corresponds to the equation

$$\mathbf{K} = \mathbb{I}[\mathbf{X}|\mathbf{X}; \text{Id}_X] \Rightarrow \mathbb{M}[\mathbf{Y}|\mathbf{XZ}; \mathbf{L}] \quad (54)$$

$$\Rightarrow \mathbf{J} = (\mathbb{I}[\mathbf{X}|\mathbf{X}; \text{Id}_X] \Rightarrow \mathbb{M}[\mathbf{Y}|\mathbf{XZ}; \mathbf{L}]) \Rightarrow \mathbb{N}[\mathbf{W}|\mathbf{XQR}; \mathbf{L}] \quad (55)$$

$$\Rightarrow \mathbf{J} = \mathbb{I}[\mathbf{X}|\mathbf{X}; \text{Id}_X] \Rightarrow (\mathbb{M}[\mathbf{Y}|\mathbf{XZ}; \mathbf{L}] \Rightarrow \mathbb{N}[\mathbf{W}|\mathbf{XQR}; \mathbf{L}]) \quad (56)$$

We can always define a new probability space so that $\mathbb{O}[\mathbf{ZXW}|\mathbf{YQR}; \mathbf{J}]$ satisfies axiom 2. We will show that if $\mathbb{M} = \mathbb{N}$ and the non-shared inputs \mathbf{Y} and \mathbf{R} are trivial then in particular we can choose $\mathbb{O} = \mathbb{M}$. \square

Definition 2.4 (marginal). Given a conditional probability $\mathbb{K}[\mathbf{XY}|\mathbf{W}]$, the marginal $\mathbb{K}[\mathbf{X}|\mathbf{W}]$ is defined as

$$\mathbb{K}[\mathbf{X}|\mathbf{W}] := \mathbf{W} \dashv_{\mathbb{K}}^{\mathbf{X}} \quad (57)$$

$$\mathbb{K}[\mathbf{X}|\mathbf{W}]_w^x = \sum_{y \in Y} \mathbb{K}[\mathbf{XY}|\mathbf{W}]_w^{xy} \quad (58)$$

Definition 2.5 (disintegration). $\mathbb{K}[\mathbf{Y}|\mathbf{XW}]$ is a disintegration of $\mathbb{K}[\mathbf{XY}|\mathbf{W}]$ if

$$\mathbb{K}[\mathbf{X}|\mathbf{W}] \Rightarrow \mathbb{K}[\mathbf{Y}|\mathbf{XW}] = \mathbb{K}[\mathbf{XY}|\mathbf{W}] \quad (59)$$

Any Markov kernel \mathbf{L} with the property

$$\mathbf{L}_{xw}^y = \frac{\mathbb{K}[\mathbf{XY}|\mathbf{W}]_w^{xy}}{\sum_{x \in X} \mathbb{K}[\mathbf{XY}|\mathbf{W}]_w^{xy}} \quad \forall w, y : \text{the denominator is positive} \quad (60)$$

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

Definition 2.6 (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 $\mathbf{K}[\mathbf{Y}|\mathbf{X}] = \mathbf{K}$.

2.10 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 .

2.11 Uniqueness of disintegrations

Every conditional probability $\mathbf{K}[X|Y]$ 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.

2.12 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.

2.13 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.

3 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 $\mathbb{S}_\alpha[D|X]$, 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 α, β . Our chosen strategy is the probability of D given X : $\mathbb{P}_\alpha[D|X] \stackrel{krn}{=} \mathbb{S}_\alpha[D|X]$. 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 α, β .

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

$$\mathbb{P}_\alpha[XDY|H] \stackrel{krn}{=} \mathbb{T}[X|H] \Rightarrow \mathbb{S}_\alpha[D|X] \Rightarrow \mathbb{T}[Y|XHD] \quad (61)$$

The proof is given in Appendix 6. Note that $\mathbb{T}[X|H]$ exists by virtue of the fact $X \perp\!\!\!\perp_{\mathbb{T}} D|H$. While this independence is what enables Equation 61, in general $X \not\perp\!\!\!\perp_{\mathbb{P}_\alpha} D|H$, so \mathbb{T} cannot be a disintegration of \mathbb{P}_α . 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 \mathbb{S}_α 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}_α , \mathbb{T} , which are mutually incompatible conditional probabilities.

3.1 See-do models and classical statistics

A *statistical model* (or *statistical experiment*) is a collection of probability distributions indexed by some set Θ . 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 $R(\theta, \mathbb{S}_\alpha) = \sum_{x \in X} \sum_{d \in D} \mathbb{P}_\theta^x(\mathbb{S}_\alpha)_x^d l(h, d)$.

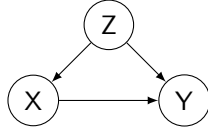
Suppose we have a see-do model $\mathbb{T}[XY|HD]$ with $Y \perp\!\!\!\perp_{\mathbb{T}} X|HD$, and suppose that the random variable Y is a “reverse utility” function taking values in \mathbb{R} for which low values are considered desirable. Then, defining a loss $l : H \times D \rightarrow \mathbb{R}$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

First we will offer some commentary

We can consider this an instance of a see-do model. To do so consistently within a modelling context \mathcal{M} , 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 $\mathbb{T}[XYX'Y'|HD]$ 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 $\mathbb{T}[X'Y'Z'|HD]_{h, x}^{x'yz} \stackrel{krn}{=} \mathbb{T}[Z|H]_h^z \delta[x]^{x'} \mathbb{T}[Y|XZH]_{h, x'z}^y$.

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

$$\mathbb{T}[XYZ|H] = \mathbb{U}[Z|H] \Rightarrow \mathbb{T}[X|H] \Rightarrow \mathbb{U}[Y|XHZ] \quad (68)$$

$$\mathbb{T}[X'Y'Z'|HD] \stackrel{krn}{=} \mathbb{U}[Z|H] \Rightarrow \mathbb{D}[X|D] \Rightarrow \mathbb{U}[Y|XHZ] \quad (69)$$

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).

5 Potential outcomes with and without counterfactuals

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 we 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.

5.1 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 5.1 (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 5.1 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 5.2 (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.

5.2 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 $\mathbb{P}[X_1X_2] = \mathbb{P}[X_2X_1]$. There are also many examples of exchangeable sequences that are not mutually independent and identically distributed – for example, if we want to use random variables Y_1 and Y_2 to model our subjective uncertainty regarding two flips of a coin of unknown fairness, we regard our initial uncertainty for each flip to be equal $\mathbb{P}[Y_1] = \mathbb{P}[Y_2]$ and we our state of knowledge of the second flip after observing only the first will be the same as our state of knowledge of the first flip after observing only the second $\mathbb{P}[Y_2|Y_1] = \mathbb{P}[Y_1|Y_2]$, then our model of subjective uncertainty is exchangeable.

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

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

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

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

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

Definition 5.3 (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 5.4 (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 5.5 (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 5.6 (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 $\mathbf{K}[Y_{\{i,j\}}|D_{\{i,j\}}Z]_{kk}^{lm} = \llbracket l = m \rrbracket \mathbf{K}[Y_i|D_iZ]_k^l$ for all i, j, k, l, m .

Theorem 5.7 (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 $y := (y_i)_{i \in |D|} \in Y^{|D|}$, $d := (d_i)_{i \in [n]} \in D^{[n]}$, $x := (x_i)_{i \in [n]} \in Y^{[n]}$.

$$\mathbf{K}[X_n|D_nZ]_{dz}^x = \sum_{y \in Y^{|D|}} \mathbf{K}[X_{[n]}|D_nZY(D)]_{dz y}^x \mathbf{K}[Y(D)|D_{[n]}Z]_{dz}^y \quad (70)$$

$$= \sum_{y \in Y^{|D|}} \prod_{i \in [n]} \delta[y_{d_i}]^{x_i} \mathbf{K}[Y(D)|D_nZ]_{dz}^y \quad (71)$$

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 $\mathbf{K}[Y_{[n]}|D_{[n]}Z]_{dz}^x = 0$ 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 $Y(S_d) = (Y_d)_{d \in S_d}$.

$$\mathbf{K}[X_{[n]} | D_{[n]} Z]_{dz}^x = \sum_{y \in Y^{|D|}} \prod_{i \in Q_d} \delta[y_{d_i}]^{x_i} \mathbf{K}[Y(D) | D_{[n]} Z]_{dz}^y \quad (72)$$

$$= \sum_{y' \in Y^{|R_d^C|}} \mathbf{K}[Y(R_d) Y(R_d^C) | D_{Q_d} D_{Q_d^C} Z]_{d_{Q_d} d_{Q_d^C} z}^{x_{Q_d} y'} \quad (73)$$

$$= \sum_{y' \in Y^{|R_d^C|}} \mathbf{K}[Y_{R_d} Y_{R_d^C} | D_{Q_d} D_{Q_d^C} Z]_{dz}^{x_{Q_d} y'} \quad (74)$$

$$= \sum_{y' \in Y^{|R_d^C|}} \mathbf{K}[Y_{[n]} | D_{[n]} Z]_{dz}^{x_{Q_d} y'} \quad (\text{using exchangeability}) \quad (75)$$

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

6 Appendix:see-do model representation

Modularise the treatment of probability

Theorem 6.1 (See-do model representation). *Suppose we have a decision problem that provides us with an observation $x \in X$, and in response to this we can select any decision or stochastic mixture of decisions from a set D ; that is we can choose a “strategy” as any Markov kernel $\mathbf{S} : X \rightarrow \Delta(D)$. 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 $\mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{DX}] = \mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{DX}]$ (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 $(\mathbf{T}, H', D', X', Y')$ such that $\mathbb{P}_{h,\mathbf{S}}[\mathbf{XDY}]^{ijk} = \mathbf{T}[X'|H']_h^i \mathbf{S}_i^j \mathbf{T}[Y'|X'H'D']_{ijk}^k$.

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

$$\mathbb{P}[\mathbf{XDY}]^{ijk} = \mathbb{P}[X]^i \mathbb{P}[D|X]_i^j \mathbb{P}[Y|XD]_{ij}^k \quad (76)$$

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

$$\mathbf{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbf{S}}[X]^k \mathbb{P}_{h,\mathbf{S}}[Y|XD]_{kj}^l \quad (77)$$

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

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

We then have

$$\mathbf{T}[X'|H'D']_{hj}^k = \mathbf{T}X'_{hj}^k \quad (78)$$

$$= \sum_l \mathbf{T}_{hj}^{kl} \quad (79)$$

$$= \mathbb{P}_{h,\mathbf{S}}[X]^k \quad (80)$$

$$= \mathbf{T}[X'|H'D']_{hj'}^k \quad (81)$$

Thus $X' \perp\!\!\!\perp_{\mathbf{T}} D'|H'$ and so $\mathbf{T}[X'|H']$ exists (section 2.10) and $(\mathbf{T}, H', D', X', Y')$ is a see-do model.

Applying Equation 76 to $\mathbb{P}_{h,\mathbf{S}}$:

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{XDY}]^{ijk} = \mathbb{P}_{h,\mathbf{S}}[X]^i \mathbb{P}_{h,\mathbf{S}}[D|X]_i^j \mathbb{P}_{h,\mathbf{S}}[Y|XD]_{ij}^k \quad (82)$$

$$= \mathbb{P}_{h,\mathbf{S}}[X]^i \mathbb{P}_{h,\mathbf{S}}[Y|XD]_{ij}^k \quad (83)$$

$$= \mathbb{P}_{h,\mathbf{S}}[D|X]_i^j \mathbf{T}[X'Y'|H'D']_{hj}^{ik} \quad (84)$$

$$= \mathbf{S}_i^j \mathbf{T}[X'Y'|H'D']_{hj}^{ik} \quad (85)$$

$$= \mathbf{S}_i^j \mathbf{T}[X'|H'D']_{hj}^i \mathbf{T}[Y'|X'H'D']_{ihj}^k \quad (86)$$

$$= \mathbf{T}[X'|H']_h^i \mathbf{S}_i^j \mathbf{T}[Y'|X'H'D']_{ihj}^k \quad (87)$$

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

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{X}]^i = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{X}]^i \quad (88)$$

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \text{ for some version of } \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}] \quad (89)$$

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

$$\mathbf{T}_{hj}^{kl} = \mathbb{P}_{h,\mathbf{Q}}[\mathbf{X}]^k \mathbb{P}_{h,\mathbf{Q}}[\mathbf{Y}|\mathbf{XD}]_{kj}^l \quad (90)$$

Then by substitution of \mathbf{Q} for \mathbf{S} in Equation 82 and working through the same steps

$$\mathbb{P}_{h,\mathbf{S}}[\mathbf{XDY}]^{ijk} = \mathbf{T}[\mathbf{X}'|\mathbf{H}']_h^i \mathbf{Q}_i^j \mathbf{T}[\mathbf{Y}'|\mathbf{X}'\mathbf{H}'\mathbf{D}']_{ihj}^k \quad (91)$$

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>.
- 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>.
- Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In Mikoaj Bojaczuk and Alex Simpson, editors, *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing, 2019. ISBN 978-3-030-17127-8.
- 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>.

- 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>.
- 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.
- 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.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England, 1950.

Appendix: