

# Causal Statistical Decision Theory|What are interventions?

David Johnston

September 25, 2020

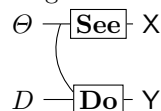
## Contents

<b>1</b>	<b>Theories of causal inference</b>	<b>2</b>
<b>2</b>	<b>Causal Questions</b>	<b>5</b>
2.1	Formalising Causal Statistical Decision Problems: See-Do Maps .	6
2.2	D-causation . . . . .	8
<b>3</b>	<b>Notation</b>	<b>9</b>
<b>4</b>	<b>What is the difference between Causal Bayesian Networks and See-Do models?</b>	<b>9</b>
4.1	Influence Diagrams vs See-Do models . . . . .	10
4.2	Influence diagrams vs Causal Bayesian Networks . . . . .	12
4.3	What is meant by “variables”? . . . . .	13
4.4	Necessary relationships . . . . .	14
4.5	Recursive Structural Causal Models . . . . .	15
4.6	Recursive Structural Causal Models with Necessary Relationships	16
4.7	Cyclic Structural Causal Models . . . . .	20
4.8	Not all variables have well-defined interventions . . . . .	22
4.8.1	Necessary relationships in cyclic SCMs . . . . .	23
<b>5</b>	<b>Definitions and key notation</b>	<b>26</b>
5.1	Standard Symbols . . . . .	27
5.2	Probability Theory . . . . .	27
5.3	Product Notation . . . . .	28
5.4	String Diagrams . . . . .	28
5.4.1	Comparison of notations . . . . .	32
5.4.2	Working With String Diagrams . . . . .	33
5.5	Random Variables . . . . .	34

## 6 Decision makers and decision models 44

### 6.1 Expected utility maximisers . . . . . 45

Note on terminology: I am trying the name “See-Do model” to describe the following:



I was calling it a “causal theory” before. Reasons for the change: I think “See-Do” helps to understand what the model does, and the name doesn’t make premature claims to explain causality. Also, it’s only two syllables which I like.

## 1 Theories of causal inference

Feedback start here

Beginning in the 1930s, a number of associations between cigarette smoking and lung cancer were established: on a population level, lung cancer rates rose rapidly alongside the prevalence of cigarette smoking. Lung cancer patients were far more likely to have a smoking history than demographically similar individuals without cancer and smokers were around 40 times as likely as demographically similar non-smokers to go on to develop lung cancer. In laboratory experiments, cells which were introduced to tobacco smoke developed *ciliastasis*, and mice exposed to cigarette smoke tars developed tumors(Proctor, 2012). Nevertheless, until the late 1950s, substantial controversy persisted over the question of whether the available data was sufficient to establish that smoking cigarettes *caused* lung cancer. Cigarette manufacturers famously argued against any possible connection (Oreskes and Conway, 2011) and Roland Fisher in particular argued that the available data was not enough to establish that smoking actually caused lung cancer (Fisher, 1958). Today, it is widely accepted that cigarettes do cause lung cancer, along with other serious conditions such as vascular disease and chronic respiratory disease (World Health Organisation, 2018; Wiblin, 2016).

The question of a causal link between smoking and cancer is a very important one. Individuals who enjoy smoking (or think they might) may wish to avoid smoking if cigarettes pose a severe health risk, so they are interested in knowing whether or not it is so. Potential investors in cigarette manufacturers want to know if the product they are backing is likely to see limited adoption due to health concerns. People holding investments in cigarette manufacturing firms want the world to be such that cigarettes do not pose a substantial health risk, as this increases the value of their investment. Governments and organisations with a responsibility for public health may see themselves as having responsibility to discourage smoking as much as possible if smoking is severely detrimental to health. The costs and benefits of poor decisions about smoking are large: 8 million annual deaths are attributed to cigarette-caused cancer and vascular disease in 2018(World Health Organisation, 2018) while global cigarette

sales were estimated at US\$711 million in 2020, while (Statista, 2020) (a figure which might be substantially larger if cigarettes were not widely believed to be harmful).

The question of whether or not cigarette smoking causes cancer illustrates two key facts about causal questions: First, having the right answers to some causal questions is of tremendous importance to huge numbers of people. Second, even when large amounts of data show unambiguous associations between phenomena of interest, it is still difficult to know when a causal conclusion is justified.

Causal conclusions are often justified on the basis of ad-hoc reasoning. For example Krittanawong et al. (2020) states:

[...] the potential benefit of increased chocolate consumption, reducing coronary artery disease (CAD) risk is not known. We aimed to explore the association between chocolate consumption and CAD.

It is not clear whether Krittanawong et. al. mean that a negative association between chocolate consumption and CAD implies that increased chocolate consumption is likely to reduce coronary artery disease, or that an association may be relevant to the question and the reader should draw their own conclusions. Whether the implication is being suggested by Krittanawong et. al. or merely imputed by naïve readers, it is being drawn on an ad-hoc basis – no argument for the implication can be found in this paper. As Pearl (2009) has forcefully argued, additional assumptions are always required to answer causal questions from associational facts, and stating these assumptions explicitly allows those assumptions to be productively scrutinised.

Theories of causal inference exist to enable formal rather than ad-hoc reasoning about causal questions. Instead of posing informal causal question and answering them based on ad-hoc reasoning, within a theory of causal inference we ask about properties of “causal models” (which are simply mathematical types defined by the theory) subject to certain assumptions we are willing to make. A successful theory of causal inference should enable causal models that “adequately represent” the original informal question, and the assumptions we invoke should be more accessible to scrutiny than ad-hoc assertions made in the course of answering the informal question.

As well as defining causal models, which represent *claims about causation*, theories of causal inference also formalise the problem of *inferring the correct causal model* - this is the problem of taking some observational data and concluding which causal models are “possible” or “appropriate to use for the given purpose”.

Defining causal models is difficult. While philosophical theories of causation are not heavily discussed in the applied literature on causal inference, the principles used to motivate the definitions of causal models are widely discussed in the philosophical literature. Applied theories of causal inference:

1. “ $X_i$  causes  $X_j$ ” means that there exist different *ideal interventions* that result in different values of  $X_i$ , hold other “causally sufficient” variables

constant, do not directly affect  $X_j$  but nonetheless entail different values of  $X_j$

2. “ $X_i$  causes  $X_j$ ” means that the *counterfactual value* of  $X_j$  would be different “if  $X_i$  had taken a different value”

In practice, most theories of causal inference seem to be based on the notion of *ideal interventions*. Even “counterfactual” theories of causal inference (such as the theory based on “potential outcomes” notation) tend to define counterfactual values as “values that a variable would have taken were it exposed to an ideal intervention”, if they are defined at all (Morgan and Winship, 2014; Rubin, 2005; Richardson and Robins, 2013). There are, however, alternative definitions of counterfactual values such as Lewis’ closest world semantics (Lewis, 1986), though these definitions have serious difficulties.

“Ideal interventions” are difficult to define. The structural model approach of Pearl (2009) defines ideal interventions in terms of “causally sufficient models”. However, this definition ends up circular:

- An  $[X_i, X_j]$ -ideal intervention is an operation whose result is determined by applying the do-calculus to a causally sufficient triple  $((\Omega, \mathcal{F}, \mathbb{Q}), \mathcal{G}, \mathbf{U})$
- A triple  $((\Omega, \mathcal{F}, \mathbb{Q}), \mathcal{G}, \mathbf{U})$  is  $[X_i, X_j]$ -causally sufficient if  $\mathbf{U}$  contains  $X_i$ ,  $X_j$  and “all intervenable variables” that *cause* (definition (1)) both  $X_i$  and  $X_j$

Circularity is a recognised problem with interventional definitions of causation (Woodward, 2016). An alternative approach is to designate certain real-world events – such as flipping coins, querying random number generators and so forth – as prototypical “ideal interventions”. The main problem with this definition is that it fails to answer any causal query that hasn’t been subject to a randomised trial, which seems to leave most causal questions unanswerable, despite the fact that many of them have apparent answers (Pearl, 2018a), and the causal questions represented under this definition are often considered to be inadequate representations of the original questions to which answers were sought (Deaton and Cartwright, 2018; Heckman, 1991). Additional challenges posed by “ideal interventions” are discussed in Section 4.

An account of ideal interventions has not yet been offered that stands up to scrutiny. This difficulty is not without precedent, and is certainly no reason to abandon theories built upon ideal interventions. It is also difficult to provide an account of what it means for data to be “distributed according to probability distribution  $\mathbb{P}$ ” (Hájek, 2019) that stands up under scrutiny, but the usefulness of probability distributions in representing data regularities is widely accepted.

However, in light of the difficulties with theories of causation, it is notable that many causal questions can be formalised without appealing to a theory of causation at all, an observation also made by Dawid (2020). Causal statistical decision theory (CSDT), introduced here, shows how the type definition

---

<sup>1</sup>Weaker conditions for causal sufficiency are possible, but they are still premised on causal relationships, so the charge of circularity stands (Shpitser and Pearl, 2008).

of “causal models” is almost completely determined by the given *causal decision problem*. Causal decision problems are problems such as: Should I smoke? Should cigarettes be taxed? Such problems are ubiquitous, and causal inferences are instrumental in answering them. As with existing theories of causal inference, determining the relation between observational probabilities and causal models is the key problem in CSDT, and existing theories of causal inference have developed many useful approaches to this. CSDT additionally offers the opportunity to understand what remains true of causal decision problems regardless of the theory of causation in use and the opportunity to consider theories of causation which may be useful in solving particular causal decision problems but need not be universally valid.

## 2 Causal Questions

Pearl and Mackenzie (2018) has proposed three types of causal question:

1. Association: How are  $X$  and  $Y$  related? How would observing  $X$  change my beliefs about  $Y$ ?
2. Intervention: What would happen if I do ... ? How can I make  $E$  happen?
3. Counterfactual: Was  $X$  responsible for  $Y$ ? What if I had done ... instead of what I actually did?

I focus on a particular type of question here: “How can I make  $E$  happen?”. Pearl calls this type of question an “interventional” question but we observe that this query is in fact a *causal decision problem* (I also wish to avoid the term “intervention” as it has a technical meaning which that I touched on above). A causal decision problem is a problem of the following form:

Given my available options  $D$  and the way I would like the world to turn out, which options are likely to have a desirable result?

When one asks “How can I make  $E$  happen?”, I suppose that the answer will be one or more elements of a set of available options  $D$ . Observing that  $E$  corresponds to “the way I would like the world to turn out”, we can see that this question is a causal decision problem.

I choose to focus on this type of problem because the problem of choosing a decision given desired results is ubiquitous, and Definition 2 is a broadly applicable prototype for this kind of problem. It’s possible that a similar approach can succeed for counterfactual queries, but obviously a different prototype problem would be required. Enough questions are raised by focusing on causal decision problems to restrict our attention to this type of problem for now.

Causal *statistical* causal decision problems extend causal decision problems by introducing data and assumptions about how the world works:

Given my available options  $D$ , data  $\mathbf{X}$  and preferences  $u$ , which options are likely to have a desirable result?

## 2.1 Formalising Causal Statistical Decision Problems: See-Do Maps

I introduce *see-do maps* which formalise causal statistical decision problems introduced in Definition 2. I don't derive see-do maps from first principles - rather, I make use of probability theory and expected utility theory because they're the standard tools used to formalise models of noisy observations and choice under uncertainty.

Section 6 (currently unfinished) considers more fundamental assumptions that might lead to see-do maps

First, we formally define *causal statistical decision problems*.

**Definition 2.1** (Causal Statistical Decision Problem). A *causal statistical decision problem* (CSDP) is a tuple  $(D, x, X, Y, u)$  where  $D$  is a set of *decisions* that may be taken,  $x \in X$  is a *data sequence* of observations,  $Y$  is the space of possible *consequences* of decisions and  $u : Y \rightarrow \mathbb{R}$  is a *utility function* where  $u(y) \geq u(y')$  implies that  $y$  is at least as desirable as  $y'$ .

We place a number of additional requirements on causal statistical decision problems:

- $D$  is a denumerable set
- The principle of *expected utility maximisation* is acceptable, so the desirability of a probability distribution  $\mathbb{P}_Y \in \Delta(\mathcal{Y})$  can be found by evaluating  $\mathbb{E}_{\mathbb{P}_Y}[u]$
- $X$  and  $Y$  are *standard measurable* ( $D$  is automatically standard measurable due to it being denumerable)
- A decision maker may select any probability distribution in  $\Delta(\mathcal{D})$  as its stochastic decision. A deterministic distribution  $\delta_d$  is equivalent to directly choosing  $d \in D$ , while deciding on a non-deterministic distribution involves consulting a source of randomness to select a particular decision in accordance with the probabilities chosen

Definition 2.1 are the elements of a CSDP that we regard as fixed.

A *decision maker* brings prior knowledge to the problem in the form of a *see-do map*, which captures the relation between the *observed probability distribution* and the *consequence map*.

**Definition 2.2** (Observed probability distribution). Given a CSDP  $(D, x, X, Y, u)$ , under hypothesis  $\theta$  the data  $x$  are “distributed according to” some *observed probability*  $\mathbb{O}_\theta \in \Delta(\mathcal{X})$ .

**Definition 2.3** (Consequence map). Given a decision set  $D$  and a consequence set  $Y$ , under hypothesis  $\theta$  a *consequence map*  $\mathbb{C}_\theta : D \rightarrow \Delta(\mathcal{Y})$  is a Markov kernel mapping decisions to probability distributions on the consequence space. For a stochastic decision  $\gamma \in \Delta(\mathcal{D})$ ,  $\gamma\mathbb{C}$  is the *consequence* of  $\gamma$  with respect to  $\mathbb{C}_\theta$ .

“Distributed according to” is a weasel phrase with no particularly satisfactory interpretation, but I'd like to avoid worrying about it at this point (Hájek, 2019).

**Definition 2.4** (Causal hypothesis). A causal hypothesis  $\theta$  on a CSDP  $(D, x, X, Y, u)$  is a pair  $(\mathbb{O}_\theta, \mathbb{C}_\theta) \in \Delta(\mathcal{X}) \times \Delta(\mathcal{Y})^D$ . A hypothesis  $(\mathbb{O}_\theta, \mathbb{C}_\theta)$  can be interpreted as the statement “the observed data are distributed according to  $\mathbb{O}_\theta$  and the consequences of decisions are given by  $\mathbb{C}_\theta$ ”

**Definition 2.5** (Causal hypothesis class). A *causal hypothesis class*  $\Theta$  is a set of causal hypotheses on  $(D, x, X, Y, u)$ .

**Definition 2.6** (See-do map). Given a CSDP  $(D, x, X, Y, u)$  and a causal hypothesis class  $\Theta$ , the *see-do* model  $\mathbb{T} : \Theta \times D \times (\mathcal{X} \otimes \mathcal{Y})$  is the map  $\mathbb{T} : (\theta, d, A \times B) \mapsto \mathbb{O}_\theta(A) \mathbb{C}_{(\theta, d)}(B)$ .

We assume that see-do maps are Markov kernels (i.e. there is a  $\sigma$ -algebra  $\theta \otimes \mathcal{D}$  such that the map  $(\theta, d) \mapsto \mathbb{O}_\theta(A)\mathbb{C}_{(\theta, d)}(B)$  is  $\theta \otimes \mathcal{D}$ -measurable). This is guaranteed if  $\Theta$  is denumerable and endowed with the discrete  $\sigma$ -algebra. If  $\Theta$  is standard measurable, this imposes restrictions on  $\Theta$ ; for example, we can't have some  $d \in D$ ,  $A \times B \in \mathcal{X} \otimes \mathcal{Y}$  such that  $\mathbb{O}_\cdot(A)\mathbb{C}_{(\cdot, d)}(B)^{-1}(\{1\})$  is a non-measurable collection of hypotheses in  $\Theta$ .

There is a generic graphical representation of see-do maps. See Section 5.4 for a thorough explanation of the notation used.

Define  $\mathbb{O} : \Theta \rightarrow \Delta(\mathcal{X})$  by  $\mathbb{O} : \theta \mapsto \mathbb{O}_\theta$  and  $\mathbb{C} : \Theta \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$  by  $\mathbb{C} : (\theta, d) \mapsto \mathbb{C}_{(\theta, d)}$ . Then we have

$$\mathbb{T} = \begin{array}{c} \theta - \boxed{\text{O}} - X \\ \quad \quad \quad \curvearrowright \\ D - \boxed{\text{C}} - Y \end{array} \quad (1)$$

A key difference between CSDT and other approaches to causal inference is that diagrams in CSDT feature two coupled maps  $\mathbb{O}$  and  $\mathbb{C}$ , while most other approaches to causal inference represent both  $\mathbb{O}$  and  $\mathbb{C}$  in one diagram. Lattimore and Rohde (2019) is the only other example I am aware of that draws both  $\mathbb{O}$  and  $\mathbb{C}$ .

A causal hypothesis class  $\Theta$  induces a binary relation between observed probability distributions  $\mathbb{O}_\theta$  and consequence maps  $\mathbb{C}_\theta$ . This approach is very agnostic about the actual relation induced – we do not even insist that the range of the observed data  $X$  is the same as the range of possible consequences  $Y$  (though we will generally limit our attention to cases where the two coincide).

In common with Heckerman and Shachter (1995), decisions (or “acts”) are primitive elements of see-do maps. In contrast to our work, Heckerman and Shachter (1995) only discuss deterministic *consequence maps*, while see-do maps represent relations between consequence maps and observed probability.

Decisions are similar to the “regime indicators” found in Dawid (2020). They coincide precisely if we suppose that the observation and consequence spaces coincide ( $X = Y$ ) and there exists an “idle” decision  $d^* \in D$  such that  $\mathbb{C}_{(\cdot, d^*)} = \mathbb{O}_{\cdot}$ . However, in general we don’t require that  $\mathbb{O}$  and  $\mathbb{C}$  are related in this manner. This assumption will be revisited in

A section I haven't written yet

## 2.2 D-causation

Feedback end here, WIP

It is possible to define a notion of  $D$ -dependent causation for consequence maps  $\mathbb{C}_\theta$ . A similar idea is discussed extensively in Heckerman and Shachter (1995). The main differences are that what we call “consequence maps” map decisions to probability distributions over possible consequences while Heckerman and Shachter work with “states” that map decisions deterministically to consequences, thus where we discuss “conditional independence” Heckerman and Shachter discuss “limited unresponsiveness”. In addition, while we define  $D$ -causation relative to a particular consequence map  $\mathbb{C}_\theta$ , Heckerman and Shachter define it with respect to a *set* of states. Dawid (2020) also introduces the key concept of “extended conditional independence” which is related to  $D$ -causation.

As discussed in Section 4, in my view it is very difficult to define a notion of “causation” without reference to a set of basic operations such as  $D$ .

See Section 5.5 for the definition of random variables.

Add definition of conditional independence, revise wire label definitions

A variable  $Y_1$   $D$ -causes another variable  $Y_2$  if knowing the effect decisions have on  $Y_1$  is enough to determine the effect of decisions on  $Y_2$ .

**Definition 2.7** ( $D$ -causation). Given a consequence map  $\mathbb{C}_\theta : D \rightarrow \Delta(\mathcal{Y})$ , random variables  $Y_1 : Y \times D \rightarrow Y_1$ ,  $Y_2 : Y \times D \rightarrow Y_2$  and domain variable  $D : Y \times D \rightarrow D$  (Definition 5.8),  $Y_1$   $D$ -causes  $Y_2$  iff  $Y_2 \perp\!\!\!\perp_{\mathbb{C}_\theta} D | Y_1$ .

For example, under most obvious consequence maps  $\mathbb{C}_\theta$ , my light switch  $D$ -causes my light to turn on given the set of decisions:

$D = \{\text{walk to the switch and press it with my thumb,}$   
     
     stay in bed}

Given some reasonable hypothesis  $\theta$ , for any  $d \in D$  the light will be on if and only if the switch ends up in the on position.

However, the light switch does *not*  $D$ -cause the light to turn on under the set of decisions :

$D = \{\text{flip the light switch without touching the fuse box,}$   
     flip the fuse box, then flip the light switch}

In this case, the state of the light depends on the decision taken beyond the position of the light switch.



### 3 Notation

- $X$  is a random variable,  $X$  is its codomain and  $\mathcal{X}$  is the  $\sigma$ -algebra on  $X$
- Bold letters  $\mathbf{X}$  may be used for product spaces,  $\mathbf{x}$  for elements of product spaces,  $\mathbf{f}$  for vector valued functions and  $\mathbf{X}$  for random variables taking values in product spaces. The absence of bold font does not imply the absence of product space structure
- Nodes in a graph are italic sans serif  $X$
- Given an indexed product space  $\prod_{i \in \mathcal{I}} X_i$ ,  $\pi_i : \prod_{i \in \mathcal{I}} X_i \rightarrow X_i$  is the projection map  $(x_1, \dots, x_i, \dots, x_n) \mapsto x_i$
- Given a product space  $X \times Y$ ,  $\pi_X : X \times Y \rightarrow X$  is the projection map  $(x, y) \mapsto x$
- Given an index set  $\mathcal{I}$ ,  $\mathbf{X}_{\mathcal{I}}$  is the indexed product space  $\prod_{i \in \mathcal{I}} X_i$
- $[n]$  is the index set  $\{1, \dots, n\}$  for  $n \in \mathbb{N}$
- Given an indexed product space  $\mathbf{X}_{[n]}$ ,  $\mathbf{X}_{< j}$  is the set  $\prod_{i=1}^j X_i$
- $\otimes$  is the coupled tensor product, see Definition 5.9
- Given a random variable  $X$ ,  $F_X$  is the associated Markov kernel, see Definition 5.1
- See section 5.4.2 for rules of string diagram manipulation
- $*_X : X \rightarrow \Delta(\{*\})$  is the discard map defined in Equation 82

### 4 What is the difference between Causal Bayesian Networks and See-Do models?

See-Do models and Causal Bayesian Networks (and related models such as SCMs) are quite different in their appearance and in the interpretation of various elements. In defining See-Do models, we assume that there is a decision problem that fixes in advance the observation space  $E$ , the space of consequences  $F$  and the set of available decisions  $D$ . By including  $D$ , See-Do models have an agent “baked in” to the definition. In contrast, Causal Bayesian Networks assume that a set of observed variables  $\mathbf{X}$  and a set of unobserved variables  $\mathbf{U}$  is fixed by nature. A subtlety here is that  $\mathbf{U}$  is generally *not* known, but it may be assumed that whatever variables actually comprise  $\mathbf{U}$ , they may be generically representable by some known set of variables  $\mathbf{U}'$  without loss. Because a set of decisions  $D$  is absent, Causal Bayesian Networks appear to model agent-independent “causal relationships”.

Despite the fact that Causal Bayesian Networks don’t seem to be built to model the consequences of an agent’s decision making, they are nonetheless considered

to be appropriate for this purpose. This is because the “do-operations” that Causal Bayesian Networks support are considered to have some relationship to any set  $D$  decisions some agent might want to consider.

The question of how this relationship might be determined in general is one that I have not seen addressed anywhere. Typically, the approach taken is “I know it when I see it”. For example, if I were a doctor and I could either A) hand a patient a prescription or B) not hand the patient the prescription, and if I had a set of observational data of past patients including a variable  $S$  representing whether or not they had received a prescription for the drug, I could consider option A to correspond to  $do(S = 1)$  and option B to correspond to  $do(S = 0)$  in some causal model (perhaps in the “true” causal model). While this might appear to be reasonable, we should be cautious: this is a completely *ad-hoc* assumption.

In fact, given a sufficiently rich set of variables  $\mathbf{U} \cup \mathbf{X}$ , I argue that  $do(S)$  will almost never even approximate the consequences of an action known in advance to fix the value of some variable  $S$ . For this reason, it is valuable to have a theory like CSDT that is concerned only with the consequences of actions.

#### 4.1 Influence Diagrams vs See-Do models

Influence diagrams are used to represent causal models in a manner similar to, but not quite the same as, Causal Bayesian Networks. Using the version found in Dawid (2002), an influence diagram is a directed acyclic graph (DAG) with two node types: “chance” nodes and “decision” nodes. For example:



Is an influence diagram with the decision node  $D$  and the chance node  $F$ . When nodes are associated with sets representing possible values, influence diagrams represent sets of Markov kernels. For example, if we associate the measurable set  $(D, \mathcal{D})$  with  $D$  and  $(F, \mathcal{F})$  with  $F$ , then we could take Diagram 2 to represent the set of all Markov kernels  $D \rightarrow \Delta(\mathcal{F})$ , or if we add some additional assumptions it might represent a particular subset  $S \subset \Delta(\mathcal{F})^D$  of these kernels that share certain properties. Compare this to a string diagram:

$$D - \boxed{\mathbb{K}} - F \quad (3)$$

This diagram represents a *particular* Markov kernel  $\mathbb{K} : D \rightarrow \Delta(\mathcal{F})$ . A set such as that represented by Diagram 2 could be constructed by creating a set  $\Theta$  and a function  $f : \Theta \rightarrow S$  that indexes each element of  $S$  with  $\theta \in \Theta$ . Then  $S = \{f(\theta) | \theta \in \Theta\}$ . If  $f$  is measurable then there is a Markov kernel  $\mathbb{T} : \Theta \times D \rightarrow \Delta(\mathcal{F})$  such that  $f(\theta) = \mathbb{T}_{\theta, \cdot}$ . In this sense, for any additional assumptions that are combined with Diagram 2 to yield the set  $S$  of Markov

notation

kernels, there exists a single Markov kernel  $\mathbb{T} : \Theta \times D \rightarrow \Delta(\mathcal{F})$  that generates  $S$ .  $\mathbb{T}$  can be drawn as

$$\begin{array}{c} \Theta \\ D \end{array} \curvearrowright \boxed{\mathbb{T}} - F \quad (4)$$

Diagram 4 has a few more elements that Diagram 2 -  $\Theta$  and  $\mathbb{T}$  in particular. If I were to define exactly what  $D$ ,  $\Theta$ ,  $\mathbb{T}$  and  $F$  were, Diagram 4 would represent a unique Markov kernel. On the other hand, Diagram 2 along with a precise definition of  $D$  and  $F$  could represent many different sets of Markov kernels, depending on whatever additional properties I want elements of  $S$  to share. Loosely, we can say that for any kernel  $\mathbb{T}$ , there is a diagram 2 along with additional assumptions that represents “the same thing”.

Influence diagrams *as typically used in causal modelling* usually additional assumptions. For example, Dawid (2002) proposes that  $D$  contains a special “do nothing” element  $o \in D$  such that the observations in state  $\theta$  are given by  $\mathbb{T}_{\theta,o}$  and consequences in state  $\theta$  are given by  $\mathbb{T}_{\theta,\cdot}$ . This corresponds to the See-Do model

$$\begin{array}{c} \Theta \\ D \end{array} \begin{array}{c} \boxed{\mathbb{T}_{\cdot,o}} \\ \boxed{\mathbb{T}} \end{array} \begin{array}{c} - F \\ - F \end{array} \quad (5)$$

Under the condition that the choice of index “doesn’t matter” somehow

notation

While this is a feature of influence diagrams in Dawid (2002), *in general* a diagram like 2 can represent an arbitrary set of typed Markov kernels. See-Do models generate  $\Theta$ -indexed sets of Markov kernels. We can therefore represent generic See-Do models with influence diagrams.

There is some kind of measurability condition needed to make the two representations coincide

Concretely, given the influence diagram

$$\begin{array}{c} \textcircled{E} \\ \textcircled{D} \rightarrow \textcircled{F} \end{array} \quad (6)$$

and any See-Do model  $\mathbb{T} : \Theta \times D \rightarrow \Delta(\mathcal{E} \otimes \mathcal{F})$  there exists a set of auxiliary conditions  $A$  such that the model  $(\mathbb{J}, A)$  is equivalent to  $\mathbb{T}$ . To illustrate how influence diagrams and string diagrams compare,  $\mathbb{T}$  can be drawn:

$$\mathbb{T} := \begin{array}{c} \Theta \\ D \end{array} \begin{array}{c} \boxed{\mathbb{H}} \\ \boxed{\mathbb{C}} \end{array} \begin{array}{c} - F \\ - E \end{array} \quad (7)$$

subject to irrelevance of index  $\Theta$

Maybe move the following to an appendix?

**Definition 4.1** (Markov kernel/influence diagram compatibility). Given a Markov kernel  $\mathbb{K} : \mathcal{E} \rightarrow \Delta(\mathcal{F})$ , an influence diagram  $\mathcal{J} = (S, A, E)$  and an injective  $f : X \cup A \rightarrow \mathcal{F} \otimes \mathcal{E}$  which assigns each node to exactly one random variable in  $\mathcal{F} \otimes \mathcal{E}$ , if for all  $X_1, X_2 \in \mathbf{X}$  we have  $X_1 \perp_{\mathcal{J}} X_2 \implies f(X_1) \perp_K f(X_2)$

still need to *define* conditional independence WRT kernels, though I think I can do so

Consider an arbitrary See-Do model  $\mathbb{T} : D \times \Theta \rightarrow \Delta(\mathcal{E}_1 \otimes \mathcal{E}_2)$  and random variables  $D := \pi_D, E := \pi_E, F := \pi_F$  on  $\mathcal{D} \otimes \mathcal{E} \otimes \mathcal{F}$ . For any  $\theta \in \Theta$ ,  $\mathbb{T}_\theta$  is compatible with the influence diagram  $\mathcal{J} = (\{E, F\}, \{D\}, \{D \rightarrow F\})$  with respect to the injective function

$$f : \begin{cases} A \mapsto D \\ E \mapsto E \\ F \mapsto F \end{cases} \quad (8)$$

There is always some  $\mathbb{H} : \Theta \rightarrow \Delta(\mathcal{E})$  and  $\mathbb{C} : \Theta \times D \rightarrow \Delta(\mathcal{E} \otimes \mathcal{F})$  such that:  $\mathbb{T}_\theta$  is equal to

$$\mathbb{T}_\theta = D - \begin{array}{c} \boxed{\mathbb{H}_\theta} \\ \boxed{\mathbb{C}_\theta} \end{array} \begin{array}{c} E \\ F \end{array} \quad (9)$$

Which implies  $E \perp_{\mathbb{T}_\theta} F$  and  $E \perp_{\mathbb{T}_\theta} D$ .  
The influence diagram

$$\mathcal{J} = \begin{array}{c} \textcircled{E} \\ \textcircled{D} \rightarrow \textcircled{F} \end{array} \quad (10)$$

Features the d-separations  $E \perp_{\mathcal{J}} F$  and  $E \perp_{\mathcal{J}} D$  (Peters et al., 2017; Woodward, 2016; Dawid, 2002). Thus  $\mathbb{T}_\theta$  is compatible with  $\mathcal{I}$  for all  $\theta \in \Theta$ .

## 4.2 Influence diagrams vs Causal Bayesian Networks

See-Do models can always be represented by influence diagrams with auxiliary assumptions. We can then learn something about how Causal Bayesian Networks compare to See-Do models by asking how they compare to influence diagrams. The key difference between Causal Bayesian Networks and influence diagrams is that the diagrams do not contain decision nodes. Instead of Diagram 2, a Causal Bayesian Network for the same system might be

$$\textcircled{F} \quad (11)$$

The difference between Diagram 2 and Diagram 11 is that the former demands a set  $D$  to be bound to the decision node  $D$  and a set  $F$  to be bound to  $F$ , while Diagram 11 demands only  $F$ . Instead of explicitly representing decisions that can be chosen, Causal Bayesian Networks *by default* feature a set of *do-interventions* on the chance nodes which seem to have a role similar to decisions in influence diagrams and See-Do models (in fact, Pearl (2009) pg 108 suggests that do-interventions and decisions are the same thing). This default set of do-interventions is what allows CBNs to avoid explicitly requiring a set  $D$ . If a set of decisions is required that is not equivalent to the set of do-interventions, this can be specified via auxiliary assumptions, although in practice influence diagrams are usually adopted such as in (Yang et al., 2018).

some more examples

Dealing with a set of decisions  $D$  can be troublesome. It can easily be the case that, for example, I might be tasked with inferring a consequence map that someone else might use and I am not privy to the decisions that they might be able to make. In this case, I'd need to pick a set of decisions  $D$  which I am pretty sure covers all the possibilities.

Alternatively, the set  $D$  might be unworkably large. The set of *all* the decisions you could in principle make at this moment in as much detail as you can - this is clearly something that's far too big to write down and work with in solving an inference problem.

Speculatively, the fact that Causal Bayesian Networks default to do-interventions might help with the problems of unknown or unworkably large decision sets. If any possible decision must be resolvable to do-interventions of some type, and the effects of do-interventions are well defined, then could provide a basis for partially solving decision problems while remaining ignorant of the particular set of decisions that will ultimately be selected from.

as in, doing the inference but not picking the best option

Beyond this, causal effects as they are informally understood *seem* to refer to universal things, not things that depend on the set of decisions one has available. While not an especially strong reason to avoid specifying  $D$  in causal models, it is a reason nonetheless.

Do-interventions, however, cannot solve these problems. If there are no restrictions on the variables that may be included in a CBN model, then as we show do-operations and equation surgery frequently produce invalid results. To ensure that any do-interventions at all are well defined, Causal Bayesian Networks require the specification of “intervenable variables” or “basic interventions”, a requirement that is analogous to the requirement of a set of decisions  $D$  in See-Do models and influence diagrams. Specifying  $D$  may be difficult, but the do-intervention paradigm provides no solutions to this difficulty; it merely sweeps it under the rug.

### 4.3 What is meant by “variables”?

Not sure where to put it, but Pearl pp 162-163 puts his models on the hook for including arbitrary variables

## 4.4 Necessary relationships

The relationship between a person’s body mass index, their weight and their height defines what body mass index is. A fundamental claim of ours is that any causal model that defines “the causal effect of body mass index” should do so without reference to any submodel that violates this definitional relationship violation of the definition. This is an important assumption, and it rests on a judgement of what causal models ought to do. I think it is quite clear that when anyone asks for a causal effect, they expect that any operations required to define the causal effect *do not change the definitions of the variables they are employing*. While theories of causality have a role in sharpening our understanding of the term *causal effect*, the thing called a “causal effect” in an SCM should still respect some of our pre-theoretic intuitions about what causal effects are or else it should be called something else. “Causal effects” that depend on redefining variables do not respect pre-theoretic intuitions about what causal effects are:

- If I ask for the “causal effect of a person’s BMI”, I do not imagine that I am asking what would happen if someone’s BMI were defined to be something other than their weight divided by their height
- If I ask for the “causal effect of a person’s weight”, I do not imagine that I am asking what would happen if someone’s weight were not equal to their volume multiplied by their density
- If I ask for the “causal effect of a person’s weight”, I also do not imagine that I am asking what would happen if their weight were not equal to the weight of fat in their body plus the weight of all non-fat parts of their body
- If I ask for the “causal effect of taking a medicine”, I do not imagine that I am asking what would happen if a person were declared to have taken a medicine independently of whatever substances have actually entered their body and how they entered

We will call relationships that have to hold *necessary relationships*. We provide the example of relationships that have to hold by definition as examples, but definitions may not be the only variety of necessary relationships. For example, one might also wish to stipulate that certain laws of physics are required to hold in all submodels.

If a causal model contains variables that are necessarily related, then an intervention on one of them must always change another variable in the relationship. If I change a person’s weight, their height or BMI must change (or both). If I change their height, their weight or BMI must change and if I change their BMI then their weight or height must change. This conflicts with the usual acyclic definition of causal models, where the proposition that *A* causes *B* rules out the possibility that *B* or any of its descendants are a cause of *A*. Thus in an acyclic model it isn’t possible for an intervention on BMI to change weight or height and interventions on weight and height to also change BMI. Theorem

4.11 formalises this conflict for recursive structural causal models: for any set of variables that are necessarily related by a cyclic relationship, at least one of them has no hard interventions defined.

## 4.5 Recursive Structural Causal Models

We begin by showing that necessary relationships are incompatible with structural causal models.

**Definition 4.2** (Recursive Structural Causal Model). A recursive structural causal model (SCM) is a tuple

$$\mathcal{M} := \langle N, M, \mathbf{X}_{[N]}, \mathbf{E}_{[M]}, \{f_i | i \in [N]\}, \mathbb{P}_{\mathcal{E}} \rangle \quad (12)$$

where

- $N \in \mathbb{N}$  is the number of *endogenous variables* in the model
- $M \in \mathbb{N}$  is the number of *exogenous variables* in the model
- $\mathbf{X}_{[N]} := \{X_i | i \in [N]\}$  where, for each  $i \in [N]$ ,  $(X_i, \mathcal{X}_i)$  is a standard measurable space taking and the codomain of the  $i$ -th endogenous variable
- $\mathbf{E}_{[M]} := \{E_j | j \in [M]\}$  where, for  $j \in [M]$ ,  $E_j$  is a standard measurable space and the codomain of the  $j$ -th exogenous variable
- $f_i : \mathbf{X}_{<i} \times \mathbf{E}_{\mathcal{J}} \rightarrow X_i$  is a measurable function which we call *the causal mechanism controlling the  $i$ -th endogenous variable*
- $\mathbb{P}_{\mathcal{E}} \in \Delta(\mathbf{E}_{\mathcal{J}})$  is a probability measure on the space of exogenous variables

**Definition 4.3** (Observable kernel). Given an SCM  $\mathcal{M}$  with causal mechanisms  $\{f_i | i \in [N]\}$ , define the *observable kernel*  $G_i : E \rightarrow \Delta(\mathbf{X}_{[i]})$  recursively:

$$G_1 = \mathbf{E}_{[M]} \begin{array}{c} \boxed{F_{f_1}} \\ \text{---} \end{array} X_1 \quad f_1 \quad (13)$$

$$G_{n+1} = \mathbf{E}_{[M]} \begin{array}{c} \boxed{G_n} \\ \text{---} \end{array} \begin{array}{c} \text{---} \mathbf{X}_{<n+1} \\ \boxed{F_{f_{n+1}}} \text{---} X_{n+1} \end{array} \quad (14)$$

**Definition 4.4** (Joint distribution on endogenous variables). The *joint distribution on endogenous variables* defined by  $\mathcal{M}$  is  $\mathbb{P}_{\mathcal{M}} := \mathbb{P}_{\mathcal{E}} G_N$  (which is the regular kernel product, see Definition 5.3). For each  $i \in [N]$  define the random variable  $\mathbf{X}_i : \mathbf{X}_{[N]} \rightarrow X_i$  as the projection map  $\pi_i : (x_1, \dots, x_i, \dots, x_N) \mapsto x_i$ . By Lemma 4.5,  $\bigotimes_{i \in [N]} \mathbf{X}_i = \text{Id}_{\mathbf{X}_{[N]}}$ , and so  $\mathbb{P}_{\mathcal{M}}$  is the joint distribution of the variables  $\{\mathbf{X}_i | i \in [N]\}$ .

I use the notation  $\mathbb{P}_{\mathcal{M}}$  rather than  $\mathbb{P}_{\mathbf{X}_{[N]}}$  to emphasize the dependence on the model  $\mathcal{M}$ .

**Lemma 4.5** (Coupled product of all random variables is the identity).  $\bigotimes_{i \in [N]} \mathbf{X}_i = \text{Id}_{\mathbf{X}_{[N]}}$

*Proof.* for any  $\mathbf{X} \in \mathbf{X}_{[N]}$ ,

$$\bigotimes_{i \in [N]} \mathbf{X}_i(\mathbf{X}) = (\pi_1(\mathbf{X}), \dots, \pi_N(\mathbf{X})) \quad (15)$$

$$= (x_1, \dots, x_n) \quad (16)$$

$$= \mathbf{X} \quad (17)$$

□

**Definition 4.6** (Hard Interventions). Let  $\mathcal{M}$  be the set of all *SCMs* sharing the indices, spaces and measure  $\langle N, M, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{[M]}, \mathbb{P}_{\mathcal{E}} \rangle$ . Note that the causal mechanisms are not fixed.

Given an SCM  $\mathcal{M} = \langle N, M, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{[M]}, \{f_i | i \in [N]\}, \mathbb{P}_{\mathcal{E}} \rangle$  and  $\mathcal{S} \subset [N]$ , a *hard intervention* on  $\mathbf{X}_{\mathcal{S}}$  is a map  $Do_{\mathcal{S}} : \mathbf{X}_{\mathcal{S}} \times \mathcal{M} \rightarrow \mathcal{M}$  such that for  $\mathbf{a} \in \mathbf{X}_{\mathcal{S}}$ ,  $Do_{\mathcal{S}}(\mathbf{a}, \mathcal{M}) = \langle N, M, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{[M]}, \{f'_i | i \in [N]\}, \mathbb{P}_{\mathcal{E}} \rangle$  where

$$f'_i = f_i \quad i \notin \mathcal{S} \quad (18)$$

$$f'_i = \pi_i(\mathbf{a}) \quad i \in \mathcal{S} \quad (19)$$

To match standard notation, we will write  $\mathcal{M}^{do(\mathbf{X}_{\mathcal{S}}=\mathbf{a})} := Do_{\mathcal{S}}(\mathbf{a}, \mathcal{M})$

## 4.6 Recursive Structural Causal Models with Necessary Relationships

Necessary relationships are extra constraints on the joint distribution on endogenous variables defined by an SCM. For example, given an SCM  $\mathcal{M}$  if the variable  $\mathbf{X}_1$  represents weight,  $\mathbf{X}_2$  represents height and  $\mathbf{X}_3$  represents BMI then we want to impose the constraint that

$$\mathbf{X}_3 = \frac{\mathbf{X}_1}{\mathbf{X}_2} \quad (20)$$

$\mathbb{P}_{\mathcal{M}}$ -almost surely.

**Definition 4.7** (Constrained Recursive Structural Causal Model (CSCM)). A CSCM  $\mathcal{M} := \langle N, M, \mathbf{X}_{[N]}, \mathbf{E}_{[M]}, \{f_i | i \in [N]\}, \{r_i | i \in [N]\}, \mathbb{P}_{\mathcal{E}} \rangle$  is an SCM along with a set of *constraints*  $r_i : \mathbf{X}_{[N]} \rightarrow X_i$ .

If  $\mathbf{X}_i = r_i(\mathbf{X}_{[N]})$   $\mathbb{P}_{\mathcal{M}}$ -almost surely then  $\mathcal{M}$  is *valid*, otherwise it is *invalid*.

We can recover regular SCMs by imposing only trivial constraints:

**Lemma 4.8** (CSCM with trivial constraints is always valid). *Let  $\mathcal{M}$  be a CSCM with the trivial constraints  $r_i = \pi_i$  for all  $i \in [N]$ . Then  $\mathcal{M}$  is valid.*

*Proof.* By definition 4.7, we require  $\mathbf{X}_i = \pi_i$ ,  $\mathbb{P}_{\mathcal{M}}$ -almost surely.  $\mathbf{X}_i(\mathbf{X}) = \pi_i(\mathbf{X})$  for all  $\mathbf{X} \in \mathbf{X}_{[N]}$  and  $P_{\mathcal{M}}(\mathbf{X}_{[N]}) = 1$ , therefore  $\mathcal{M}$  is valid. □



Call a constraint  $r_i$  *cyclic* if  $\mathbf{X}_i = r_i(\mathbf{X}_{[N]})$  implies there exists an index set  $O \subset [N]$ ,  $O \ni i$ , such that for each  $j \in O$ ,  $\mathbf{b} \in \mathbf{X}_{O \setminus \{j\}}$  there exists  $a \in X_j$  such that

$$\mathbf{X}_{O \setminus \{j\}} = \mathbf{b} \quad (21)$$

$$\implies \mathbf{X}_j = a \quad (22)$$

BMI is an example of a cyclic constraint if we insist that weight and height are always greater than 0. If  $\mathbf{X}_3 = \frac{\mathbf{X}_1}{\mathbf{X}_2}$  then we have:

$$[\mathbf{X}_1, \mathbf{X}_2] = [b_1, b_2] \quad (23)$$

$$\implies \mathbf{X}_3 = \frac{b_1}{b_2} \quad (24)$$

$$[\mathbf{X}_2, \mathbf{X}_3] = [b_2, b_3] \quad (25)$$

$$\implies \mathbf{X}_1 = b_2 b_3 \quad (26)$$

$$[\mathbf{X}_1, \mathbf{X}_3] = [b_1, b_3] \quad (27)$$

$$\implies \mathbf{X}_2 = \frac{b_1}{b_3} \quad (28)$$

The following is a generally useful lemma that should probably be in basic definitions of Markov kernel spaces

**Lemma 4.9** (Projection and selectors). *Given an indexed product space  $\mathbf{X} := \prod_{i \in \mathcal{I}} X_i$  with ordered finite index set  $\mathcal{I} \ni i$ , let  $\pi_i : \mathbf{X} \rightarrow X_i$  be the projection of the  $i$ -indexed element of  $\mathbf{X} \in \mathbf{X}$ .*

*Let  $F_{\pi_i} : \mathbf{X} \rightarrow \Delta(\mathcal{X}_i)$  be the Markov kernel associated with the function  $\pi_i$ ,  $F_{\pi_i} : \mathbf{X} \mapsto \delta_{\pi_i(\mathbf{X})}$ . Given  $O \subset \mathcal{I}$ , define the selector  $S_i^O$ :*

$$S_i^O = \begin{cases} \text{Id}_{X_i} & i \in O \\ *_{X_i} & i \notin O \end{cases} \quad (29)$$

*Then  $\underline{\otimes}_{i \in O} F_{\pi_i} = \otimes_{i \in \mathcal{I}} S_i^O$ .*

*Proof.* Suppose  $O$  is the empty set. Then the empty tensor product  $\otimes_{i \in \emptyset} S_i$  and the empty coupled tensor product  $\underline{\otimes}_{i \in \emptyset} F_{\pi_i}$  are both equal to  $*_{\mathbf{X}}$ .

By definition of  $F_{\pi_i}$ ,  $F_{\pi_i} = \otimes_{i \in \mathcal{I}} S_i^{\{i\}}$ .

Suppose for  $P \subsetneq O$  with greatest element  $k$  we have  $\underline{\otimes}_{i \in P} F_{\pi_i} = \otimes_{i \in \mathcal{I}} S_i^P$ , and suppose that  $j$  is the next element of  $O$  not in  $P$ .

$$(\otimes_{i \in P} F_{\pi_i}) \otimes F_{\pi_j} = \begin{array}{c} \mathbf{X} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in P} F_{\pi_i}} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \boxed{F_{\pi_j}} \quad} X_j \end{array} \end{array} \quad (30)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in P} F_{\pi_i}} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \boxed{F_{\pi_j}} \quad} X_j \end{array} \end{array} \quad (31)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in P} F_{\pi_i}} \quad} \mathbf{X}_P \\ \xrightarrow{\quad * \quad} X_j \\ \xrightarrow{\quad * \quad} X_j \end{array} \end{array} \quad (32)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in P} F_{\pi_i}} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \quad \quad} X_j \end{array} \end{array} \quad (33)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in \mathcal{I}} S_i^P} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \quad \quad} X_j \end{array} \end{array} \quad (34)$$

Because all elements of  $P$  are less than  $j$ , the selector  $S_k^P$  resolves to the discard map for  $k > j$ :

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i < j} S_i^P} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \quad \quad} X_j \end{array} \end{array} \quad (35)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i < j} S_i^P} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \quad \quad} X_j \\ \xrightarrow{\quad \quad \quad} * \end{array} \end{array} \quad (36)$$

$$= \begin{array}{c} \mathbf{X}_{<j} \quad \mathbf{X}_j \quad \mathbf{X}_{>j} \begin{array}{l} \xrightarrow{\quad \boxed{\otimes_{i \in \mathcal{I}} S_i^{P \cup \{j\}}} \quad} \mathbf{X}_P \\ \xrightarrow{\quad \quad \quad} X_j \end{array} \end{array} \quad (37)$$

Where 37 follows from the definition of the selector  $S_i^{P \cup \{j\}}$ .

The proof follows by induction on the elements of  $O$ .

□

**Lemma 4.10** (Hard interventions do not affect the joint distributions of earlier variables). *Given a CSCM  $\mathcal{M} = \langle N, M, \mathbf{X}_{[N]}, \mathbf{E}_{[M]}, \{f_i | i \in [N]\}, \{r_i | i \in$*

$[N]\}, \mathbb{P}_{\mathcal{E}}\rangle$ , any  $k \in [N]$  and any  $O \subset [k-1]$ ,  $P_{\mathcal{M}}(\mathbf{X}_O) = P_{\mathcal{M}}^{do(\mathbf{X}_k)=a}(\mathbf{X}_O)$  for all  $a \in X_k$ .

*Proof.* Let  $G_i^{\mathcal{Q}}$ ,  $i \in [N]$  be the  $i$ -th iteration of the kernel defined in Equations 13 and 14 with respect to model  $\mathcal{Q}$ . Note that from Equation 14

$$\mathbf{E}_{[M]} - \boxed{G_i^{\mathcal{Q}}} - \mathbf{X}_{<i} \xrightarrow{*} = \mathbf{E}_{[M]} - \boxed{G_{i-1}^{\mathcal{Q}}} - \mathbf{X}_{<i} \xrightarrow{F_{f_i}} \xrightarrow{*} \quad (38)$$

$$= G_{i-1}^{\mathcal{Q}} \quad (39)$$

It follows that

$$\mathbf{E}_{[M]} - \boxed{G_N} - \mathbf{X}_{<i} \xrightarrow{*} = G_{i-1} \quad (40)$$

Because  $f_i = f_i^{do(\mathbf{X}_k=a)}$  for  $i < k$ , we have

$$G_i^{\mathcal{M}} = G_i^{\mathcal{M}^{do(\mathbf{X}_k=a)}} \quad (41)$$

for all  $i < k$ . By lemma 4.9, for any  $O \subset [k-1]$  we have  $F_{\mathbf{X}_O} = \otimes_{i \in [N]} S_i^O$ . As there are no elements of  $O$  greater than or equal to  $k$ , the selector  $S_i^O$  resolves to the discard for all  $i \geq k$ . Thus  $F_{\mathbf{X}_O} = (\otimes_{i \in [k-1]} S_i^O) \otimes \ast_{\mathbf{X}_{[N] \setminus [k-1]}}$ . Defining  $S_{[k-1]}^O := \otimes_{i \in [k-1]} S_i^O$ , we have:

$$F_{\mathbf{X}_O} = \mathbf{X}_{[N] \setminus [k-1]} \xrightarrow{\mathbf{X}_{[k-1]} \xrightarrow{\boxed{S_{[k-1]}^O}} \mathbf{X}_O} \xrightarrow{*} \quad (42)$$

Thus

$$\mathbb{P}_{\mathcal{M}}(\mathbf{X}_O) = \mathbb{P}_{\mathcal{E}} G_N^{\mathcal{M}} F_{\mathbf{X}_O} \quad (43)$$

$$\stackrel{42}{=} \triangleleft \mathbb{P}_{\mathcal{E}} \boxed{G_N^{\mathcal{M}}} \boxed{S_{[k-1]}^O} \mathbf{X}_O \xrightarrow{*} (\mathbf{X}_{[N] \setminus [k-1]}) \quad (44)$$

$$\stackrel{40}{=} \triangleleft \mathbb{P}_{\mathcal{E}} \boxed{G_{k-1}^{\mathcal{M}}} \boxed{S_{[k-1]}^O} \mathbf{X}_O \quad (45)$$

$$\stackrel{41}{=} \triangleleft \mathbb{P}_{\mathcal{E}} \boxed{G_{k-1}^{\mathcal{M}^{do(\mathbf{X}_k=a)}}} \boxed{S_{[k-1]}^O} \mathbf{X}_O \quad (46)$$

$$\stackrel{40}{=} \triangleleft \mathbb{P}_{\mathcal{E}} \boxed{G_N^{\mathcal{M}^{do(\mathbf{X}_k=a)}}} \boxed{S_{[k-1]}^O} \mathbf{X}_O \xrightarrow{*} (\mathbf{X}_{[N] \setminus [k-1]}) \quad (47)$$

$$= P_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_O) \quad (48)$$

□

**Theorem 4.11** (Undefined hard interventions with cyclic constraints). *Consider a CSCM  $\mathcal{M} = \langle N, M, \mathbf{X}_{[N]}, \mathbf{E}_{[M]}, \{f_i | i \in [N]\}, \{r_i | i \in [N]\}, \mathbb{P}_{\mathcal{E}} \rangle$  with  $r_i$  a cyclic constraint with respect to  $O \subset [N]$  and the rest of the constraints trivial:  $r_j = \pi_j$ ,  $j \neq i$ , and suppose  $\mathcal{M}$  is valid.*

*If for each  $k \in O$ ,  $\exists A \in \mathcal{X}_i$  such that  $0 < \mathbb{P}_{\mathcal{M}}(\mathbf{X}_i \in A) < 1$  then for at least one  $k \in O$  all models given by a hard intervention on  $\mathbf{X}_k$  are invalid.*

*Proof.* Choose  $k$  to be the maximum element of  $O$ . By the assumption  $\mathcal{M}$  is valid, we have  $\mathbf{X}_i = r_i(\mathbf{X})$ ,  $\mathbb{P}_{\mathcal{M}}$ -almost surely. Let  $B^A = \{\mathbf{b} \in \mathbf{X}_{O \setminus k} | \mathbf{X}_{O \setminus k} = \mathbf{b} \implies \mathbf{X}_k \in A\}$  and  $B^{A^C} = \{\mathbf{b} \in \mathbf{X}_{O \setminus k} | \mathbf{X}_{O \setminus k} = \mathbf{b} \implies \mathbf{X}_k \notin A\}$ .

$r_i$  holds on a set of measure 1, and wherever it holds  $\mathbf{X}_{O \setminus \{k\}}$  is either in  $B^A$  or  $B^{A^C}$ . Thus  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A \cup B^{A^C}) = 1$ .

$B^A$  and  $B^{A^C}$  are disjoint.

By construction,  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A \ \& \ \mathbf{X}_k \in A) = \mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A)$  and  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C} \ \& \ \mathbf{X}_k \in A^C) = \mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C})$ .

By additivity,  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A \ \& \ \mathbf{X}_k \in A) + \mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \notin B^A \ \& \ \mathbf{X}_k \in A) = P_{\mathcal{M}}(\mathbf{X}_k \in A)$ .

By additivity again

$$\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \notin B^A \ \& \ \mathbf{X}_k \in A) = \mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C} \ \& \ \mathbf{X}_k \in A) \quad (49)$$

$$+ \mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in (B^{A^C} \cup B^A)^C \ \& \ \mathbf{X}_k \in A) \quad (50)$$

$$\leq 0 + P_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in (B^{A^C} \cup B^A)^C) \quad (51)$$

$$= 0 \quad (52)$$

Thus  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A \ \& \ \mathbf{X}_k \in A) = P_{\mathcal{M}}(\mathbf{X}_k \in A) = P_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^A)$  and by an analogous argument  $\mathbb{P}_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C}) = P_{\mathcal{M}}(\mathbf{X}_k \in A^C)$ .

Choose some  $a \in A$ , and consider the hard intervention  $\mathcal{M}^{do(\mathbf{X}_k=a)}$ . Suppose  $\mathcal{M}^{do(\mathbf{X}_k=a)}$  is also valid. Then, as before,  $\mathbb{P}_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C}) = \mathbb{P}_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_k \in A^C)$ .

By definition of hard interventions,  $f_k^{do(\mathbf{X}_k=a)} = a$ . Thus  $G_N^{\mathcal{M}^{do(\mathbf{X}_k=a)}} F_{\mathbf{X}_k}$  is the kernel  $\mathbf{X} \mapsto \delta_a$  and it follows that  $\mathbb{P}_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_k) = \delta_a$ .

By lemma 4.10,  $\mathbb{P}_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C}) = P_{\mathcal{M}}(\mathbf{X}_{O \setminus \{k\}} \in B^{A^C}) = P_{\mathcal{M}}(\mathbf{X}_k \in A^C) > 0$ . But  $P_{\mathcal{M}^{do(\mathbf{X}_k=a)}}(\mathbf{X}_k \in A^C) = \delta_z(\mathbf{X}_k \in A^C) = 0$ , contradicting the assumption of validity of  $\mathcal{M}^{do(\mathbf{X}_k=a)}$ .

An analogous argument shows that all hard interventions  $a' \in A^C$  are also invalid.  $\square$

## 4.7 Cyclic Structural Causal Models

It is not very surprising that acyclic causal models cannot accommodate cyclic constraints. Can cyclic causal models do so? While Bongers et al. (2016) has

develope a theory of cyclic causal models, cyclic are generally far less well understood than acyclic models. I show that the theory of cyclic models that Bongers has developed also fails to define hard interventions on variables subject to cyclic constraints. This does not rule out the possibility that there is some other way to define cyclic causal models that do handle these constraints, but I have not taken it upon myself to develop such a theory.

Haven't done any work from here on

We adopt the framework of cyclic structural causal models to make our arguments, adapted from Bongers et al. (2016). This is somewhat non-standard, but allows us to make a stronger argument for the impossibility of modelling arbitrary sets of variables using structural interventional models.

**Definition 4.12** (Structural Causal Model). A structural causal model (SCM) is a tuple

$$\mathcal{M} := \langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}_{\mathcal{I}}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle \quad (53)$$

where

- $\mathcal{I}$  is a finite index set of *endogenous variables*
- $\mathcal{J}$  is a finite index set of *exogenous variables*
- $\mathbf{X}_{\mathcal{I}} := \{X_i\}_{\mathcal{I}}$  where, for each  $i \in \mathcal{I}$ ,  $(X_i, \mathcal{X}_i)$  is a standard measurable space taking and the codomain of the  $i$ -th endogenous variable
- $\mathbf{E}_{\mathcal{J}} := \{E_j\}_{\mathcal{J}}$  where, for  $j \in \mathcal{J}$ ,  $E_j$  is a standard measurable space and the codomain of the  $j$ -th endogenous variable
- $\mathbf{f}_{\mathcal{I}} = \otimes_{i \in \mathcal{I}} f_i$  is a measurable function, and  $f_i : \mathbf{X}_{\mathcal{I}} \times \mathbf{E}_{\mathcal{J}} \rightarrow X_i$  is the causal mechanism controlling  $X_i$
- $\mathbb{P}_{\mathcal{E}} \in \Delta(\mathbf{E}_{\mathcal{J}})$  is a probability measure on the space of exogenous variables
- $\mathbf{E}_{\mathcal{J}} = \otimes_{j \in \mathcal{J}} E_j$  is the set of exogenous variables, with  $\mathbb{P}_{\mathcal{E}} = \mathbf{E}_{\mathcal{J}\#} P_{\mathcal{E}}$  and  $E_j$  is the  $j$ -th exogenous variable with marginal distribution given by  $E_{j\#} \mathbb{P}_{\mathcal{E}}$

If for  $\mathbb{P}_{\mathcal{E}}$ -almost every  $\mathbf{e} \in \mathbf{E}_{\mathcal{J}}$  there exists  $\mathbf{X} \in \mathbf{X}_{\mathcal{I}}$  such that

$$\mathbf{X} = \mathbf{f}_{\mathcal{I}}(\mathbf{X}, \mathbf{e}) \quad (54)$$

Then an SCM  $\mathcal{M}$  induces a unique probability space  $(\mathbf{X}_{\mathcal{I}} \times \mathbf{E}_{\mathcal{J}}, \mathcal{X}_{\mathcal{I}} \otimes \mathcal{E}_{\mathcal{J}}, \mathbb{P}_{\mathcal{M}})$  (Bongers et al., 2016). If no such solution exists then we will say an SCM is invalid, as it imposes mutually incompatible constraints on the endogenous variables. It may be also the case that multiple solutions exist.

If an SCM induces a unique probability space then there exist random variables  $\{X_i\}_{i \in \mathcal{I}}$  such that,  $P_{\mathcal{M}}$  almost surely Bongers et al. (2016):

$$X_i = f_i(\mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}) \quad (55)$$

Where  $\mathbf{X}_{\mathcal{I}} = \bigotimes_{i \in \mathcal{I}} \mathbf{X}_i$ .

A structural causal model can be transformed by *mechanism surgery*. Given  $\mathcal{S} \subset \mathcal{I}$  and a set of new functions  $\mathbf{f}'_{\mathcal{S}} : \mathbf{X}_{\mathcal{S}} \times \mathbf{E}_{\mathcal{J}} \rightarrow \mathbf{X}_{\mathcal{S}}$ , mechanism surgery “replaces” the corresponding parts of  $\mathbf{f}_{\mathcal{I}}$  with  $\mathbf{f}'_{\mathcal{S}}$ .

**Definition 4.13** (Mechanism surgery). Let  $\mathcal{M}$  be the set of SCMs with elements  $\langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \_, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$  (note that the causal mechanisms are unspecified). Mechanism surgery is an operation  $I : \mathbf{X}_{\mathcal{I}}^{\mathbf{X}_{\mathcal{I}} \times \mathbf{E}_{\mathcal{J}}} \times \mathcal{M} \rightarrow \mathcal{M}$  that takes a causal model  $\mathcal{M}$  with arbitrary causal mechanisms and a set of causal mechanisms  $\mathbf{f}'_{\mathcal{I}}$  and maps it to a model  $\mathcal{M}' = \langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}'_{\mathcal{I}}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$ .

If  $\mathcal{M}$  has causal mechanisms  $\mathbf{f}_{\mathcal{I}}$  and  $\mathcal{O} \subset \mathcal{I}$  is the largest set such that  $\pi_{\mathcal{O}} \circ \mathbf{f}_{\mathcal{I}} = \pi_{\mathcal{O}} \circ \mathbf{f}'_{\mathcal{I}}$  then we say  $I$  is an *intervention* on  $\mathcal{L} := \mathcal{I} \setminus \mathcal{O}$ . We will use the special notation  $\mathcal{M}^{I(\mathcal{L}), \mathbf{f}'_{\mathcal{L}}} := I(\mathcal{M}, \mathbf{f}'_{\mathcal{L}})$  to denote an SCM related to  $\mathcal{M}$  by intervention on a subset of  $\mathcal{I}$ .

If furthermore  $\pi_{\mathcal{L}} \mathbf{f}'_{\mathcal{I}}$  is a constant function equal to  $\mathbf{a}$ , then we say  $I$  is a *hard intervention* on  $\mathcal{L}$ . We use the special notation  $\mathcal{M}^{Do(\mathcal{L})=\mathbf{a}} := I(\mathcal{M}, \mathbf{f}'_{\mathcal{L}})$  to denote SCMs related to  $\mathcal{M}$  by hard interventions. We also say that the *causal effect* of  $\mathcal{L}$  is the set of SCNs  $\{\mathcal{M}^{Do(\mathcal{L})=\mathbf{a}} | \mathbf{a} \in \mathbf{X}_{\mathcal{L}}\}$ .

We say a *causal model* is any kind of model that defines causal effects. An SCM  $\mathcal{M}$  in combination with hard interventions defines causal effects, so an SCM is a causal model. Call each interventional model  $\mathcal{M}^{do(\mathbf{X}_i=x)}$  a *submodel* of  $\mathcal{M}$ .

Strictly, the random variables  $\mathbf{X}_i$  depend on the probability space induced by a particular model  $\mathcal{M}$ , they are intended to refer to “the same variable” across different models that are related by mechanism surgery. We will abuse notation and use  $\mathbf{X}_i$  to refer to the *family* of random variables induced by a set of models related by mechanism surgery, and rely on explicitly noting the measure  $\mathbb{P} \dots$  (...) to specify exactly which random variables we are talking about.

In practice, we typically specify a “small” SCM containing a few endogenous variables  $\mathcal{I}$  (called a “marginal SCM” by Bongers et al. (2016)) which is understood to summarise the relevant characteristics of a “large” SCM containing many variables  $\mathcal{I}^*$ . We will argue that without restrictions on the large set of variables  $\mathcal{I}^*$ , surgically transformed SCMs will usually be invalid.

Incidentally, this messiness with random variables can be solved if we use See-Do models.

## 4.8 Not all variables have well-defined interventions

A long-running controversy about causal inference concerns the question of “the causal effect of body mass index on mortality”. On the one hand, Hernán and Taubman (2008) and others claim that there is no well-defined causal effect of a person’s body mass index (BMI), defined as their weight divided by their height, and their risk of death. Pearl claims, in defence of Causal Bayesian Networks, that the causal effect of *obesity* is well-defined, though it is not clear whether he defends the proposition that BMI itself has a causal effect:

That BMI is merely a coarse proxy of obesity is well taken; obesity should ideally be described by a vector of many factors, some are easy

to measure and others are not. But accessibility to measurement has no bearing on whether the effect of that vector of factors on morbidity is “well defined” or whether the condition of consistency is violated when we fail to specify the interventions used to regulate those factors. (Pearl, 2018b)

We argue that BMI does *not* have a well-defined causal effect, and without further assumptions neither does any variable.

#### 4.8.1 Necessary relationships in cyclic SCMs

If an SCM contains variables that are necessarily related, we wish to impose the additional restriction that these necessary relationships hold for every submodel. This can be done by extending the previous definition:

**Definition 4.14** (SCM with necessary relationships). An SCM with necessary relationships (SCNM) is a tuple  $\mathcal{M} := \langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}_{\mathcal{I}}, \mathbf{g}_{\mathcal{I}}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$ , which is an SCM with the addition of a vector function of *necessary relationships*  $\mathbf{g}_{\mathcal{I}} := \otimes_{i \in \mathcal{I}} g_i$  where each  $g_i : \mathbf{X}_{\mathcal{I}} \rightarrow X_i$  is a necessary relationship involving  $\mathbf{X}_i$ .

An SCM with necessary induces a unique probability space if for  $\mathbb{P}_{\mathcal{E}}$ -almost every  $e \in \mathcal{E}$  there exists a unique  $\mathbf{X} \in \mathbf{X}_{\mathcal{I}}$  such that simultaneously

$$\mathbf{X} = \mathbf{f}_{\mathcal{I}}(\mathbf{X}, \mathbf{e}) \quad (56)$$

$$\mathbf{X} = \mathbf{g}_{\mathcal{I}}(\mathbf{X}) \quad (57)$$

If no such  $\mathbf{X}$  exists then an SCNM is invalid.

Mechanism surgery for SCNMs involves modification of  $\mathbf{f}_{\mathcal{I}}$  only, just like SCMs.

If we wish to stipulate that a particular variable  $\mathbf{X}_i$  has no causal relationships or necessary relationships we can specify this with the trivial mechanisms  $f_i : (\mathbf{X}, \mathbf{e}) \mapsto x_i$  and  $g_i : \mathbf{X} \mapsto x_i$  respectively. An SCNM  $\mathcal{M}$  with the trivial necessary relationship  $\mathbf{g}_{\mathcal{I}} : \mathbf{X} \mapsto \mathbf{X}$  induces the equivalent probability spaces as the SCM obtained by removing  $\mathbf{g}_{\mathcal{I}}$  from  $\mathcal{M}$ .

Because BMI is always equal height/weight, given some SCNM  $\mathcal{M}$  containing endogenous variables  $\mathbf{X}_h$ ,  $\mathbf{X}_w$  and  $\mathbf{X}_b$  representing height, weight and BMI respectively, it should be possible to construct a more “primitive” SCNM  $\mathcal{M}^p$  containing every variable  $\mathcal{M}$  does except  $\mathbf{X}_b$  that agrees with  $\mathcal{M}$  on all interventions except those on  $\mathbf{X}_b$ .

**Definition 4.15** (Marginal model). Given an SCNM

$$\mathcal{M} = \langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}_{\mathcal{I}}, \mathbf{g}_{\mathcal{I}}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$$

a marginal model over  $\mathcal{L} \subset \mathcal{I}$  is an SCNM

$$\mathcal{M}^{*_{\mathcal{L}}} = \langle \mathcal{O}, \mathcal{J}, \mathbf{X}_{\mathcal{O}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}_{\mathcal{O}}^{\mathcal{L}}, \mathbf{g}_{\mathcal{O}}^{\mathcal{L}*}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$$

such that  $(\mathbb{P}_{\mathcal{M}})^*_{\mathcal{L}} = \mathbb{P}_{(\mathcal{M}^*_{\mathcal{L}})}$  and for all interventions  $\mathbf{f}'_{\mathcal{O}}$  on  $\mathcal{O} := \mathcal{I} \setminus \mathcal{L}$  that do not depend on  $\mathcal{L}$

$$(\mathbb{P}_{\mathcal{M}^{I(\mathcal{O}), \mathbf{f}'_{\mathcal{O}}}})^*_{\mathcal{L}} = \mathbb{P}_{(\mathcal{M}^*_{\mathcal{L}}, I(\mathcal{O}), \mathbf{f}'_{\mathcal{O}} \circ \pi_{\mathcal{O}})}$$

A *primitive model* is a special case of a marginal model where any intervention that depended only on endogenous variables in the original model can be replicated with some intervention that depends only on endogenous variables in the marginal model. If the endogenous variables represent *observed* variables, then the plausible intervention operations may only be allowed to depend on these variables. In general, there may be interventions that are possible in the original model that are not possible in the marginal model.

**Definition 4.16** (Primitive model). A *primitive model*  $\mathcal{M}^p$  is a marginal model of  $\mathcal{M}$  with respect to some  $\mathcal{L}$  such that for all interventions  $\mathbf{f}'_{\mathcal{O}}$  that do not depend on  $\mathcal{J}$  there exists some  $\mathbf{g}'_{\mathcal{O}} : \mathbf{X}_{\mathcal{O}} \times \mathbf{E}_{\mathcal{J}} \rightarrow \mathbf{X}_{\mathcal{O}}$  that does not depend on  $\mathcal{J}$  such that

$$(\mathbb{P}_{\mathcal{M}^{I(\mathcal{O}), \mathbf{f}'_{\mathcal{O}}}})^*_{\mathcal{L}} = \mathbb{P}_{(\mathcal{M}^*_{\mathcal{L}}, I(\mathcal{O}), \mathbf{g}'_{\mathcal{O}})}$$

We claim that given any SCNM  $\mathcal{M}$  containing endogenous variables  $\mathbf{X}_h$ ,  $\mathbf{X}_w$  and  $\mathbf{X}_b$  representing height, weight and BMI there should be a primitive model  $\mathcal{M}^p$  of  $\mathcal{M}$  with respect to  $\{p\}$ .

**Lemma 4.17** (Primitive models).  $\mathcal{M}^p$  is a primitive model of  $\mathcal{M}$  with respect to  $\mathcal{L} \subset \mathcal{I}$  iff  $S(\pi_{\mathcal{O}} \mathbf{f}_{\mathcal{I}}) \stackrel{a.s.}{=} S(\mathbf{f}_{\mathcal{O}}^p)$  for  $\mathcal{O} := \mathcal{I} \setminus \mathcal{L}$  and for all  $\mathbf{X} \in \mathbf{X}_{\mathcal{I}}$ ,  $\mathbf{g}$

However, as Theroem 4.19 shows, if an SCNM with height, weight and BMI can be derived from an SCNM containing just height and weight then there are no valid hard interventions on BMI.

**Definition 4.18** (Derived model). Given a SCNM  $\mathcal{M} := \langle \mathcal{I}, \mathcal{J}, \mathbf{X}_{\mathcal{I}}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}_{\mathcal{I}}, \mathbf{g}_{\mathcal{I}}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$ , say  $\mathcal{M}' = \langle \mathcal{I}', \mathcal{J}, \mathbf{X}_{\mathcal{I}'}, \mathbf{E}_{\mathcal{J}}, \mathbf{f}'_{\mathcal{I}'}, \mathbf{g}'_{\mathcal{I}'}, \mathbb{P}_{\mathcal{E}}, \mathbf{E}_{\mathcal{J}} \rangle$  is *derived* from  $\mathcal{M}$  if there exists some additional index/variable/relationships  $i' \notin \mathcal{I}$ ,  $X_{i'}$  such that

$$\mathcal{I}' = \mathcal{I} \cup \{i'\} \tag{58}$$

$$\mathbf{X}_{\mathcal{I}'} = \mathbf{X}_{\mathcal{I}} \cup X_{i'} \tag{59}$$

and, defining  $\pi_{\mathcal{I}' \setminus i'} : \mathbf{X}_{\mathcal{I}'} \rightarrow \mathbf{X}_{\mathcal{I}}$  as the projection map that “forgets”  $X_{i'}$ , for any  $\mathbf{e} \in \mathbf{E}_{\mathcal{J}}$  we have

$$\mathbf{X}' = \mathbf{f}'_{\mathcal{I}'}(\mathbf{X}', \mathbf{e}) \tag{60}$$

$$\text{and } \mathbf{X}' = \mathbf{g}'_{\mathcal{I}'}(\mathbf{X}') \implies \pi_{\mathcal{I}' \setminus i'}(\mathbf{X}') = \mathbf{f}_{\mathcal{I}}(\pi_{\mathcal{I}' \setminus i'}(\mathbf{X}'), \mathbf{e}) \tag{61}$$

$$\text{and } \pi_{\mathcal{I}' \setminus i'}(\mathbf{X}') = \mathbf{g}'_{\mathcal{I}'}(\pi_{\mathcal{I}' \setminus i'}(\mathbf{X}')) \tag{62}$$

**Theorem 4.19** (Interventions and necessary relationships don’t mix). If  $\mathcal{M}'$  is derived from  $\mathcal{M}$  with the additional elements  $i'$ ,  $X_{i'}$ ,  $f_{i'}$ ,  $g_{i'}$  and both  $\mathcal{M}$  and  $\mathcal{M}'$  are uniquely solvable and  $\mathbb{P}_{\mathcal{X}' \otimes \mathcal{E}}(X_{i'})$  is not single valued then no hard interventions on  $X_{i'}$  are possible.



*Proof.* Because  $\mathcal{M}$  is uniquely solvable, for  $\mathbb{P}_{\mathcal{E}}$  almost every  $\mathbf{e}$  there is a unique  $\mathbf{X}^e$  such that

$$\mathbf{X}^e = \mathbf{f}_{\mathcal{I}}(\mathbf{X}^e, \mathbf{e}) \quad (63)$$

$$\mathbf{X}^e = \mathbf{g}_{\mathcal{I}}(\mathbf{X}^e) \quad (64)$$

Because  $\mathcal{M}'$  is also uniquely solvable, for  $\mathbb{P}_{\mathcal{E}}$  almost every  $\mathbf{e}$  we have  $\mathbf{X}'^e \in \mathbf{X}_{\mathcal{I}'}$  such that  $\pi_{\mathcal{I}' \setminus i'}(\mathbf{X}')'^e = \mathbf{X}^e$  and

$$x_{i'}'^e = \mathbf{g}_{i'}(\mathbf{X}'^e) \quad (65)$$

Because  $\mathbb{P}_{\mathcal{X}' \otimes \mathcal{E}}(\mathbf{X}_{i'})$  is not single valued there are non-null sets  $A, B \in \mathcal{E}$  such that  $e_a \in A$ ,  $e_b \in B$  implies

$$\mathbf{g}_{i'}(\mathbf{X}'^{e_a}) \neq \mathbf{g}_{i'}(\mathbf{X}'^{e_b}) \quad (66)$$

Therefore there exists no  $a \in \mathcal{X}_{i'}$  that can simultaneously satisfy 65 for almost every  $\mathbf{e}$ . However, any hard intervention  $\mathcal{M}', do(\mathbf{X}_{i'}=a)$  requires such an  $a$  in order to be solvable.  $\square$

**Corollary 4.20.** *Either there are no hard interventions defined on BMI or there is no SCNM containing height and weight with a unique solution from which an SCNM containing height, weight and BMI can be derived.*

I can formalise the following, but I'm just writing it out so I can get to the end for now

The problem posed by Theorem 4.19 can be circumvented to some extent by joint interventions. Consider the variables  $\mathbf{X}_1$  and  $\mathbf{X}_2$  where  $\mathbf{X}_1 = -\mathbf{X}_2$  necessarily. While Theorem 4.19 disallows interventions on  $\mathbf{X}_2$  individually (supposing we can obtain a unique model featuring only  $\mathbf{X}_1$ ), it does not disallow interventions that jointly set  $\mathbf{X}_1$  and  $\mathbf{X}_2$  to permissible values. In this case, this is unproblematic as the only joint intervention that sets  $\mathbf{X}_1$  to 1 must also set  $\mathbf{X}_2$  to  $-1$ .

If we have non-invertible necessary relationships such as  $\mathbf{X}_1 = \mathbf{X}_2 + \mathbf{X}_3$ , however, there are now *multiple* joint interventions on  $\mathbf{X}_1$  that can be performed. I regard this as the most plausible solution to the difficulties raised so far: for variables that are in non-invertible necessary relationships, there is a set of operations associated with the “intervention” that sets  $\mathbf{X}_1 = 1$ .

However, we still need to make sure the interventions that we have supposed comprise the operations associated with setting  $\mathbf{X}_1 = 1$  exist themselves. It is sufficient that the SCNM with  $\mathbf{X}_1$  is derived from a higher order *uniquely solvable SCM* with  $\mathbf{X}_2$  and  $\mathbf{X}_3$  only .

And necessary? There might be “degenerate” necessary relationships that don't harm the possibility of defining interventions, and I'd need to show an equivalence to an SCM in this case

because interventions are defined in uniquely solvable SCMs and derivation preserves interventions on the old variables

If any variables are included in a causal model that are necessarily related to other variables (and honestly, is there any variable that isn't?), it is not enough to suppose that the model being used is a marginalisation of some larger causal model. Rather, it must be obtained by derivation and marginalisation from some model that represents the basic interventions that are possible, which we call the *atomic model*.

**Definition 4.21** (Atomic model). Given an SCNM  $\mathcal{M}$ , the *atomic model*  $\mathcal{M}_{\text{atom}}$  is a uniquely solvable SCM such that there exists a model  $\mathcal{M}$  is derived from of  $\mathcal{M}_{\text{atom}}$ .

Typically, in order to get an actually usable model you'll also need to marginalize, but I think this complication can be avoided

**Definition 4.22** (Causal universality hypothesis). There exists a uniquely solvable SCM  $\mathcal{M}_{\text{atom}}$  which is the atomic model that correctly represents all decision problems

what does that mean?

I don't know how to define "correctly represents" or "causal problem", but it seems like something like the universality hypothesis is necessary if you want to define "the causal effect of X" independent of any atomic model

or causal problems?

Relate decisions to interventions on atomic model. Decisions  $\rightarrow$  atomic model is straightforward, but the reverse direction is not so obvious

Causal effects are uniquely defined via atoms iff they are defined via decisions

Are there any plausible ways to construct atomic models?

## 5 Definitions and key notation

We use three notations for working with probability theory. The "elementary" notation makes use of regular symbolic conventions (functions, products, sums, integrals, unions etc.) along with the expectation operator  $\mathbb{E}$ . This is the most flexible notation which comes at the cost of being verbose and difficult to read. Secondly, we use a semi-formal string diagram notation extending the formal diagram notation for symmetric monoidal categories Selinger (2010). Objects in this diagram refer to stochastic maps, and by interpreting diagrams as symbols we can, in theory, be just as flexible as the purely symbolic approach. However, we avoid complex mixtures of symbols and diagrams elements, and fall back to symbolic representations if it is called for. Finally, we use a matrix-vector product convention that isn't particularly expressive but can compactly express some common operations.

## 5.1 Standard Symbols

Symbol	Meaning
$[n]$	The natural numbers $\{1, \dots, n\}$
$f : a \mapsto b$	Function definition, equivalent to $f(a) := b$
Dots appearing in function arguments: $f(\cdot, \cdot, z)$	The “curried” function $(x, y) \mapsto f(x, y, z)$
Capital letters: $A, B, X$	sets
Script letters: $\mathcal{A}, \mathcal{B}, \mathcal{X}$	$\sigma$ -algebras on the sets $A, B, X$ respectively
Script $\mathcal{G}$	A directed acyclic graph made up of nodes $V$ and edges
Greek letters $\mu, \xi, \gamma$	Probability measures
$\delta_x$	The Dirac delta measure: $\delta_x(A) = 1$ if $x \in A$ and 0 otherwise
Capital delta: $\Delta(\mathcal{E})$	The set of all probability measures on $\mathcal{E}$
Bold capitals: $\mathbf{A}$	Markov kernel $\mathbf{A} : X \times \mathcal{Y} \rightarrow [0, 1]$ (stochastic map)
Subscripted bold capitals: $\mathbf{A}_x$	The probability measure given by the curried Markov kernel $\mathbf{A}_x$
$A \rightarrow \Delta(\mathcal{B})$	Markov kernel signature, treated as equivalent to $A \times \mathcal{B}$
$\mathbf{A} : x \mapsto \nu$	Markov kernel definition, equivalent to $\mathbf{A}(x, B) = \nu(B)$ for all $B \in \mathcal{B}$
Sans serif capitals: $A, X$	Measurable functions; we will also call them random variables
$\mathbf{F}_X$	The Markov kernel associated with the function $X$ : $\mathbf{F}_X \equiv \mathbf{A}_X$
$\mathbf{N}_{A B}$	The conditional probability (disintegration) of $\mathbf{A}$ given $B$
$\nu \mathbf{F}_X$	The marginal distribution of $X$ under $\nu$

## 5.2 Probability Theory

Given a set  $A$ , a  $\sigma$ -algebra  $\mathcal{A}$  is a collection of subsets of  $A$  where

- $A \in \mathcal{A}$  and  $\emptyset \in \mathcal{A}$
- $B \in \mathcal{A} \implies B^C \in \mathcal{A}$
- $\mathcal{A}$  is closed under countable unions: For any countable collection  $\{B_i | i \in \mathbb{N}\}$  of elements of  $\mathcal{A}$ ,  $\cup_{i \in \mathbb{N}} B_i \in \mathcal{A}$

A measurable space  $(A, \mathcal{A})$  is a set  $A$  along with a  $\sigma$ -algebra  $\mathcal{A}$ . Sometimes the sigma algebra will be left implicit, in which case  $A$  will just be introduced as a measurable space.

**Common  $\sigma$  algebras** For any  $A$ ,  $\{\emptyset, A\}$  is a  $\sigma$ -algebra. In particular, it is the only sigma algebra for any one element set  $\{*\}$ .

For countable  $A$ , the power set  $\mathcal{P}(A)$  is known as the discrete  $\sigma$ -algebra.

Given  $A$  and a collection of subsets of  $B \subset \mathcal{P}(A)$ ,  $\sigma(B)$  is the smallest  $\sigma$ -algebra containing all the elements of  $B$ .

Let  $T$  be all the open subsets of  $\mathbb{R}$ . Then  $\mathcal{B}(\mathbb{R}) := \sigma(T)$  is the *Borel  $\sigma$ -algebra* on the reals. This definition extends to an arbitrary topological space  $A$  with topology  $T$ .

A *standard measurable set* is a measurable set  $A$  that is isomorphic either to a discrete measurable space  $A$  or  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ . For any  $A$  that is a complete separable metric space,  $(A, \mathcal{B}(A))$  is standard measurable.

Given a measurable space  $(E, \mathcal{E})$ , a map  $\mu : \mathcal{E} \rightarrow [0, 1]$  is a *probability measure* if

- $\mu(E) = 1, \mu(\emptyset) = 0$
- Given countable collection  $\{A_i\} \subset \mathcal{E}$ ,  $\mu(\cup_i A_i) = \sum_i \mu(A_i)$

Write by  $\Delta(\mathcal{E})$  the set of all probability measures on  $\mathcal{E}$ .

Given a second measurable space  $(F, \mathcal{F})$ , a *stochastic map* or *Markov kernel* is a map  $\mathbf{M} : E \times \mathcal{F} \rightarrow [0, 1]$  such that

- The map  $\mathbf{M}(\cdot; A) : x \mapsto \mathbf{M}(x; A)$  is  $\mathcal{E}$ -measurable for all  $A \in \mathcal{F}$
- The map  $\mathbf{M}_x : A \mapsto \mathbf{M}(x; A)$  is a probability measure on  $F$  for all  $x \in E$

Extending the subscript notation above, for  $\mathbf{C} : X \times Y \rightarrow \Delta(\mathcal{Z})$  and  $x \in X$  we will write  $\mathbf{C}_x$  for the “curried” map  $y \mapsto \mathbf{C}_{x,y}$ .

The map  $x \mapsto \mathbf{M}_x$  is of type  $E \rightarrow \Delta(\mathcal{F})$ . We will abuse notation somewhat to write  $\mathbf{M} : E \rightarrow \Delta(\mathcal{F})$ , which captures the intuition that a Markov kernel maps from elements of  $E$  to probability measures on  $\mathcal{F}$ . Note that we “reverse” this idea and consider Markov kernels to map from elements of  $\mathcal{F}$  to measurable functions  $E \rightarrow [0, 1]$ , an interpretation found in Clerc et al. (2017), but (at this stage) we don’t make use of this interpretation here.

Given an indiscrete measurable space  $(\{*\}, \{\{*\}, \emptyset\})$ , we identify Markov kernels  $\mathbf{N} : \{*\} \rightarrow \Delta(\mathcal{E})$  with the probability measure  $\mathbf{N}_*$ . In addition, there is a unique Markov kernel  $*$  :  $E \rightarrow \Delta(\{\{*\}, \emptyset\})$  given by  $x \mapsto \delta_*$  for all  $x \in E$  which we will call the “discard” map.

### 5.3 Product Notation

We can use a notation similar to the standard notation for matrix-vector products to represent operations with Markov kernels. Probability measures  $\mu \in \Delta(\mathcal{X})$  can be read as row vectors, Markov kernels as matrices and measurable functions  $\mathbf{T} : Y \rightarrow T$  as column vectors. Defining  $\mathbf{M} : X \rightarrow \Delta(\mathcal{Y})$  and  $\mathbf{N} : Y \rightarrow \Delta(\mathcal{Z})$ , the measure-kernel product  $\mu \mathbf{A}(G) := \int \mathbf{A}_x(G) d\mu(x)$  yields a probability measure  $\mu \mathbf{A}$  on  $\mathcal{Z}$ , the kernel-kernel product  $\mathbf{M} \mathbf{N}(x; H) = \int_Y \mathbf{B}(y; H) d\mathbf{A}_x$  yields a kernel  $\mathbf{M} \mathbf{N} : X \rightarrow \Delta(\mathcal{Z})$  and the kernel-function product  $\mathbf{A} \mathbf{T}(x) := \int_Y \mathbf{T}(y) d\mathbf{A}_x$  yields a measurable function  $X \rightarrow T$ . Kernel products are associative (Çinlar, 2011).

The tensor product  $(\mathbf{M} \otimes \mathbf{N})(x, y; G, H) := \mathbf{M}(x; G) \mathbf{N}(y; H)$  yields a kernel  $(\mathbf{M} \otimes \mathbf{N}) : X \times Y \rightarrow \Delta(\mathcal{Y} \otimes \mathcal{Z})$ .

### 5.4 String Diagrams

Some constructions are unwieldy in product notation; for example, given  $\mu \in \Delta(\mathcal{E})$  and  $\mathbf{M} : E \rightarrow (\mathcal{F})$ , it is not straightforward to construct a measure  $\nu \in \Delta(\mathcal{E} \otimes \mathcal{F})$  that captures the “joint distribution” given by  $A \times B \mapsto \int_A \mathbf{M}(x; B) d\mu$ .

Such constructions can, however, be straightforwardly captured with string diagrams, a notation developed for category theoretic probability. Cho and Jacobs (2019) also provides an extensive introduction to the notation discussed here.

Some key ideas of string diagrams:

- Basic string diagrams can always be interpreted as a mixture of kernel-kernel products and tensor products of Markov kernels
  - Extended string diagrams can be interpreted as a mixture of kernel-kernel products, kernel-function products, tensor products of kernels and functions and scalar products
- String diagrams are the subject of a coherence theorem: taking a string diagram and applying a planar deformation yields a string diagram that represents the same kernel (Selinger, 2010). This also holds for a number of additional transformations detailed below

A kernel  $\mathbf{M} : X \rightarrow \Delta(\mathcal{Y})$  is written as a box with input and output wires, probability measures  $\mu \in \Delta(\mathcal{X})$  are written as triangles “closed on the left” and measurable functions (which are only elements of the “extended” notation)  $T : Y \rightarrow T$  as triangles “closed on the right”. For this introduction we will label wires with the names of their corresponding spaces, but in practice we will usually name them with corresponding *random variables*, though additional care is required when using random variables as labels (see paragraph ??).

For  $\mathbf{M} : X \rightarrow \Delta(\mathcal{Y})$ ,  $\mu \in \Delta(\mathcal{X})$  and  $f : X \rightarrow W$ :

$$X \text{ --- } \boxed{\mathbf{M}} \text{ --- } Y \quad \triangleleft_{\mu} \text{ --- } X \quad X \text{ --- } \triangleright_f \quad (67)$$

**Elementary operations** We can compose Markov kernels with appropriate spaces - the equivalent operation of the “matrix products” of product notation. Given  $\mathbf{M} : X \rightarrow \Delta(\mathcal{Y})$  and  $\mathbf{N} : Y \rightarrow \Delta(\mathcal{Z})$ , we have

$$\mathbf{MN} := X \text{ --- } \boxed{\mathbf{M}} \text{ --- } \boxed{\mathbf{N}} \text{ --- } Z \quad (68)$$

Probability measures are distinguished in that that they only admit “right composition” while functions only admit “left composition”. For  $\mu \in \Delta(\mathcal{E})$ ,  $h : F \rightarrow X$ :

$$\mu \mathbf{M} := \triangleleft_{\mu} \text{ --- } \boxed{\mathbf{M}} \text{ --- } Z \quad (69)$$

$$\mathbf{M} f := X \text{ --- } \boxed{\mathbf{M}} \text{ --- } \triangleright_f \quad (70)$$

A diagram that is closed on the right and the left is an expectation:

$$\mathbb{E}_{\mu\mathbf{M}}(f) = \mu\mathbf{M}f \quad (71)$$

$$:= \triangleleft \mu \text{---} \boxed{\mathbf{M}} \text{---} f \triangleright \quad (72)$$

We can also combine Markov kernels using tensor products, which we represent with vertical juxtaposition. For  $\mathbf{O} : Z \rightarrow \Delta(\mathcal{W})$ :

$$\mathbf{M} \otimes \mathbf{N} := \begin{array}{c} X \text{---} \boxed{\mathbf{M}} \text{---} Y \\ Z \text{---} \boxed{\mathbf{O}} \text{---} W \end{array} \quad (73)$$

Product spaces can be represented either by two parallel wires or a single wire:

$$X \times Y \cong \text{Id}_X \otimes \text{Id}_Y := \begin{array}{c} X \text{---} X \\ Y \text{---} Y \end{array} \quad (74)$$

$$= X \times Y \text{---} X \times Y \quad (75)$$

Because a product space can be represented by parallel wires, a kernel  $\mathbf{L} : X \rightarrow \Delta(\mathcal{Y} \otimes \mathcal{Z})$  can be written using either two parallel output wires or a single output wire:

$$X \text{---} \boxed{\mathbf{L}} \text{---} \begin{array}{c} Y \\ Z \end{array} \quad (76)$$

$$\equiv \quad (77)$$

$$X \text{---} \boxed{\mathbf{L}} \text{---} Y \times Z \quad (78)$$

**Probability measures, Markov kernels and functions** One has to exercise special care when including functions in diagrammatic notation. While any diagram that includes only probability measures (triangles pointing to the left) and Markov kernels (rectangles) is automatically a Markov kernel itself, while diagrams that include functions (triangles pointing to the right) only represent Markov kernels if they are correctly normalised, which is not a property that can be checked just by looking at the shape of the diagram.

**Markov kernels with special notation** A number of Markov kernels are given special notation distinct from the generic “box” representation above. These special representations facilitate intuitive graphical interpretations.

The identity kernel  $\mathbf{Id} : X \rightarrow \Delta(X)$  maps a point  $x$  to the measure  $\delta_x$  that places all mass on the same point:

$$\mathbf{Id}_x : x \mapsto \delta_x \equiv X \text{ --- } X \quad (79)$$

The copy map  $\Upsilon : X \rightarrow \Delta(\mathcal{X} \times \mathcal{X})$  maps a point  $x$  to two identical copies of  $x$ :

$$\Upsilon : x \mapsto \delta_{(x,x)} \equiv X \text{ --- } \begin{array}{c} X \\ X \end{array} \quad (80)$$

The swap map  $\sigma : X \times Y \rightarrow \Delta(\mathcal{Y} \otimes \mathcal{X})$  swaps its inputs:

$$\sigma := (x, y) \mapsto \delta_{(y,x)} \equiv \begin{array}{c} Y \\ X \end{array} \text{ --- } \begin{array}{c} X \\ Y \end{array} \quad (81)$$

The discard map  $*$  :  $X \rightarrow \Delta(\{*\})$  maps every input to  $\delta_*$ . Note that the only non-empty event in  $\{\emptyset, \{*\}\}$  must have probability 1.

$$* : x \mapsto \delta_* \equiv X \text{ --- } * \quad (82)$$

Any measurable function  $F \rightarrow X$  has an associated Markov kernel  $F \rightarrow \Delta(\mathcal{X})$ . The Markov kernel associated with a function is different to the function itself - while the product of a probability measure  $\mu$  with a function  $f$  is an expectation  $\mu f$  (see Definition 72), the product of a probability measure with the associated Markov kernel is the pushforward measure  $f_{\#}\mu$ .

**Definition 5.1** (Function induced kernel). Given a measurable function  $g : F \rightarrow X$ , define the function induced kernel  $\mathbf{F}_g : F \rightarrow \Delta(\mathcal{X})$  to be the Markov kernel  $a \mapsto \delta_{g(a)}$  for all  $a \in X$ .

**Definition 5.2** (Pushforward kernel). Given a kernel  $\mathbf{M} : E \rightarrow \Delta(\mathcal{F})$  and a measurable function  $g : F \rightarrow X$ , the *pushforward kernel*  $g_{\#}\mathbf{M} : E \rightarrow \Delta(\mathcal{X})$  is the kernel such that  $g_{\#}\mathbf{M}(a; B) = \mathbf{M}(a; g^{-1}(B))$ .

If  $E$  is the one element space  $\{*\}$ , then  $\mathbf{M} : \{*\} \rightarrow \Delta(\mathcal{F})$  can be identified with the probability measure  $\mathbf{M}_*$  and the pushforward kernel  $g_{\#}\mathbf{M}$  identified with the pushforward measure  $g_{\#}\mathbf{M}_*$ , so pushforward kernels reduce to pushforward measures.

**Lemma 5.3** (Pushforward kernels are functional kernel products). *Given a kernel  $\mathbf{M} : E \rightarrow \Delta(\mathcal{F})$  and a measurable function  $g : F \rightarrow X$ , the pushforward  $g_{\#}\mathbf{M} = \mathbf{M}\mathbf{F}_g$ .*

*Proof.*

$$\mathbf{MF}_g(a; B) = \int_F \delta_{g(y)}(B) d\mathbf{M}_a(y) \quad (83)$$

$$= \int_F \delta_y(g^{-1}(B)) d\mathbf{M}_a(y) \quad (84)$$

$$= \int_{g^{-1}(B)} d\mathbf{M}_a(y) \quad (85)$$

$$= g_{\#} \mathbf{M}(a; B) \quad (86)$$

□

#### 5.4.1 Comparison of notations

We are in a position to compare the three introduced notations using a few examples. Given  $\mu \in \Delta(X)$ ,  $\mathbf{A} : X \rightarrow \Delta(Y)$  and  $A \in \mathcal{X}$ ,  $B \in \mathcal{Y}$ , the following correspondences hold, where we express the same object in elementary notation, product notation and string notation respectively:

$$\nu := A \times B \mapsto \int_A A(x; B) d\mu(x) \equiv \mu \curlyvee (\mathbf{Id}_X \otimes \mathbf{A}) \equiv \begin{array}{c} \text{---} X \\ \swarrow \quad \searrow \\ \triangleleft \mu \quad \boxed{\mathbf{A}} \text{---} Y \end{array} \quad (87)$$

Where the resulting object is a probability measure  $\nu \in \Delta(\mathcal{X} \otimes \mathcal{Y})$ . Note that the elementary notation requires a function definition here, while the product and string notations can represent the measure without explicitly addressing its action on various inputs and outputs. Cho and Jacobs (2019) calls this construction “integrating  $\mathbf{A}$  with respect to  $\mu$ ”.

Define the marginal  $\nu_Y \in \Delta(\mathcal{Y}) : B \mapsto \nu(X \times B)$  for  $B \in \mathcal{Y}$  and similarly for  $\nu_X$ . We can then express the result of marginalising 87 over  $X$  in our three separate notations as follows:

$$\nu_Y(B) = \nu(X \times B) = \int_X A(x; B) d\mu(x) \quad (88)$$

$$\nu_Y = \mu \mathbf{A} = \mu \curlyvee (\mathbf{Id}_X \otimes \mathbf{A}) (* \otimes \mathbf{Id}_Y) \quad (89)$$

$$\nu_Y = \begin{array}{c} \text{---} * \\ \swarrow \quad \searrow \\ \triangleleft \mu \quad \boxed{\mathbf{A}} \text{---} Y \end{array} = \begin{array}{c} \triangleleft \mu \text{---} \boxed{\mathbf{A}} \text{---} Y \end{array} \quad (90)$$

The elementary notation 88 makes the relationship between  $\nu_Y$  and  $\nu$  explicit and, again, requires the action on each event to be defined. The product notation 89 is, in my view, the least transparent but also the most compact in the form  $\mu \mathbf{A}$ , and does not demand the explicit definition of how  $\nu_Y$  treats every event. The graphical notation is the least compact in terms of space taken



up on the page, but unlike the product notation it shows a clear relationship to the graphical construction in 87, and displays a clear graphical logic whereby marginalisation corresponds to “cutting off branches”. Like product notation, it also allows for the definition of derived measures such as  $\nu_Y$  without explicit definition of the handling of all events. It also features a much smaller collection of symbols than does elementary notation.

String diagrams often achieve a good balance between being ease of understanding at a glance and expressive power. On the downside, they can be time consuming to typeset, and formal reasoning with them takes some practice.

#### 5.4.2 Working With String Diagrams

todo:

- Functional generalisation
- Conditioning
- Infinite copy map
- De Finetti’s representation theorem

There are a relatively small number of manipulation rules that are useful for string diagrams. In addition, we will define graphically analogues of the standard notions of *conditional probability*, *conditioning*, and infinite sequences of exchangeable random variables.

**Axioms of Symmetric Monoidal Categories** For the following, we either omit labels or label diagrams with their domain and codomain spaces, as we are discussing identities of kernels rather than identities of components of a conditional probability space. Recalling the unique Markov kernels defined above, the following equivalences, known as the *commutative comonoid axioms*, hold among string diagrams:

$$\begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} := \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} \quad (91)$$

$$\begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} = \text{---} \quad (92)$$

$$\begin{array}{c} \text{X} \text{---} \text{---} \text{---} \\ \text{X} \text{---} \text{---} \text{---} \\ \text{X} \text{---} \text{---} \text{---} \end{array} = \begin{array}{c} \text{X} \text{---} \text{---} \text{---} \\ \text{X} \text{---} \text{---} \text{---} \\ \text{X} \text{---} \text{---} \text{---} \end{array} \quad (93)$$

The discard map  $*$  can “fall through” any Markov kernel:

$$\text{---} \boxed{\mathbf{A}} \text{---} * = \text{---} * \quad (94)$$

Combining 92 and 94 we can derive the following: integrating  $\mathbf{A} : X \rightarrow \Delta(\mathcal{Y})$  with respect to  $\mu \in \Delta(\mathcal{X})$  and then discarding the output of  $\mathbf{A}$  leaves us with  $\mu$ :

$$\quad (95)$$

In elementary notation, this is equivalent to the fact that, for all  $B \in \mathcal{X}$ ,  $\int_B \mathbf{A}(x; B) d\mu(x) = \mu(B)$ .

The following additional properties hold for  $*$  and  $\curlyvee$ :

$$X \times Y \text{---} * = \begin{matrix} X \text{---} * \\ Y \text{---} * \end{matrix} \quad (96)$$

$$\quad (97)$$

A key fact that *does not* hold in general is

$$\quad (98)$$

In fact, it holds only when  $\mathbf{A}$  is a *deterministic* kernel.

**Definition 5.4** (Deterministic Markov kernel). A *deterministic* Markov kernel  $\mathbf{A} : E \rightarrow \Delta(\mathcal{F})$  is a kernel such that  $\mathbf{A}_x(B) \in \{0, 1\}$  for all  $x \in E$ ,  $B \in \mathcal{F}$ .

**Theorem 5.5** (Copy map commutes for deterministic kernels (Fong, 2013)). Equation 98 holds iff  $\mathbf{A}$  is deterministic.

## 5.5 Random Variables

The summary of this section is:

- Random variables are usually defined as measurable functions on a *probability space*

- It's possible to define them as measurable functions on a *Markov kernel space* instead
- It is useful to labelling wires with random variable names instead of names of spaces

Probability theory is primarily concerned with the behaviour of *random variables*. This behaviour can be analysed via a collection of probability measures and Markov kernels representing joint, marginal and conditional distributions of random variables of interest. In the framework developed by Kolmogorov, this collection of joint, marginal and conditional distributions is modeled by a single underlying *probability space*, and random variables by measurable functions on the probability space.

We use the same approach here, with a couple of additions. We are interested in variables whose outcomes depend both on random processes and decisions. These variables are better modelled by a Markov kernels than probability measure - *given* a particular decision, they inherit a particular probability distribution. Thus, variables in our work are modeled by an underlying Markov kernel rather than a probability measure; we call this a *Markov kernel space*.

In addition to following standard conventions regarding the use of random variables, we can motivate their introduction with a graphical example. Suppose we have some  $\mu \in \Delta(\mathcal{X} \otimes \mathcal{X})$ ,  $\mathbf{K} : \mathcal{X} \rightarrow \Delta(\mathcal{X})$  such that the following holds:

$$\triangleleft_{\mu} \frac{X}{X} = \triangleleft_{\mu} \text{---} \overline{\text{---} \boxed{\mathbf{K}} \text{---}} \frac{X}{X} \quad (99)$$

This implies, roughly, that  $\mathbf{K}$  is the probability of the lower wire of  $\mu$  conditional on the upper (it is a *disintegration* of  $\mu$ , defined later). However, it is very cumbersome to write 99 and define  $\mathbf{K}$  in terms of the geometry of the diagrams. Instead, it would be nice to have a system where we can unambiguously assign *names* to wires:

$$\triangleleft_{\mu} \frac{X_1}{X_2} \quad (100)$$

Once the wires of  $\mu$  have names, we can define a convention such that:

$$X_1 \text{---} \overline{\text{---} \boxed{\gamma|X} \text{---}} Y_2 \quad (101)$$

Is a Markov kernel  $X \rightarrow \Delta(\mathcal{X})$  that satisfies Equation 99 when substituted for  $\mathbf{K}$ . We take wire names to stand for random variables on a *Markov kernel space*.

**Definition 5.6** (Probability space, Markov kernel space). A *probability space*  $(\mathbb{P}, \Omega, \mathcal{F})$  is a probability measure  $\mathbb{P}$ , which we call the *ambient measure*, along with the *sample space*  $\Omega$  and the *events*  $\mathcal{F}$ .

A *Markov kernel space*  $(\mathbb{K}, \Omega, \mathcal{F}, D, \mathcal{D})$  is a Markov kernel  $\mathbb{K} : D \rightarrow \Delta(\mathcal{D} \otimes \mathcal{F})$ , called the *ambient kernel*, along with the sample space  $(\Omega, \mathcal{F})$  and the domain  $(D, \mathcal{D})$ . We suppose that  $\mathbb{K}$  is such that there exists a *fundamental kernel*  $\mathbb{K}_0$  satisfying

$$\mathbb{K} := \text{---} \boxed{\mathbb{K}_0} \text{---} \quad (102)$$

It is in general much more practical to work with  $\mathbb{K}$  than  $\mathbb{K}_0$ .

Is this sufficient to make kernel spaces a special case of conditional probability spaces?

**Definition 5.7** (Random variable). Given a sample space  $\Omega$  and domain  $D$ , a random variable  $X$  is a measurable function  $\Omega \times D \rightarrow E$  for arbitrary measurable  $E$ .

**Definition 5.8** (Domain variable). Given a sample space  $\Omega$  and an domain  $D$ , the *domain variable*  $D : \Omega \times D \rightarrow D$  is the random variable given by  $D : (x, d) \mapsto d$ .

Unlike random variables on probability spaces, random variables on Markov kernel spaces do not in general have unique marginal distributions. An analogous operation of *marginalisation* can be defined, but the result is generally a Markov kernel.

**Definition 5.9** (Coupled tensor product  $\otimes$ ). Given two Markov kernels  $\mathbf{M}$  and  $\mathbf{N}$  or functions  $f$  and  $g$  with shared domain  $E$ , let  $\mathbf{M} \otimes \mathbf{N} := \vee(\mathbf{M} \otimes \mathbf{N})$  and  $f \otimes g := \vee(f \otimes g)$  where these expressions are interpreted using standard product notation. Graphically:

$$\mathbf{M} \otimes \mathbf{N} := \begin{array}{c} E \text{---} \boxed{\mathbf{M}} \text{---} X \\ \quad \quad \quad \boxed{\mathbf{N}} \text{---} Y \end{array} \quad (103)$$

$$f \otimes g := \begin{array}{c} E \text{---} \triangle f \\ \quad \quad \quad \triangle g \end{array} \quad (104)$$

The operation denoted by  $\otimes$  is associative (Lemma 5.10), so we can without ambiguity write  $f \otimes g \otimes h = (f \otimes g) \otimes h = f \otimes (g \otimes h)$  for finite groups of functions or Markov kernels sharing a domain.

The notation  $\otimes_{i \in [N]} f_i$  is taken to mean  $f_1 \otimes f_2 \otimes \dots \otimes f_N$ .

**Lemma 5.10** ( $\otimes$  is associative). For Markov kernels  $\mathbf{L} : E \rightarrow \delta(\mathcal{F})$ ,  $\mathbf{M} : E \rightarrow \delta(\mathcal{G})$  and  $\mathbf{N} : E \rightarrow \delta(\mathcal{H})$ ,  $(\mathbf{L} \otimes \mathbf{M}) \otimes \mathbf{N} = \mathbf{L} \otimes (\mathbf{M} \otimes \mathbf{N})$ .

*Proof.*

$$\mathbf{L} \underline{\otimes} (\mathbf{M} \underline{\otimes} \mathbf{N}) = \begin{array}{c} \begin{array}{c} E \text{ --- } \begin{array}{c} \boxed{\mathbf{L}} \text{ --- } F \\ \boxed{\mathbf{M}} \text{ --- } G \\ \boxed{\mathbf{N}} \text{ --- } H \end{array} \end{array} \end{array} \quad (105)$$

$$= \begin{array}{c} \begin{array}{c} E \text{ --- } \begin{array}{c} \boxed{\mathbf{L}} \text{ --- } F \\ \boxed{\mathbf{M}} \text{ --- } G \\ \boxed{\mathbf{N}} \text{ --- } H \end{array} \end{array} \end{array} \quad (106)$$

$$= (\mathbf{L} \underline{\otimes} \mathbf{M}) \underline{\otimes} \mathbf{N} \quad (107)$$

This follows directly from Equation 91.  $\square$

**Definition 5.11** (Marginal distribution, marginal kernel). Given  $\mathbb{P} \in \Delta(\mathcal{F})$ , random variable  $\mathbf{X} : \Omega \rightarrow G$  the *marginal distribution* of  $\mathbf{X}$   $\mathbb{P}_{\mathbf{X}} \in \Delta(\mathcal{G})$  of  $\mathbf{X}$  is the product measure  $\mathbb{P}\mathbf{F}_{\mathbf{X}}$ .

See Lemma 5.3 for the proof that this matches the usual definition of marginal distribution.

Given  $\mathbb{K} : D \rightarrow \Delta(\mathcal{F})$  and random variable  $\mathbf{X} : \Omega \rightarrow G$ , the *marginal kernel* is  $\mathbb{K}_{\mathbf{X}|D} := \mathbb{K}\mathbf{F}_{\mathbf{X}}$ .

**Definition 5.12** (Joint distribution, joint kernel). Given  $\mathbb{P} \in \Delta(\mathcal{F})$ ,  $\mathbf{X} : \Omega \rightarrow G$  and  $\mathbf{Y} : \Omega \rightarrow H$ , the *joint distribution*  $\mathbb{P}_{\mathbf{X}\mathbf{Y}} \in \Delta(\mathcal{G} \otimes \mathcal{H})$  of  $\mathbf{X}$  and  $\mathbf{Y}$  is the marginal distribution of  $\mathbf{X} \underline{\otimes} \mathbf{Y}$ .

This is identical to the definition in Çinlar (2011) if we note that the random variable  $(\mathbf{X}, \mathbf{Y}) : \omega \mapsto (\mathbf{X}(\omega), \mathbf{Y}(\omega))$  (Çinlar's definition) is precisely the same thing as  $\mathbf{X} \underline{\otimes} \mathbf{Y}$ .

Analogously, the joint kernel  $\mathbb{K}_{\mathbf{X}\mathbf{Y}|D}$  is the product  $\mathbb{K}\mathbf{F}_{\mathbf{X} \underline{\otimes} \mathbf{Y}}$ .

Joint distributions and kernels have a nice visual representation, as a result of Lemma 5.13 which follows.

**Lemma 5.13** (Dual representation of coupled products of functions). *Given two functions, the kernel associated with their coupled product is equal to the coupled product of the kernels associated with each function.*

*Given  $\mathbf{X} : \Omega \rightarrow G$  and  $\mathbf{Y} : \Omega \rightarrow H$ ,  $\mathbf{F}_{\mathbf{X} \underline{\otimes} \mathbf{Y}} = \mathbf{F}_{\mathbf{X}} \underline{\otimes} \mathbf{F}_{\mathbf{Y}}$*

*Proof.* For  $a \in \Omega$ ,  $B \in \mathcal{G}$ ,  $C \in \mathcal{H}$ ,

$$\mathbf{F}_{\mathbf{X} \underline{\otimes} \mathbf{Y}}(a; B \times C) = \delta_{\mathbf{X}(a), \mathbf{Y}(a)}(B \times C) \quad (108)$$

$$= \delta_{\mathbf{X}(a)}(B) \delta_{\mathbf{Y}(a)}(C) \quad (109)$$

$$= (\delta_{\mathbf{X}(a)} \otimes \delta_{\mathbf{Y}(a)})(B \times C) \quad (110)$$

$$= \mathbf{F}_{\mathbf{X}} \underline{\otimes} \mathbf{F}_{\mathbf{Y}} \quad (111)$$

Equality follows from the monotone class theorem.  $\square$

**Corollary 5.14.** *Given a Markov kernel space  $(\Omega, D, \mathbb{K})$  and random variables  $X : \Omega \times D \rightarrow X$ ,  $Y : \Omega \times D \rightarrow Y$ , the following holds:*

$$D - \boxed{\mathbb{K}_{XY|D}} - \begin{array}{c} X \\ Y \end{array} = D - \boxed{\mathbb{K}} - \left( \begin{array}{c} \boxed{\mathbf{F}_X} - X \\ \boxed{\mathbf{F}_Y} - Y \end{array} \right) \quad (112)$$

We will now define wire labels for “output” wires.

**Definition 5.15** (Wire labels - joint kernels). Suppose we have a Markov kernel space  $(\Omega, D, \mathbb{K})$  and random variables  $X : \Omega \times D \rightarrow X$ ,  $Y : \Omega \times D \rightarrow Y$ , and a Markov kernel  $\mathbf{L} : D \rightarrow \Delta(\mathcal{X} \times \mathcal{Y})$ .

Relative to  $(\Omega, D, \mathbb{K})$ , the wires terminating on a free end on the right of a diagram of  $\mathbf{L}$  may be labelled with  $\mathbf{X}$  and  $\mathbf{Y}$  as follows:

$$D - \boxed{\mathbf{L}} - \begin{array}{c} \mathbf{X} \\ \mathbf{Y} \end{array} \quad (113)$$

iff

$$\mathbf{L} = \mathbb{K}_{XY|D} \quad (114)$$

and

$$D - \boxed{\mathbf{L}} - \begin{array}{c} \mathbf{X} \\ * \end{array} = \mathbb{K}_{X|D} \quad (115)$$

and

$$D - \boxed{\mathbf{L}} - \begin{array}{c} * \\ \mathbf{Y} \end{array} = \mathbb{K}_{Y|D} \quad (116)$$

The second and third conditions are nontrivial: suppose  $X$  takes values in some product space  $\text{Range}(X) = W \times Z$ , and  $Y$  takes values in  $Y$ . Then we could have  $\mathbf{L} = \mathbb{K}_{XY|D}$  and draw the diagram

$$D - \boxed{\mathbf{L}} - \begin{array}{c} W \\ Z \times Y \end{array} \quad (117)$$

For *this* diagram, properties 115 and 116 do not hold, even though 114 does.

I need to prove that if 114 holds and the spaces match the codomains of the random variables, then labels can be assigned

Having defined output wire labels, I will now proceed to use them without special colouring.

**Definition 5.16** (Disintegration). Given a probability space  $(\mathbb{P}, \Omega, \mathcal{F})$ , random variables  $X$  and  $Y$  and joint probability measure  $\mathbb{P}_{XY} \in \Delta(\mathcal{E} \otimes \mathcal{F})$ , we say that  $\mathbf{M} : E \rightarrow \Delta(\mathcal{F})$  is a *Y on X disintegration* of  $\mu$  iff

$$\triangleleft \mu \begin{array}{c} X \\ Y \end{array} = \triangleleft \mu \begin{array}{c} X \\ * \mathbf{M} \\ Y \end{array} \quad (118)$$

$\mathbf{M}$  is a version of  $\mathbb{P}_{Y|X}$ , “the probability of  $Y$  given  $X$ ”. Let  $\mathbb{P}_{\{Y|X\}}$  be the set of all kernels that satisfy 118 and  $\mathbb{P}_{Y|X}$  an arbitrary member of  $\mathbb{P}_{Y|X}$ .

Given a Markov kernel space  $(\mathbb{K}, \Omega, D)$  and random variables  $X : \Omega \times D \rightarrow X$ ,  $Y : \Omega \times D \rightarrow Y$ ,  $\mathbf{M} : D \times E \rightarrow \Delta(\mathcal{F})$  is a *Y on DX disintegration* of  $\mathbb{K}_{YX|D}$  iff

$$\begin{array}{c} X \\ \mathbb{K}_{YX|D} \\ Y \end{array} = \begin{array}{c} X \\ \mathbb{K}_{YX|D} * \mathbf{M} \\ Y \end{array} \quad (119)$$

Write  $\mathbb{K}_{\{Y|XD\}}$  for the set of kernels satisfying 119 and  $\mathbb{K}_{Y|XD}$  for an arbitrary member of  $\mathbb{K}_{\{Y|XD\}}$ .

Note that for any variable  $X : \Omega \times D \rightarrow X$  and the domain variable  $D : \Omega \times D \rightarrow D$  we have by definition of  $\mathbb{K}$ :

$$\begin{array}{c} X \\ \mathbb{K}_{XD|D} \\ D \end{array} = \begin{array}{c} \mathbb{K}_0 \\ \mathbf{F}_X \\ \mathbf{F}_D \end{array} \begin{array}{c} X \\ D \end{array} \quad (120)$$

$$= \begin{array}{c} \mathbb{K}_0 \\ \mathbf{F}_X \end{array} \begin{array}{c} X \\ D \end{array} \quad (121)$$

$$= \begin{array}{c} \mathbb{K}_0 \\ \mathbf{F}_X \end{array} \begin{array}{c} X \\ D \end{array} \quad (122)$$

$$= \begin{array}{c} \mathbb{K} \\ \mathbf{F}_X \end{array} \begin{array}{c} X \\ D \end{array} \quad (123)$$

$$= \begin{array}{c} \mathbb{K}_{X|D} \\ \mathbf{F}_X \end{array} \begin{array}{c} X \\ D \end{array} \quad (124)$$

That is, any joint kernel including the variable  $D$  can be drawn such that the wire labeled  $D$  is copied from the input wire. Conversely, if we have a joint kernel  $\mathbb{K}_{X|D}$  and add a wire copied from the input, we now have  $\mathbb{K}_{XD|D}$ . We use this insight to give names to input wires: if copying an input wire to the output allows us to label the *output* wire with  $X$ , then the *input* wire will also be labeled  $X$ .

Warning: all work from here on out requires another pass of editing as of 17/08/2020

**Definition 5.17** (Wire labels - disintegrations). Given a conditional probability space with ambient kernel  $\mathcal{K} : D \rightarrow \Delta(\mathcal{F})$  (or a probability space with measure  $\mathbb{P}$ ),

Note that  $\mathbb{P}^*$  is simply  $\mathbb{P}$  for a probability space

Recall that  $D$  is the domain variable. Given two collections of random variables  $c_1 = [X_1, X_2, \dots]$  and  $c_2 = [Y_1, Y_2]$ , we adopt the convention that any diagram with the input wires labeled with  $c_1$  and the output wires labeled with  $c_2$  is an element of  $\mathcal{K}_{Y_1 Y_2 \dots | X_1 X_2 \dots}^*$ .

That is, by this convention, the diagram

$$\begin{array}{c} X \\ D \end{array} \text{---} \boxed{M} \text{---} Y \quad (125)$$

implies that  $M \in \mathcal{K}_{Y|XD}$ . Note further that by Theorem 5.19, we can rely on the existence of disintegrations such as  $M$  that are conditional on the global conditioning variable  $D$  provided we have countable  $D$  and standard measurable  $(Y, \mathcal{Y})$ .

If we have some version  $M$  of  $\mathcal{K}_{Y|XD}$  that does not depend on the value of  $D$  - i.e.  $M_{(x,d)} = M_{(x,d')}$  for all  $x \in X$ ,  $d, d' \in D$ , then there exists some  $M'$  such that:

$$\begin{array}{c} X \\ D \end{array} \text{---} \boxed{M} \text{---} Y = \begin{array}{c} X \\ D \end{array} \text{---} \boxed{M'} \text{---} Y \quad (126)$$

Under these circumstances, we will abuse notation to say  $M' = \mathcal{K}_{Y|X}$ .

We can't expect Equation 126 to hold in an arbitrary conditional probability spaces. For a very simple example, take  $\mathcal{K} : \{0, 1\} \rightarrow \Delta(\{0, 1\})$  where  $\mathcal{K}_0 = \mathcal{K}_1 = \text{Bernoulli}(0.5)$ , and let  $X : (x, d) \mapsto x$  - i.e. the random variable projecting the output of  $\mathcal{K}$ . Then there is no disintegration  $\mathcal{K}_{D|X}$  - we can't recover the input  $D$  from  $X$ .

Under some (strong) regularity conditions, disintegrations of conditional probability spaces do exist.

**Theorem 5.18** (Disintegration existence - probability space). *Given a probability measure  $\mu \in \Delta(\mathcal{E} \otimes \mathcal{F})$ , if  $(F, \mathcal{F})$  is standard then a disintegration  $K : E \rightarrow \Delta(\mathcal{F})$  exists (Çinlar, 2011).*

**Theorem 5.19** (Disintegration existence - conditional probability space). *Given a kernel  $L : D \rightarrow \Delta(\mathcal{E} \otimes \mathcal{F})$ , define  $L^*$ :*

$$\begin{array}{c} \text{---} \boxed{L} \text{---} \\ \text{---} \end{array} \quad (127)$$



If  $D$  is countable and  $(F, \mathcal{F})$  is standard, then there is a disintegration  $\mathbf{M} : D \times E \rightarrow \Delta(\mathcal{F})$  of  $\mathbf{L}^*$ .

*Proof.* By Theorem 5.18, for each  $d \in D$  we have a disintegration  $\mathbf{K}^{(d)} : E \rightarrow \Delta(\mathcal{F})$  of  $\mathbf{L}_d$ . Define  $\mathbf{M} : D \times E \rightarrow \Delta(\mathcal{F})$  by  $\mathbf{M}(d, e; A) = \mathbf{K}^{(d)}(e; A)$  for  $d \in D$ ,  $e \in E$ ,  $A \in \mathcal{F}$ . Clearly  $\mathbf{M}_{(d,e)}$  is a probability measure. Furthermore, for  $B \in \mathcal{B}(\mathbb{R})$ ,  $\mathbf{M}^{-1}(\cdot; A)(B) = \cup_{d \in D} \{d\} \times \mathbf{K}^{(d)-1}(\cdot; A)(B)$ , which is a countable union of measurable sets and therefore measurable.  $\square$

As an aside, Hájek (2003) pointed out that in general there are many Markov kernels that satisfy the definition of conditional probability for a given probability measure and random variables. While it is interesting that for a given Markov kernel space  $(\mathbb{K}, \Omega, \mathcal{F}, D, \mathcal{D})$  there is in general no probability measure on  $\Omega \times D$  such that  $\mathbb{K}$  is uniquely defined as a disintegration of  $\mu$ . By limiting  $D$  to countable sets, we avoid this possibility, and limit ourselves to Markov kernel spaces that can be uniquely defined as disintegrations.

From here on out, we will assume whether explicitly stated or not that any global conditioning space is countable and any other measurable space is standard, guaranteeing the existence of disintegrations.

In general, we don't want to spent time explicitly setting up conditional probability spaces. Rather, we will specify key marginals and disintegrations from which a conditional probability space can be constructed - call these marginals and conditional "components". Clearly we cannot build a conditional probability space from two kernels that represent the same component but disagree with each other on a non-negligible set. Also, in general, for an arbitrary collection of components there may be many ambient kernels from which we can extract these components. There is no particular problem if we have multiple ambient kernels over undefined random variable; if we are only interested in  $\mathbf{X}$  then the possibility of many joint kernels over  $\mathbf{X}$  and  $\mathbf{Y}$  is no cause for concern. We do, however, want to avoid ambient kernels supporting non-negligibly distinct marginals or disintegrations over the random variables that have been defined.

**Example 5.20** (Implicit conditional probability space). Suppose we have labeled Markov kernels

$$D \text{ -- } \boxed{\mathbf{L}} \text{ -- } \mathbf{X} \quad \mathbf{X} \text{ -- } \boxed{\mathbf{M}} \text{ -- } \mathbf{Y} \quad (128)$$

We want to define a conditional probability space  $(\mathcal{K}, \Omega, D)$  supporting random variables  $D$ ,  $\mathbf{X}$  and  $\mathbf{Y}$  yielding the above kernels as the relevant marginals and disintegrations. Strictly:

- $\mathbf{L} = \mathcal{K}_{\mathbf{X}|D}$
- $\mathbf{M} \otimes \mathbf{I}_D \in \mathcal{K}_{\mathbf{Y}|XD}$  ("informally",  $\mathbf{M} \in \mathcal{K}_{\mathbf{Y}|\mathbf{X}}$ )

Take  $\Omega = W \times X \times Y \times Z$  and define  $\mathcal{K}$  such that

$$\mathcal{K}^* = \begin{array}{c} \text{D} - \boxed{\text{L}} - \boxed{\text{M}} - \text{Y} \\ \quad \quad \quad \quad \quad \quad \text{X} \\ \quad \quad \quad \quad \quad \quad \text{D} \end{array} \quad (129)$$

Where  $\mathcal{K}^*$  is the copy map composed with  $\mathcal{K}$  as in previous definitions.  $\mathcal{K}$  is the unique Markov kernel  $D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$  supporting the two criteria above, assuming finite  $D$  and standard measurable  $X, Y$ .

*Proof.* By assumption, for any suitable  $\mathcal{K}' : D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$  we have

$$D - \boxed{\mathcal{K}'} -_X^* = D - \boxed{\mathbf{L}} - X \quad (130)$$

and by the fact that  $\mathbf{M} \otimes_D^*$  is by assumption a disintegration of  $\mathcal{K}^*$ :

$$\begin{array}{c} \text{D} \\ | \\ \boxed{\mathcal{K}'} \\ | \\ \begin{matrix} X \\ Y \\ D \end{matrix} \end{array} = \begin{array}{c} \text{D} \\ | \\ \boxed{\mathcal{K}'} \xrightarrow{*} \boxed{\mathbf{M}} \xrightarrow{*} \\ | \quad \quad \quad | \\ \begin{matrix} X \\ Y \\ D \end{matrix} \end{array} \quad (131)$$

$$= \begin{array}{c} D - \boxed{\text{L}} - \boxed{\text{M}} - Y \\ \quad \quad \quad \quad \quad \quad X \\ \quad \quad \quad \quad \quad \quad D \end{array} \quad (132)$$

Finally, if  $\mathcal{K}^* = \mathcal{K}'^*$ , then at least if  $D$  is countable we must have  $\mathcal{K} = \mathcal{K}'$  as they must agree on all points in  $D$ .  $\square$

This example was chosen to illustrate a peculiarity of our notation of conditional probability spaces. Consider a problem that appears similar: find an ambient measure  $\mathbb{P}$  decomposing into the following marginal and conditionals:

$$\triangleleft^{\mu} - \text{D} \quad \text{D} - \boxed{\text{L}} - \text{X} \quad \text{X} - \boxed{\text{M}} - \text{Y} \quad (133)$$

Here there are many choices of  $\mathbb{P}$  that satisfy our conditions arising from different choices of  $\mathbb{P}_{Y|X\mathcal{D}}$ . This is not possible in the conditional probability space because  $\mathcal{K}_{Y|X}$  only exists if  $\mathcal{K}_{Y|X\mathcal{D}}$  is independent of  $\mathcal{D}$ . That is, in a conditional probability space every disintegration is conditional on  $\mathcal{D}$ , but we may not explicitly write this if it does not actually depend on  $\mathcal{D}$ .

A sufficient condition for the construction of a unique ambient kernel from a collection of components  $\{C_1, \dots, C_n\}$  is if there is some ordering of components  $\{i_1, i_2, \dots, i_n\}$  such that the input labels of  $C_{i_{k+1}}$  is the union of the inputs and outputs of  $C_{i_1}, \dots, C_{i_j}$ . This can be shown by repeated application of Theorem 5.19.

In general, diagram labels are “well behaved” with regard to the application of any of the special Markov kernels: identities 79, swaps 81, discards 82 and copies 80 as well as with respect to the coherence theorem of the CD category. They are not “well behaved” with respect to composition.

**Lemma 5.21** (Diagrammatic consequences of labels). *Fix some conditional probability space  $(\mathcal{K}, \Omega, D)$  and random variables  $X, Y, Z$  taking values in arbitrary spaces.  $\text{Sat} :$  indicates that a labeled diagram satisfies definitions 5.15 and 5.17 with respect to  $(\mathcal{K}, \Omega, D)$  and  $X, Y, Z$ . The following always holds:*

$$\text{Sat} : X - X \quad (134)$$

and the following implications hold:

$$\text{Sat} : Z - \boxed{\mathbf{K}} - \begin{array}{c} X \\ \diagdown \\ Y \end{array} \implies \text{Sat} : Z - \boxed{\mathbf{K}} - * \begin{array}{c} X \\ \diagdown \\ Y \end{array} \quad (135)$$

$$\text{Sat} : Z - \boxed{\mathbf{K}} - \begin{array}{c} X \\ \diagdown \\ Y \end{array} \implies \text{Sat} : Z - \boxed{\mathbf{K}} - \begin{array}{c} X \\ \diagdown \\ Y \end{array} \quad (136)$$

$$\text{Sat} : Z - \boxed{\mathbf{L}} - X \implies \text{Sat} : Z - \boxed{\mathbf{L}} - \begin{array}{c} X \\ \diagdown \\ X \end{array} \quad (137)$$

$$\text{Sat} : Z - \boxed{\mathbf{K}} - Y \implies \text{Sat} : \begin{array}{c} Z \\ \diagdown \\ \boxed{\mathbf{K}} - Y \end{array} \quad (138)$$

*Proof.* •  $\text{Id}_X$  is a version of  $\mathbb{P}_{X|X}$  for all  $\mathbb{P}$ ;  $\mathbb{P}_X \text{Id}_X = \mathbb{P}_X$

- $\mathbf{K} \text{Id} \otimes * (w; A) = \int_{X \times Y} \delta_x(A) \mathbb{1}_Y(y) d\mathbf{K}_w(x, y) = \mathbf{K}_w(A \times Y) = \mathbb{P}_{X|Z}(w; A)$
  - $\int_{X \times Y} \delta_{\text{swap}(x,y)}(A \times B) d\mathbf{K}_w(x, y) = \mathbb{P}_{YX|Z}(w; A \times B)$
  - $\mathbf{K}^\vee(w; A \times B) = \int_X \delta_{x,x}(A \times B) d\mathbf{K}_w(x) = \mathbb{P}_{XX|Z}(w; A \times B)$
- 138: Suppose  $\mathbf{K}$  is a version of  $\mathbb{P}_{Y|Z}$ . Then

$$\mathbb{P}_{ZY} = \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \text{---} \boxed{\mathbf{K}} \text{---} \begin{array}{c} Z \\ \diagdown \\ Y \end{array} \end{array} \quad (139)$$

$$\mathbb{P}_{ZZY} = \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \text{---} \boxed{\mathbf{K}} \text{---} \begin{array}{c} Z \\ \diagdown \\ Z \\ \diagdown \\ Y \end{array} \end{array} \quad (140)$$

$$= \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \text{---} \boxed{\mathbf{K}} \text{---} \begin{array}{c} Z \\ \diagdown \\ Z \\ \diagdown \\ Y \end{array} \end{array} \quad (141)$$

Therefore  $\vee(\text{Id}_X \otimes \mathbf{K})$  is a version of  $\mathbb{P}_{ZY|Z}$  by ??  $\square$

The following property, on the other hand, does *not* generally hold:

$$\text{Sat} : Z \dashv \boxed{\mathbf{K}} \dashv Y, Y \dashv \boxed{\mathbf{L}} \dashv X \implies \text{Sat} : Z \dashv \boxed{\mathbf{K}} \dashv \boxed{\mathbf{L}} \dashv X \quad (142)$$

Consider some ambient measure  $\mathbb{P}$  with  $Z = X$  and  $\mathbb{P}_{Y|X} = x \mapsto \text{Bernoulli}(0.5)$  for all  $z \in Z$ . Then  $\mathbb{P}_{Z|Y} = y \mapsto \mathbb{P}_Z$ ,  $\forall y \in Y$  and therefore  $\mathbb{P}_{Y|Z}\mathbb{P}_{Z|Y} = x \mapsto \mathbb{P}_Z$  but  $\mathbb{P}_{Z|X} = x \mapsto \delta_x \neq \mathbb{P}_{Y|Z}\mathbb{P}_{Z|Y}$ .

This is an attempt to understand which assumptions lead to see-do models being represented by probability distributions.

## 6 Decision makers and decision models

We will consider a very general type of decision maker to motivate the construction of decision models. A decision maker takes some kind of evidence - for example, data, prior knowledge or a problem specification - and determines the consequences it believes each decision will have via some *inference rule*. It then examines the consequences of each decision and chooses a preferred decision via some *choice rule*.

More formally, let the set  $A$  be a set of possible pieces of evidence a decision maker could obtain,  $E$  be the set of consequences it may experience and  $D$  the set of decisions it may take. Call any function  $C : D \rightarrow E$  a *consequence map*. A *inference rule* is a function  $f : A \rightarrow E^D$  that takes a piece of evidence and returns a consequence map. A *choice rule*  $\text{ch} : E^D \rightarrow D$  is a function that takes a consequence map and returns a preferred decision. A decision maker is specified by a particular choice of  $f$  and  $\text{ch}$  and implements the *decision function*  $h := \text{ch} \circ f$ .

Suppose  $D = \{0, 1\}$  represents the decision to turn some switch to the on or off positions. Some examples of consequence maps follow:

- $E = \{0, 1\}$  is the state of a light connected to the switch and  $\mathbf{C} = \text{Id}_{\{0,1\}}$  is the consequence function that implies the light is always in the same state as the switch
- $E = \{\{0\}, \{1\}, \{0, 1\}\}$  represents sets of possible outcomes and  $\mathbf{C} : 0 \mapsto \{0\}, 1 \mapsto \{0, 1\}$  is the consequence function that implies the light is always off when the switch is off, but may take either state with the switch on (i.e. "the bulb may be out")

If  $A = \{0, 1\}^{2N}$  is a sequence of joint states of the switch and lightbulb, then an inference function could be the map

$$f(s_1, l_1, \dots, s_N, l_N) = \begin{cases} 0 \mapsto \llbracket \frac{\sum_{i \in N} l_i(1-s_i)}{\sum_{i \in N} (1-s_i)} > 0.5 \rrbracket \\ 1 \mapsto \llbracket \frac{\sum_{i \in N} l_i s_i}{\sum_{i \in N} s_i} > 0.5 \rrbracket \end{cases} \quad (143)$$

That is,  $f$  returns the consequence function that maps each switch position to the lightbulb state more commonly observed in conjunction with that switch position in the given data.

Some examples of decision rules, where in each case  $\mathbf{C} : D \rightarrow E$  is some consequence map:

- $E = \{\text{off}, \text{on}\}$  are states of the light bulb,  $u : \text{off} \mapsto 0, \text{on} \mapsto 1$  is a utility function and  $\text{ch}(\mathbf{C}) = \arg \min_{d \in D} u(\mathbf{C}(d))$
- $E = [0, 1]$  is a set of losses and  $\text{ch}(\mathbf{C}) = \arg \min_{d \in D} (\mathbf{C}(d))$
- $E = \{\{0\}, \{1\}, \{0, 1\}\}$  are sets of possible outcomes and  $\text{ch}$  is the minimax operator:  $\text{ch}(\mathbf{C}) = \arg \min_{d \in D} \max_{e \in \mathbf{C}(d)} (e)$

Another example of a decision maker is a learning algorithm that produces binary classifiers. A binary classifier takes some piece of data  $X$  and returns a class in  $\{0, 1\}$ . That is, the set of available decisions  $D$  is the set of functions from  $X \rightarrow \{0, 1\}$ . Suppose also that the given information is a set of  $N$  labeled examples  $A = (X \times \{0, 1\})^N$ . Then the inference function  $f$  might be an *empirical risk* functional that takes data  $a \in A$  and a possible classifier  $d \in D$  and returns the number of misclassified examples when  $d$  is applied to  $A$ . In such a case,  $\text{ch}$  might be the  $\arg \min$  functional, making our decision maker an *empirical risk minimiser*. This is not necessarily a good inference function - typically we are interested in the number of misclassified examples a classifier will produce on unseen data, not the number of misclassified examples it will produce on the data that has already been seen.

## 6.1 Expected utility maximisers

Expected utility maximisation (henceforth “EU”) necessitates that the class of results models  $f$  have a particular signature. The EU choice rule compares a set of probability distributions over states of the world  $E$  and returns a decision associated with the preferred probability distribution. Thus, denoting by  $\Delta(\mathcal{E})$  the set of probability distribution over  $E$  (now equipped with  $\sigma$ -algebra  $\mathcal{E}$ ),  $f$  must have the signature  $f : A \rightarrow \Delta(\mathcal{E})^D$ . In particular  $f$  returns stochastic functions of the type  $D \rightarrow \Delta(\mathcal{E})$ .

So far we have not made any structural assumptions about the set  $D$ . One such assumption that we do want to make is that stochastic decisions are possible. That is, if we can choose  $d_1 \in D$  and we can choose  $d_2 \in D$ , then it is also possible to choose  $d_1$  with probability  $\alpha$  and  $d_2$  with probability  $(1 - \alpha)$  for  $0 \leq \alpha \leq 1$ . We will make this assumption as follows:  $D$  represents an underlying set of “elementary” decisions, and the true set of decisions is the set  $\Delta(D)$  of probability measures on  $D$  (equipped with some  $\sigma$ -algebra  $\mathcal{D}$ ). Then, for example, if  $D = \{d_1, d_2\}$  we can “choose”  $d_1$  by selecting the measure  $\delta_{d_1}$ , and we can choose a mixture of  $d_1$  and  $d_2$  by selecting the measure  $\alpha\delta_{d_1} + (1 - \alpha)\delta_{d_2}$ . Technically, all that we have done so far is assume that the set of decisions is isomorphic to  $\Delta(D)$  - which, if  $D$  is standard Borel, ensures that the set of decisions is convex closed. In many cases, we also want the following to hold:

$$f_a(\alpha\delta_{d_1} + (1 - \alpha)\delta_{d_2}) = \alpha f_a(\delta_{d_1}) + (1 - \alpha)f_a(\delta_{d_2}) \quad (144)$$

That is, we assume that the results model  $f$  is *additive*. Informally, choosing to do  $d_1$  with probability  $\alpha$  and  $d_2$  with probability  $1 - \alpha$  will yield the result of  $d_1$  with probability  $\alpha$  and the result of  $d_2$  with probability  $1 - \alpha$ . In addition,  $f_a$  must be continuous with respect to the total variation norm.

For arbitrary  $\gamma \in \Delta(\mathcal{D})$ , we can write for all  $B \in \mathcal{D}$

$$\gamma(B) = \int_D \delta_d(B) d\gamma(d) \quad (145)$$

Can we also say

$$\gamma = \int_D \delta_d d\gamma(d) \quad (146)$$

For  $G \in \mathcal{E}$ , write  $f_a(\delta_d; C)$  for the probability measure  $f_a(\delta_d)$  evaluated at  $C$ , and  $f_a(\cdot; G)$  for the map  $d \mapsto f_a(d; C)$ . If  $f_a(\cdot; C)$  is measurable for all  $C$  and additive, then it is a property of the Lebesgue integral that for all  $\gamma \in \Delta(\mathcal{D})$

$$\int f_a(\gamma; C) = \int f_a\left(\int_D \delta_{d'} d\gamma(d'); C\right) \quad (147)$$

$$= \int_D f_a(\delta_{d'}; C) d\gamma(d') \quad (148)$$

Define  $\mathcal{C}_a : D \rightarrow \Delta(\mathcal{E})$  by  $\mathcal{C}_a : d \mapsto f_a(\delta_d; \cdot)$ . Then  $\mathcal{C}_a(\cdot; B)$  is measurable and  $\mathcal{C}_a(d; \cdot)$  is a probability measure - that is,  $\mathcal{C}_a$  is a Markov kernel. We can therefore adopt the product notation defined elsewhere to write the consequences of choosing a distribution  $\gamma$  as  $\gamma \mathcal{C}_a$ .

Have I made any assumptions here besides additivity, convex closedness of  $D$ ?

As we are chiefly interested in the set of elementary decisions  $D'$ , we will henceforth simply use  $(D, \mathcal{D})$  for this set.

Assumption 144 still permits some unexpected behaviour with respect to randomised decisions. Consider the following variations of a problem:

**Example 6.1** (Fighting children problem). A decision maker is mediating a dispute over a toy between two children. If he chooses to give the toy to the first child, the second child will cry, and if he chooses to give the toy to the second child, the first child will cry (the underlying decision set  $D' = \{\text{give to first child, give to second child}\}$ ). If he chooses a random procedure:

**Version 1** These children strongly object to randomness, and both will cry no matter the outcome

**Version 2** These are children with a strong insistence on a version of procedural fairness which no doubt make a great deal of sense to them. In particular, if a procedure is used which gives the first child a chance  $p$  of receiving the toy, the first child will cry with probability  $p$  irrespective of the outcome. The second child, meanwhile, will cry iff the first child does not

**Version 3** Whether or not the children cry depends only on whether or not they receive the toy

In version 1, assumption 144 is violated - the outcome when choosing a random procedure is not simply a mixture of the outcomes of each deterministic decision. In version 3, assumption 144 holds, and the results depend only on the outcome of randomisation, as is desired. In version 2, assumption 144 *also* holds - one child cries iff the other does not, and each child cries with the same probability as the probability that they receive the toy. However, the consequences of randomising and then giving the first child the toy are not the same as simply giving the first child the toy. To discount the possibility of this kind of behaviour, we need additional assumptions.

In the text of this example, we have spoken as if our decision had multiple consequences - whether or not each kid gets a toy, whether or not each kid cries - but the compliance of version 2 with 144 depends on considering only the children crying to be a consequence of our decision. To support the intuition that only version 3 is the kind of problem we want to deal with, we introduce the additional assumption that if we choose an extreme decision  $\delta_d$ , then one consequence of this decision is that we will have chosen  $d$  with probability 1.

Formally, we assume the sample space can be written as  $E \times D$  for some  $E$ , and, defining the random variable  $D : E \times D \rightarrow D$  by the projection  $D : (e, d) \mapsto d$ , suppose

$$(\mathcal{C}_{(a,d)})_D = \delta_d \quad (149)$$

We find in Joyce (2000) the assumption of a “supposition function”  $\text{prob}(\cdot||\cdot)$  that  $\text{prob}(C||C) = 1$  where  $C$  is “some distinguished condition”. Assumption 149 along with assumption 144 implies an analogous condition.

**Theorem 6.2** ([Decision determinism implies decision certainty]).

That is, for any stochastic decision  $\gamma$ , the consequence conditional on  $D$  is the consequence map  $f_a$  itself. This assumption rules out both versions 1 and 2 of the fighting children problem above - it assumes that if we randomise and the result of randomisation is some  $d \in D$ , this is the same as having chosen  $d$  to begin with.

**Lemma 6.3** (Conditional agreement implies randomised decisions). *If we have*

$$(\gamma f_a)_D = f_a \quad (150)$$

Then for all  $\gamma \in \Delta(\mathcal{D})$ ,

$$\gamma f_a(G) = \int_D f_a(d'; G) d\gamma(d') \quad (151)$$

*Proof.* By definition, for  $\gamma \in \Delta(\mathcal{D})$  □

A consequence of assumption ?? is that for fixed  $a \in A$ , every stochastic decision  $\gamma \in \Delta(\mathcal{D})$  leads to the same conditional probability  $(\gamma f_a)_{|D}$ . In fact, as a result of the assumption that the codomain of  $f$  is a set of probability measures as well as assumptions 149 and ?? we find that the tuple  $\langle E \times D, \mathcal{E} \otimes \mathcal{D}, \mathcal{D}, f_a \rangle$  for each  $a \in A$  forms a *conditional probability space* (Rényi, 1955).

For our causal analysis we work with *subjunctive probability spaces*, a generalisation of the more familiar probability spaces. The subjunctive mood is used to describe hypothetical or supposed states of the world, and subjunctive probability spaces are models that ask for some hypothetical state and give us back a probability space. Subjunctive probability spaces are also different to *conditional probability spaces* Rényi (1956) as *hypothesising* or *supposing* (that is, those things described by the subjunctive mood) are different to *conditioning*, which is more closely analogous to *focussing your attention*.

We use subjunctive probability spaces because *supposition* is a core part of decision problems, one we cannot get away from. In contrast, modelling supposition with conditional probability (which would be necessary if we were to insist on using probability spaces) adds additional structure to our models which isn't clearly warranted and is potentially confusing.

Any decision problem must involve the comparison of different decisions. This comparison takes the form

- *Suppose* I choose the first decision - then the consequence would be  $X$
- *Suppose* I choose the second decision - then the consequence would be  $Y$
- Etc.

If we prefer  $X$  to  $Y$ , then perhaps we should choose the first decision. We may be ambivalent as to what type of thing  $X$  and  $Y$  represent, how we ought to determine what values  $X$  and  $Y$  take or how we ought to determine what is preferable, but it is hard to do away with supposing that different decisions were taken and that these decisions come with their own consequences and still have what could reasonably be considered a decision problem. This process of supposition implicitly invokes a “subjunctive function” - we provide a hypothetical decision, and we are given a consequence of that hypothetical. Of particular interest are models that, for each hypothetical choice, return a probability distribution over space  $\Omega$ . Such models are called *supposition functions* by Joyce (2000), but we will call them *consequence maps* in this work. We consider this particular type of subjunctive function - from possible decisions to probability distributions - to be attractive because we consider probability to be a sound choice for modelling uncertainty and stochasticity.

Does this definition need to be here? I didn't know what this word meant before I looked it up



Formally, a consequence map  $\mathcal{C}$  is a stochastic function or *Markov kernel* from  $D \rightarrow \Delta(\Omega)$ , where  $\Delta(\Omega)$  represents the set of all probability distributions on  $\Omega$ . One might recall that, given two random variables  $Y : \Omega \rightarrow Y$  and  $X : \Omega \rightarrow X$  on some probability space  $(\mathbb{P}, \Omega, \mathcal{F})$ , the probability of  $Y$  conditional on  $X$  is a Markov kernel  $X \rightarrow \Delta(\mathcal{Y})$ . It is possible, in general, to define a probability distribution on the expanded space  $\Omega \times D$  such that  $\mathcal{C}$  is a conditional probability. However, there are good reasons to keep the concepts of consequence maps and conditional probability separate. Firstly, there are technical issues such as the fact that it is not always possible to find a distribution on  $\Omega \times D$  that yields  $\mathcal{C}$  as a *unique* conditional probability (Hájek, 2003) (though this requires an uncountable set of possible decisions, which is not a problem we consider in this work). Secondly, it is not clear that the way we ought to handle consequence maps is identical to the way we handle conditional probabilities. Consider an expanded set of options:

1. Suppose I choose the first decision  $d_1$  - then the consequence would be  $P_1$
2. Suppose I choose the second decision  $d_2$  - then the consequence would be  $P_2$
3. Suppose I choose either the first or second decision  $d_1$  or  $d_2$  - then the consequence would be ???
4. Suppose I choose  $d_1$  with probability  $0 \leq q \leq 1$  and  $d_2$  otherwise - then the consequence would be ???

If we regard the consequence map as a conditional probability then it would be natural to consider the result of the third option to be a unique probability distribution obtained by conditioning on  $\{d_1, d_2\}$ , equal to  $\alpha P_1 + (1 - \alpha) P_2$  for some fixed  $0 \leq \alpha \leq 1$ . However, there is no obvious reason that these should be mixed in some fixed proportion  $\alpha$  - it seems more appropriate to me to say that if  $d_1$  or  $d_2$  is chosen then the result will be  $P_1$  or  $P_2$ .

On the other hand, the result of the fourth option seems like it should be given by  $q P_1 + (1 - q) P_2$ , at least for most ordinary problems. If we were confronted with a mind reader who could tell the difference between us having chosen  $d_1$  and us having randomised between  $d_1$  and  $d_2$  but come up with  $d_1$  anyway then we might have reason to revise this assumption, but we will generally proceed under the assumption that such mind readers are absent. For any ambient probability measure, we can choose  $q$  not to be equal to  $\alpha$  (as defined in the previous paragraph), and so the result will not be equal to the ambient measure conditioned on  $\{d_1, d_2\}$ .

In general, choosing the consequence map to be a conditional probability requires there to exist some joint distribution over the decision  $D$  and all other random variables in the problem. However, when we adopt hypotheses about what decisions we might make, we throw away key parts of this joint distribution (for example, we throw away the marginal distribution of  $D$ ) and replace it with whatever we want to suppose instead. Instead of adopting a full joint

distribution and then throwing away some parts to get hold of the consequence map, as we would if we were working with a probability space, a subjunctive probability space only supplies those parts which we intend to keep - i.e. in only supplies the consequence map.

## References

- Stephan Bongers, Jonas Peters, Bernhard Schölkopf, and Joris M. Mooij. Theoretical Aspects of Cyclic Structural Causal Models. *arXiv:1611.06221 [cs, stat]*, November 2016. URL <http://arxiv.org/abs/1611.06221>. arXiv: 1611.06221.
- Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- Kenta Cho and Bart Jacobs. Disintegration and Bayesian inversion via string diagrams. *Mathematical Structures in Computer Science*, 29(7):938–971, August 2019. ISSN 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488.
- Florence Clerc, Fredrik Dahlqvist, Vincent Danos, and Ilias Garner. Pointless learning. *20th International Conference on Foundations of Software Science and Computation Structures (FoSSaCS 2017)*, March 2017. doi: 10.1007/978-3-662-54458-7\_21. URL [https://www.research.ed.ac.uk/portal/en/publications/pointless-learning\(694fb610-69c5-469c-9793-825df4f8ddec\).html](https://www.research.ed.ac.uk/portal/en/publications/pointless-learning(694fb610-69c5-469c-9793-825df4f8ddec).html).
- A. P. Dawid. Influence Diagrams for Causal Modelling and Inference. *International Statistical Review*, 70(2):161–189, 2002. ISSN 1751-5823. doi: 10.1111/j.1751-5823.2002.tb00354.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1751-5823.2002.tb00354.x>.
- A. Philip Dawid. Decision-theoretic foundations for statistical causality. *arXiv:2004.12493 [math, stat]*, April 2020. URL <http://arxiv.org/abs/2004.12493>. arXiv: 2004.12493.
- Angus Deaton and Nancy Cartwright. Understanding and misunderstanding randomized controlled trials. *Social Science & Medicine*, 210:2–21, August 2018. ISSN 0277-9536. doi: 10.1016/j.socscimed.2017.12.005. URL <http://www.sciencedirect.com/science/article/pii/S0277953617307359>.
- Ronald A. Fisher. Cancer and Smoking. *Nature*, 182(4635):596–596, August 1958. ISSN 1476-4687. doi: 10.1038/182596a0. URL <https://www.nature.com/articles/182596a0>. Number: 4635 Publisher: Nature Publishing Group.
- Brendan Fong. Causal Theories: A Categorical Perspective on Bayesian Networks. *arXiv:1301.6201 [math]*, January 2013. URL <http://arxiv.org/abs/1301.6201>. arXiv: 1301.6201.

- D. Heckerman and R. Shachter. Decision-Theoretic Foundations for Causal Reasoning. *Journal of Artificial Intelligence Research*, 3:405–430, December 1995. ISSN 1076-9757. doi: 10.1613/jair.202. URL <https://www.jair.org/index.php/jair/article/view/10151>.
- James J. Heckman. Randomization and Social Policy Evaluation. SSRN Scholarly Paper ID 995151, Social Science Research Network, Rochester, NY, July 1991. URL <https://papers.ssrn.com/abstract=995151>.
- M. A. Hernán and S. L. Taubman. Does obesity shorten life? The importance of well-defined interventions to answer causal questions. *International Journal of Obesity*, 32(S3):S8–S14, August 2008. ISSN 1476-5497. doi: 10.1038/ijo.2008.82. URL <https://www.nature.com/articles/ijo200882>.
- Alan Hájek. What Conditional Probability Could Not Be. *Synthese*, 137(3): 273–323, December 2003. ISSN 1573-0964. doi: 10.1023/B:SYNT.0000004904.91112.16. URL <https://doi.org/10.1023/B:SYNT.0000004904.91112.16>.
- Alan Hájek. Interpretations of Probability. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, fall 2019 edition, 2019. URL <https://plato.stanford.edu/archives/fall2019/entries/probability-interpret/>.
- James M. Joyce. Why We Still Need the Logic of Decision. *Philosophy of Science*, 67:S1–S13, 2000. ISSN 0031-8248. URL [www.jstor.org/stable/188653](http://www.jstor.org/stable/188653).
- Chayakrit Krittanawong, Bharat Narasimhan, Zhen Wang, Joshua Hahn, Hafeez Ul Hassan Virk, Ann M. Farrell, HongJu Zhang, and WH Wilson Tang. Association between chocolate consumption and risk of coronary artery disease: a systematic review and meta-analysis. *European Journal of Preventive Cardiology*, July 2020. doi: 10.1177/2047487320936787. URL <http://journals.sagepub.com/doi/10.1177/2047487320936787>. Publisher: SAGE PublicationsSage UK: London, England.
- Finnian Lattimore and David Rohde. Replacing the do-calculus with Bayes rule. *arXiv:1906.07125 [cs, stat]*, December 2019. URL <http://arxiv.org/abs/1906.07125>. arXiv: 1906.07125.
- David K Lewis. Causation. *Journal of Philosophy*, 1986.
- Stephen L. Morgan and Christopher Winship. *Counterfactuals and Causal Inference: Methods and Principles for Social Research*. Cambridge University Press, New York, NY, 2 edition edition, November 2014. ISBN 978-1-107-69416-3.
- Naomi Oreskes and Erik M. Conway. *Merchants of Doubt: How a Handful of Scientists Obscured the Truth on Issues from Tobacco Smoke to Climate Change: How a Handful of Scientists ... Issues from Tobacco Smoke to Global Warming*. Bloomsbury Press, New York, NY, June 2011. ISBN 978-1-60819-394-3.

- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2 edition, 2009.
- Judea Pearl. Challenging the hegemony of randomized controlled trials: A commentary on Deaton and Cartwright. *Social Science & Medicine*, 2018a.
- Judea Pearl. Does Obesity Shorten Life? Or is it the Soda? On Non-manipulable Causes. *Journal of Causal Inference*, 6(2), 2018b. 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>.
- Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, New York, 1 edition edition, May 2018. ISBN 978-0-465-09760-9.
- Jonas Peters, Dominik Janzing, and Bernard Schölkopf. *Elements of Causal Inference*. MIT Press, 2017.
- Robert N. Proctor. The history of the discovery of the cigarettelung cancer link: evidentiary traditions, corporate denial, global toll. *Tobacco Control*, 21(2):87–91, March 2012. ISSN 0964-4563, 1468-3318. doi: 10.1136/tobaccocontrol-2011-050338. URL <https://tobaccocontrol.bmj.com/content/21/2/87>. Publisher: BMJ Publishing Group Ltd Section: The shameful past.
- Thomas S Richardson and James M Robins. Single world intervention graphs (swigs): A unification of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social Sciences, University of Washington Series. Working Paper*, 128(30):2013, 2013.
- Donald B. Rubin. Causal Inference Using Potential Outcomes. *Journal of the American Statistical Association*, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880. URL <https://doi.org/10.1198/016214504000001880>.
- A. Rényi. On Conditional Probability Spaces Generated by a Dimensionally Ordered Set of Measures. *Theory of Probability & Its Applications*, 1(1):55–64, January 1956. ISSN 0040-585X. doi: 10.1137/1101005. URL <https://epubs.siam.org/doi/abs/10.1137/1101005>.
- Alfréd Rényi. On a new axiomatic theory of probability. *Acta Mathematica Academiae Scientiarum Hungarica*, 6(3):285–335, September 1955. ISSN 1588-2632. doi: 10.1007/BF02024393. URL <https://doi.org/10.1007/BF02024393>.
- 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.

- Ilya Shpitser and Judea Pearl. Complete Identification Methods for the Causal Hierarchy. *Journal of Machine Learning Research*, 9(Sep):1941–1979, 2008. ISSN 1533-7928. URL <https://www.jmlr.org/papers/v9/shpitser08a.html>.
- Statista. Cigarettes - worldwide | Statista Market Forecast, 2020. URL <https://www.statista.com/outlook/50010000/100/cigarettes/worldwide>.
- Robert Wiblin. Why smoking in the developing world is an enormous problem and how you can help save lives, 2016. URL <https://80000hours.org/problem-profiles/tobacco/>.
- James Woodward. Causation and Manipulability. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2016 edition, 2016. URL <https://plato.stanford.edu/archives/win2016/entries/causation-mani/>.
- World Health Organisation. Tobacco Fact sheet no 339, 2018. URL <https://www.webcitation.org/6gUXrCDKA>.
- Karren Yang, Abigail Katoff, and Caroline Uhler. Characterizing and Learning Equivalence Classes of Causal DAGs under Interventions. In *International Conference on Machine Learning*, pages 5537–5546, July 2018. URL <http://proceedings.mlr.press/v80/yang18a.html>.

## Appendix: