

Causal Statistical Decision Theory|What are interventions?

David Johnston

April 22, 2022

Contents

1	Introduction	5
1.1	Theories of causal inference	5
2	Technical Prerequisites	9
2.1	Conventions	10
2.2	Probability Theory	10
2.2.1	Standard Probability Theory	10
2.3	String Diagrams	15
2.3.1	Elements of string diagrams	15
2.3.2	Special maps	17
2.3.3	Commutative comonoid axioms	18
2.3.4	Manipulating String Diagrams	18
2.4	Probability Sets	20
2.4.1	Almost sure equality	22
2.4.2	Extended conditional independence	23
2.4.3	Examples	25
2.4.4	Maximal probability sets and valid conditionals	27
2.4.5	Existence of conditional probabilities	29
3	Models with choices and consequences	35
3.1	What is the point of causal inference?	36
3.1.1	Modelling decision problems	37
3.1.2	Formal definitions	38
3.2	Representation theorems for decision problems	39
3.2.1	von Neumann-Morgenstern utility	40
3.2.2	Savage decision theory	40
3.2.3	Jeffrey's decision theory	43
3.2.4	Causal decision theory	44
3.2.5	Statistical decision theory	45
3.3	Variables	49
3.3.1	Variables and measurement procedures	51
3.3.2	Measurement procedures	51
3.3.3	Observable variables	53
3.3.4	Model variables	54

3.3.5	Variable sequences and partial order	54
3.3.6	Decision procedures	54
4	Repeatable Response Functions	57
4.1	When do response functions exist?	57
4.1.1	Relevance to previous work	58
4.1.2	Causal contractibility	59
4.1.3	Existence of response conditionals	61
4.1.4	Elaborations and examples	63
4.1.5	Assessing causal contractibility	64
4.1.6	Body mass index revisited	68
4.2	Allowing dependence on observations	69
4.2.1	Combs	69
4.2.2	Response conditionals in models with history dependence	71
4.2.3	Validity	72
4.2.4	Combs are the output of the “fix” operation	72
4.3	Weaker assumptions than causal contractibility	73
5	Statistical Decision Theory	81
5.1	Summary	81
5.2	Modelling observations, choices and consequences	82
5.2.1	Modelling observations with statistical models	82
5.2.2	Modelling choices and consequences with two-player sta- tistical models	84
6	See-do models, interventions and counterfactuals	87
6.1	How do see-do models relate to other approaches to causal infer- ence?	87
6.2	Interpretations of the choice set	87
6.3	Causal Bayesian Networks as see-do models	87
6.4	Unit Potential Outcomes models	88
6.4.1	D-causation	91
6.4.2	D-causation vs Limited Unresponsiveness	92
6.4.3	Properties of D-causation	95
6.4.4	Decision sequences and parallel decisions	95
6.5	Existence of counterfactuals	96

Chapter 1

Introduction

1.1 Theories of causal inference

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 to many different people. 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. Additionally, some may desire reassurance that their habit is not too risky, whether or not this is true. Potential and actual investors in cigarette manufacturers may see health concerns as a barrier to adoption, and also may personally want to avoid supporting products that harm many people. Like smokers, such people might have some interest in knowing the truth of this question, and a separate interest in hearing that cigarettes are not too risky, whether or not this is true. 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 billion in 2020 (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 causal questions is of tremendous importance to huge numbers of people. Second, confusion over causal questions can persist even when a great deal of data and facts relevant to the question are agreed upon.

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

[...] 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 (which is suggested by the word “benefit”), 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.

For causal questions that are controversial or difficult, it is tremendously advantageous to be able to address them transparently. Theories of causation enable this; given a theory of causation and a set of assumptions, if anyone claims that some conclusion follows it is publicly verifiable whether or not it actually does so. If the deduction is correct, then any remaining disagreement must be in the assumptions or in the theory. For people who are interested in understanding what is true, pinpointing disagreement can be enlightening. Someone could learn, for example, that there are assumptions that they find plausible that permit conclusions they did not initially believe. Alternatively, if a motivated conclusion follows only from implausible assumptions, hearing these assumptions explicitly might make the conclusion less attractive.

Theories of causation help us to answer causal questions, which means that before we have any theory, we already have causal questions we want to answer. If potential outcomes notation and causal graphical models had never been invented there would still be just as many people who want to the answer to questions something like “does smoking causes cancer?”, even if on-one could

say what exactly they meant by “causes” and even if many people actually want answers to slightly different questions. Theories exist to serve our need for transparent answers to causal questions.

Potential outcomes and causal graphical models are prominent examples of “practical theories” of causation. I call them “practical theories” because most of the time we encounter them they are being used to answer “practical” questions like “Does smoking cause cancer?”, or “In general, when does data allow us to conclude that X causes Y ?” It is less common to see the “fundamental questions” addressed, like “Does the theory of causal graphical models offer an adequate account of what ‘cause’ means?”, which is more often found in the field of philosophy. Spirtes et al. (2000) explain their motivation to study what I call “practical theories of causation” as follows:

One approach to clarifying the notion of causation – the philosophers approach ever since Plato – is to try to define “causation” in other terms, to provide necessary and sufficient and noncircular conditions for one thing, or feature or event or circumstance, to cause another, the way one can define “bachelor” as “unmarried adult male human.” Another approach to the same problem – the mathematicians approach ever since Euclid – is to provide axioms that use the notion of causation without defining it, and to investigate the necessary consequences of those assumptions. We have few fruitful examples of the first sort of clarification, but many of the second [...]

I think what Spirtes, Glymour and Scheines (henceforth: SGS) mean here is that they *define* a notion of causation – because causal graphical models do define a notion of causation – without interrogating whether it means the same thing as the word “causation”. Incidentally, since publication of this paragraph, the notion of causation defined by causal graphical models has been subject to substantial interrogation by philosophers (Woodward, 2016).

I am sympathetic to the argument that it does not matter a great deal whether “causal-graphical-models-causation” and “causation” mean the same thing in everyday language. It is common for words to have somewhat different meanings when used by specialists to when they are used by laypeople, and this isn’t because the specialists are ignorant or confused about their subject. However, I think it matters a lot which causal questions can be transparently answered by “causal-graphical-models-causation”, and so I believe that the notions of causation adopted by practical theories do warrant scrutiny.

I think one reason that SGS are keen to avoid dwelling on the definition of causation is that satisfactory definitions of causation are difficult. For example, causal graphical models depend on the notion of *causal relationships* between variables. These may be defined as follows:

X_i is a *cause* of X_j if there is an *ideal intervention* on X_i that changes the value X_j

This definition is incomplete without a definition of “ideal interventions”. Ideal interventions may be defined by their action in “causally sufficient models”:

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

While I don’t offer a definition of the *do-calculus* in this introduction, it can be rigorously defined, see for example Pearl (2009). The problem is that the definition of a *causally sufficient* model itself invokes the word *cause*, which is what the original definition was trying to address. Circularity is a recognised problem with interventional definitions of causation (Woodward, 2016). In Section ??, I further show models with ideal interventions generally have counterintuitive properties. The purpose of a theory of causation like causal graphical models is to support transparent reasoning about causal questions, and a circular definition that leads to counterintuitive conclusions undermines this purpose.

As with Euclid’s parallel postulate, I think it is reasonable to ask if the notion of ideal interventions and other causal definitions can be modified or avoided. Causal statistical decision theory (CSDT) is a theory of causation that is motivated by the problem of *what is generally needed to answer causal questions* rather than *what does “causation” mean?* Along similar lines to CSDT, Dawid (2020) has observed that the problem of deciding how to act in light of data can be formalised without appeal to theories of causation. We develop this in substantial detail, showing how both *interventional models* and *counterfactual models* arise as special cases of CSDT.

A key feature of CSDT is what I call the *option set*. This is the set of decisions, acts or counterfactual propositions under consideration in a given problem. A causal graphical model and a potential outcomes model will both implicitly define an option set as a result of their basic definitions of causation, but CSDT demands that this is done explicitly. I argue that this is a key strength of CSDT, on the basis of the following claims which I defend in the following chapters:

- Causal questions are not well-posed without an option set in the same way a function is not well-defined without its domain
- The option set need not correspond in any fixed manner to the set of observed variables
- The nature of the option set can affect the difficulty of causal inference questions

I commented out an additional section about potential outcomes and closest world counterfactuals, which is a second example of “opaque causal definitions”. I’m interested if any readers think it would be good to have a second example

¹Weaker conditions for causal sufficiency are possible, but they don’t avoid circularity (Shpitser and Pearl, 2008)

I want to revisit the claims about what I actually show, hopefully to add to it

Chapter 2

Technical Prerequisites

Our approach to causal inference is (like most other approaches) based on probability theory. Many results and conventions will be familiar to readers, and these are collected in Section 2.2.1.

Less likely to be familiar to readers is the string diagram notation we use to represent probabilistic functions. This is a notation created for reasoning about abstract Markov categories, and is somewhat different to existing graphical languages. The main difference is that in our notation wires represent variables and boxes (which are like nodes in directed acyclic graphs) represent probabilistic functions. Standard directed acyclic graphs annotate nodes with variable names and represent probabilistic functions implicitly. The advantage of explicitly representing probabilistic functions is that we can write equations involving graphics. It is introduced in Section 2.3.

We also extend the theory of probability to a theory of probability sets, which we introduce in Section 3.1.2. This section goes over some ground already trodden by Section 2.2.1; this structure was chosen so that people familiar with the Section 2.2.1 can skip to Section 3.1.2 for relevant generalisations to probability sets. Two key ideas introduced here are *uniform conditional probability*, similar but not identical to conditional probability, and *extended conditional independence* as introduced by Constantinou and Dawid (2017), similar but not identical to regular conditional independence.

We finally introduce the assumption of *validity*, which ensures that probability sets constructed by “assembling” collections of uniform conditionals are non-empty.

This is a reference chapter – a reader who is already quite familiar with probability theory may skip to Chapter 3. Where necessary, references back to theorems and definitions in this chapter are given. In Chapter 4, we will introduce one additional probabilistic primitive: *combs*, as we feel that additional context is helpful for understanding them.

2.1 Conventions

One of the unusual conventions in this thesis is the notation of uniform conditional probability. Given a set of probability distributions $\mathbb{P}_C := \{\mathbb{P}_\alpha | \alpha \in C\}$ on a common sample space (Ω, \mathcal{F}) with variables $X : \Omega \rightarrow X$ and $Y : \Omega \rightarrow Y$, $\mathbb{P}_C^{Y|X}$ represents a Markov kernel $X \rightarrow Y$ that satisfies the definition of the distribution of Y given X (Definition 2.2.16) for every $\alpha \in C$, while $\mathbb{P}_\alpha^{Y|X}$ is a conditional distribution only for α . There are two unusual features: firstly, it is more common to write a conditional distribution $\mathbb{P}(Y|X)$ and secondly, the subscript indicating the “domain of validity” of the conditional probability is unusual.

Because this thesis uses sets of probability measures rather than single probability measures, in general a conditional distribution may be valid only for some subset of the probability measures, and always including a subscript indicating which subset or element for which a conditional distribution is valid avoids any ambiguity about this. Avoiding notation of the form $\mathbb{P}(Y|X)$ is an aesthetic preference; writing a conditional distribution like this suggests $\mathbb{P}(Y|X)$ is the result of function composition between \mathbb{P} and some function denoted “ $Y|X$ ”. However, conditional probabilities are not given by composition of functions like this.

Name	notation	meaning
Iverson bracket	$\llbracket \cdot \rrbracket$	Function equal to 1 if \cdot is true, false otherwise
Identity function	idf_X	Identity function $X \rightarrow X$
Identity kernel	id_X	Kernel associated with the identity function $X \rightarrow X$

2.2 Probability Theory

2.2.1 Standard Probability Theory

σ -algebras

Definition 2.2.1 (Sigma algebra). Given a set A , a σ -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} , $\cup_{i \in \mathbb{N}} B_i \in \mathcal{A}$

Definition 2.2.2 (Measurable space). A measurable space (A, \mathcal{A}) is a set A along with a σ -algebra \mathcal{A} .

Definition 2.2.3 (Sigma algebra generated by a set). Given a set A and an arbitrary collection of subsets $U \subset \mathcal{P}(A)$, the σ -algebra generated by U , $\sigma(U)$, is the smallest σ -algebra containing U .

Common σ algebras For any A , $\{\emptyset, A\}$ is a σ -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 σ -algebra.

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

If A is a topological space with open sets T , $\mathcal{B}(\mathbb{R}) := \sigma(T)$ is the *Borel σ -algebra* on A .

If A is a separable, complete topological space, then $(A, \mathcal{B}(A))$ is a *standard measurable set*. All standard measurable sets are isomorphic to either $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ or $(C, \mathcal{P}(C))$ for denumerable C (Çinlar, 2011, Chap. 1).

Probability measures and Markov kernels

Definition 2.2.4 (Probability measure). 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 $\{A_i\} \subset \mathcal{E}$, $\mu(\cup_i A_i) = \sum_i \mu(A_i)$

Nxample 2.2.5 (Set of all probability measures). The set of all probability measures on (E, \mathcal{E}) is written $\Delta(E)$.

Definition 2.2.6 (Markov kernel). Given measurable spaces (E, \mathcal{E}) and (F, \mathcal{F}) , a *Markov kernel* or *stochastic function* is a map $\mathbb{M} : E \times \mathcal{F} \rightarrow [0, 1]$ such that

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

Nxample 2.2.7 (Signature of a Markov kernel). Given measurable spaces (E, \mathcal{E}) and (F, \mathcal{F}) and $\mathbb{M} : E \times \mathcal{F} \rightarrow [0, 1]$, we write the signature of $\mathbb{M} : E \rightarrow F$, read “ \mathbb{M} maps from E to probability measures on F ”.

Definition 2.2.8 (Deterministic Markov kernel). A *deterministic* Markov kernel $\mathbb{A} : E \rightarrow \Delta(\mathcal{F})$ is a kernel such that $\mathbb{A}_x(B) \in \{0, 1\}$ for all $x \in E, B \in \mathcal{F}$.

Common probability measures and Markov kernels

Definition 2.2.9 (Dirac measure). The *Dirac measure* $\delta_x \in \Delta(X)$ is a probability measure such that $\delta_x(A) = \mathbb{I}[x \in A]$

Definition 2.2.10 (Markov kernel associated with a function). Given measurable $f : (X, \mathcal{X}) \rightarrow (Y, \mathcal{Y})$, $\mathbb{F}_f : X \rightarrow Y$ is the Markov kernel given by $x \mapsto \delta_{f(x)}$

Definition 2.2.11 (Markov kernel associated with a probability measure). Given (X, \mathcal{X}) , a one-element measurable space $(\{*\}, \{\{*\}, \emptyset\})$ and a probability measure $\mu \in \Delta(X)$, the associated Markov kernel $\mathbb{Q}_\mu : \{*\} \rightarrow X$ is the unique Markov kernel $* \mapsto \mu$

Lemma 2.2.12 (Products of functional kernels yield function composition). *Given measurable $f : X \rightarrow Y$ and $g : Y \rightarrow Z$, $\mathbb{F}_f \mathbb{F}_g = \mathbb{F}_{g \circ f}$.*

Proof.

$$(\mathbb{F}_f \mathbb{F}_g)_x(A) = \int_X (\mathbb{F}_g)_y(A) d(\mathbb{F}_f)_x(y) \quad (2.1)$$

$$= \int_X \delta_{g(y)}(A) d\delta_{f(x)}(y) \quad (2.2)$$

$$= \delta_{g(f(x))}(A) \quad (2.3)$$

$$= (\mathbb{F}_{g \circ f})_x(A) \quad (2.4)$$

□

Variables, conditionals and marginals

Definition 2.2.13 (Variable). Given a measurable space (Ω, \mathcal{F}) and a measurable space of values (X, \mathcal{X}) , an *X-valued variable* is a measurable function $X : (\Omega, \mathcal{F}) \rightarrow (X, \mathcal{X})$.

Definition 2.2.14 (Sequence of variables). Given a measurable space (Ω, \mathcal{F}) and two variables $X : (\Omega, \mathcal{F}) \rightarrow (X, \mathcal{X})$, $Y : (\Omega, \mathcal{F}) \rightarrow (Y, \mathcal{Y})$, $(X, Y) : \Omega \rightarrow X \times Y$ is the variable $\omega \mapsto (X(\omega), Y(\omega))$.

Definition 2.2.15 (Marginal distribution). Given a probability space $(\mu, \Omega, \mathcal{F})$ and a variable $X : \Omega \rightarrow (X, \mathcal{X})$, the *marginal distribution* of X with respect to μ , $\mu^X : \mathcal{X} \rightarrow [0, 1]$ by $\mu^X(A) := \mu(X^{-1}(A))$ for any $A \in \mathcal{X}$.

Definition 2.2.16 (Conditional distribution). Given a probability space $(\mu, \Omega, \mathcal{F})$ and variables $X : \Omega \rightarrow X$, $Y : \Omega \rightarrow Y$, the *conditional distribution* of Y given X is any Markov kernel $\mu^{Y|X} : X \rightarrow Y$ such that

$$\mu^{XY}(A \times B) = \int_A \mu^{Y|X}(B|x) d\mu^X(x) \quad \forall A \in \mathcal{X}, B \in \mathcal{Y} \quad (2.5)$$

$$\iff \quad (2.6)$$

$$\mu^{XY} = \begin{array}{c} \text{---} X \\ \text{---} \mu^X \text{---} \bullet \text{---} \mu^{Y|X} \text{---} Y \end{array} \quad (2.7)$$

Markov kernel product notation

Three pairwise *product* operations involving Markov kernels can be defined: measure-kernel products, kernel-kernel products and kernel-function products. These are analagous to row vector-matrix products, matrix-matrix products and matrix-column vector products respectively.

, $\mathbb{T} : Y \rightarrow T$, $\mathbb{M} : X \rightarrow \Delta(\mathcal{Y})$ and $\mathbb{N} : Y \rightarrow \Delta(\mathcal{Z})$

Definition 2.2.17 (Measure-kernel product). Given $\mu \in \Delta(\mathcal{X})$ and $\mathbb{M} : X \rightarrow Y$, the *measure-kernel product* $\mu\mathbb{M} \in \Delta(Y)$ is given by

$$\mu\mathbb{M}(A) := \int_X \mathbb{M}(A|x)\mu(dx) \quad (2.8)$$

for all $A \in \mathcal{Y}$.

Definition 2.2.18 (Kernel-kernel product). Given $\mathbb{M} : X \rightarrow Y$ and $\mathbb{N} : Y \rightarrow Z$, the *kernel-kernel product* $\mathbb{M}\mathbb{N} : X \rightarrow Z$ is given by

$$\mathbb{M}\mathbb{N}(A|x) := \int_Y \mathbb{N}(A|y)\mathbb{M}(dy|x) \quad (2.9)$$

for all $A \in \mathcal{Z}$, $x \in X$.

Definition 2.2.19 (Kernel-function product). Given $\mathbb{M} : X \rightarrow Y$ and $f : Y \rightarrow Z$, the *kernel-function product* $\mathbb{M}f : X \rightarrow Z$ is given by

$$\mathbb{M}f(x) := \int_Y f(y)\mathbb{N}(dy|x) \quad (2.10)$$

for all $x \in X$.

Definition 2.2.20 (Tensor product). Given $\mathbb{M} : X \rightarrow Y$ and $\mathbb{L} : W \rightarrow Z$, the tensor product $\mathbb{M} \otimes \mathbb{L} : X \times W \rightarrow Y \times Z$ is given by

$$(\mathbb{M} \otimes \mathbb{L})(A \times B|x, w) := \mathbb{M}(A|x)\mathbb{L}(B|w) \quad (2.11)$$

For all $x \in X$, $w \in W$, $A \in \mathcal{Y}$ and $B \in \mathcal{Z}$.

All products are associative (Çinlar, 2011, Chapter 1).

One application of the product notation is that marginal distributions can be alternatively defined in terms of a kernel product, as shown in Lemma 2.2.21.

Lemma 2.2.21 (Marginal distribution as a kernel product). *Given a probability space $(\mu, \Omega, \mathcal{F})$ and a variable $\mathbf{X} : \Omega \rightarrow (X, \mathcal{X})$, define $\mathbb{F}_{\mathbf{X}} : \Omega \rightarrow X$ by $\mathbb{F}_{\mathbf{X}}(A|\omega) = \delta_{\mathbf{X}(\omega)}(A)$, then*

$$\mu^{\mathbf{X}} = \mu\mathbb{F}_{\mathbf{X}} \quad (2.12)$$

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

$$\mu\mathbb{F}_{\mathbf{X}}(A) = \int_{\Omega} \delta_{\mathbf{X}(\omega)}(A) d\mu(\omega) \quad (2.13)$$

$$= \int_{\mathbf{X}^{-1}(A)} d\mu(\omega) \quad (2.14)$$

$$= \mu^{\mathbf{X}}(A) \quad (2.15)$$

□

Semidirect product

Given a marginal μ^X and a conditional $\mu^{Y|X}$, the product of the two yields the marginal distribution of Y : $\mu^Y = \mu^X \mu^{Y|X}$. We define another product – the *semidirect* product \odot – as the product that yields the joint distribution of (X, Y) : $\mu^{XY} = \mu^X \odot \mu^{Y|X}$. The semidirect product is associative (Lemma 2.2.23)

Definition 2.2.22 (Semidirect product). Given $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : Y \times X \rightarrow Z$, the semidirect product $\mathbb{K} \odot \mathbb{L} : X \rightarrow Y \times Z$ is given by

$$(\mathbb{K} \odot \mathbb{L})(A \times B|x) = \int_A \mathbb{L}(B|y, x) \mathbb{K}(dy|x) \quad \forall A \in \mathcal{Y}, B \in \mathcal{Z} \quad (2.16)$$

Lemma 2.2.23 (Semidirect product is associative). *Given $\mathbb{K} : X \rightarrow Y$, $\mathbb{L} : Y \times X \rightarrow Z$ and $\mathbb{M} : Z \times Y \times X \rightarrow W$*

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (2.17)$$

$$(2.18)$$

Proof.

$$(\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M} = \begin{array}{c} \text{Diagram showing the composition of kernels } \mathbb{K}, \mathbb{L}, \text{ and } \mathbb{M} \text{ for } (\mathbb{K} \odot \mathbb{L}) \odot \mathbb{M}. \end{array} \quad (2.19)$$

$$= \begin{array}{c} \text{Diagram showing the composition of kernels } \mathbb{K}, \mathbb{L}, \text{ and } \mathbb{M} \text{ for } \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}). \end{array} \quad (2.20)$$

$$= \mathbb{K} \odot (\mathbb{L} \odot \mathbb{M}) \quad (2.21)$$

□

The semidirect product can be used to define a notion of almost sure equality: two kernels $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are μ -almost surely equal if $\mu \odot \mathbb{K} = \mu \odot \mathbb{L}$. This is identical to the notion of almost sure equality in Cho and Jacobs (2019), who shows that under the assumption that (Y, \mathcal{Y}) is countably generated, $\mathbb{K} \stackrel{\mu}{\cong} \mathbb{L}$ if and only if $\mathbb{K} = \mathbb{L}$ μ -almost everywhere.

Definition 2.2.24 (Almost sure equality). Two Markov kernels $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are almost surely equal $\stackrel{\mathbb{P}_C}{\cong}$ with respect to a probability space (μ, X, \mathcal{X}) , written $\mathbb{K} \stackrel{\mu}{\cong} \mathbb{L}$ if

$$\mu \odot \mathbb{K} = \mu \odot \mathbb{L} \quad (2.22)$$

Theorem 2.2.25. *Given (μ, X, \mathcal{X}) , $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$, $\mathbb{K} \stackrel{\mu}{\cong} \mathbb{L}$ if and only if, defining $U := \{x | \exists A \in \mathcal{Y} : \mathbb{K}(A|x) \neq \mathbb{L}(A|x)\}$, $\mu(U) = 0$.*

Proof. Cho and Jacobs (2019) proposition 5.4. \square

We often want to talk about almost sure equality of two different versions \mathbb{K} and \mathbb{L} of a conditional distribution $\mathbb{P}^{Y|X}$ with respect to some ambient probability space $(\mathbb{P}, \Omega, \mathcal{F})$. This simply means \mathbb{K} and \mathbb{L} satisfy Definition 2.2.16 with respect to \mathbb{P} , X and Y , and they are almost surely equal with respect to the marginal \mathbb{P}^X . The relevant variables are usually obvious from the context and we leave them implicit and we will write $\mathbb{K} \stackrel{\mathbb{P}}{\cong} \mathbb{L}$. If the relevant marginal is ambiguous, we will instead write $\mathbb{K} \stackrel{\mathbb{P}^X}{\cong} \mathbb{L}$.

Definition 2.2.26 (Almost sure equality with respect to a pair of variables). Given $(\mathbb{P}, \Omega, \mathcal{F})$ and $X : \Omega \rightarrow X$, $Y : \Omega \rightarrow Y$, two Markov kernels $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are X -almost surely equal with respect to \mathbb{P} , written $\mathbb{K} \stackrel{\mathbb{P}}{\cong} \mathbb{L}$, if they are almost surely equal with respect to the marginal \mathbb{P}^X .

2.3 String Diagrams

We make use of string diagram notation for probabilistic reasoning. Graphical models are often employed in causal reasoning, and string diagrams are a kind of graphical notation for representing Markov kernels. The notation comes from the study of Markov categories, which are abstract categories that represent models of the flow of information. For our purposes, we don't use abstract Markov categories but instead focus on the concrete category of Markov kernels on standard measurable sets.

A coherence theorem exists for string diagrams and Markov categories. Applying planar deformation or any of the commutative comonoid axioms to a string diagram yields an equivalent string diagram. The coherence theorem establishes that any proof constructed using string diagrams in this manner corresponds to a proof in any Markov category (Selinger, 2011). More comprehensive introductions to Markov categories can be found in Fritz (2020); Cho and Jacobs (2019).

2.3.1 Elements of string diagrams

In the string, Markov kernels are drawn as boxes with input and output wires, and probability measures (which are Markov kernels with the domain $\{*\}$) are represented by triangles:

$$\mathbb{K} := \boxed{\mathbb{K}} \quad (2.23)$$

$$\mu := \triangleleft \mathbb{P} \quad (2.24)$$

Given two Markov kernels $\mathbb{L} : X \rightarrow Y$ and $\mathbb{M} : Y \rightarrow Z$, the product $\mathbb{L}\mathbb{M}$ is represented by drawing them side by side and joining their wires:

$$\mathbb{L}\mathbb{M} := X \begin{array}{|c|} \hline \mathbb{K} \\ \hline \mathbb{M} \\ \hline \end{array} Z \quad (2.25)$$

Given kernels $\mathbb{K} : W \rightarrow Y$ and $\mathbb{L} : X \rightarrow Z$, the tensor product $\mathbb{K} \otimes \mathbb{L} : W \times X \rightarrow Y \times Z$ is graphically represented by drawing kernels in parallel:

$$\mathbb{K} \otimes \mathbb{L} := \begin{array}{c} W \begin{array}{|c|} \hline \mathbb{K} \\ \hline \end{array} Y \\ X \begin{array}{|c|} \hline \mathbb{L} \\ \hline \end{array} Z \end{array} \quad (2.26)$$

Given $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : Y \times X \rightarrow Z$, the semidirect product is graphically represented by connecting \mathbb{K} and \mathbb{L} and keeping an extra copy

$$\mathbb{K} \odot \mathbb{L} := \text{copy}_X(\mathbb{K} \otimes \text{id}_X)(\text{copy}_Y \otimes \text{id}_X)(\text{id}_Y \otimes \mathbb{L}) \quad (2.27)$$

$$= X \begin{array}{c} \bullet \\ \downarrow \end{array} \begin{array}{|c|} \hline \mathbb{K} \\ \hline \end{array} \begin{array}{c} \bullet \\ \downarrow \end{array} \begin{array}{|c|} \hline \mathbb{L} \\ \hline \end{array} \begin{array}{c} Y \\ Z \end{array} \quad (2.28)$$

A space X is identified with the identity kernel $\text{id}^X : X \rightarrow \Delta(\mathcal{X})$. A bare wire represents the identity kernel:

$$\text{Id}^X := X \text{ ————— } X \quad (2.29)$$

Product spaces $X \times Y$ are identified with tensor product of identity kernels $\text{id}^X \otimes \text{id}^Y$. These can be represented either by two parallel wires or by a single wire representing the identity on the product space $X \times Y$:

$$X \times Y \cong \text{Id}^X \otimes \text{Id}^Y := \begin{array}{c} X \text{ — } X \\ Y \text{ — } Y \end{array} \quad (2.30)$$

$$= X \times Y \text{ ————— } X \times Y \quad (2.31)$$

A kernel $\mathbb{L} : X \rightarrow \Delta(\mathcal{Y} \otimes \mathcal{Z})$ can be written using either two parallel output wires or a single output wire, appropriately labeled:

$$X \text{ — } \begin{array}{|c|} \hline \mathbb{L} \\ \hline \end{array} \begin{array}{c} Y \\ Z \end{array} \quad (2.32)$$

$$\equiv \quad (2.33)$$

$$X \text{ — } \begin{array}{|c|} \hline \mathbb{L} \\ \hline \end{array} \text{ — } Y \times Z \quad (2.34)$$

We read diagrams from left to right (this is somewhat different to Fritz (2020); Cho and Jacobs (2019); Fong (2013) but in line with Selinger (2011)), and any diagram describes a set of nested products and tensor products of Markov kernels. There are a collection of special Markov kernels for which we can replace the generic “box” of a Markov kernel with a diagrammatic elements that are visually suggestive of what these kernels accomplish.

2.3.2 Special maps

Definition 2.3.1 (Identity map). The identity map $\text{id}_X : X \rightarrow X$ defined by $(\text{id}_X)(A|x) = \delta_x(A)$ for all $x \in X$, $A \in \mathcal{X}$, is represented by a bare line.

$$\text{id}_X := X \text{---} X \quad (2.35)$$

Definition 2.3.2 (Erase map). Given some 1-element set $\{*\}$, the erase map $\text{del}_X : X \rightarrow \{*\}$ is defined by $(\text{del}_X)(*|x) = 1$ for all $x \in X$. It “discards the input”. It looks like a lit fuse:

$$\text{del}_X := \text{---} * X \quad (2.36)$$

Definition 2.3.3 (Swap map). The swap map $\text{swap}_{X,Y} : X \times Y \rightarrow Y \times X$ is defined by $(\text{swap}_{X,Y})(A \times B|x, y) = \delta_x(B)\delta_y(A)$ for $(x, y) \in X \times Y$, $A \in \mathcal{X}$ and $B \in \mathcal{Y}$. It swaps two inputs and is represented by crossing wires:

$$\text{swap}_{X,Y} := \begin{array}{c} \diagup \diagdown \\ \diagdown \diagup \end{array} \quad (2.37)$$

Definition 2.3.4 (Copy map). The copy map $\text{copy}_X : X \rightarrow X \times X$ is defined by $(\text{copy}_X)(A \times B|x) = \delta_x(A)\delta_x(B)$ for all $x \in X$, $A, B \in \mathcal{X}$. It makes two identical copies of the input, and is drawn as a fork:

$$\text{copy}_X := X \text{---} \begin{array}{c} \swarrow \\ \searrow \end{array} \begin{array}{c} X \\ X \end{array} \quad (2.38)$$

Definition 2.3.5 (n -fold copy map). The n -fold copy map $\text{copy}_X^n : X \rightarrow X^n$ is given by the recursive definition

$$\text{copy}_X^1 = \text{copy}_X \quad (2.39)$$

$$\text{copy}_X^n = \begin{array}{c} \text{---} \boxed{\text{copy}_X^{n-1}} \text{---} \\ \bullet \diagdown \end{array} \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \quad n > 1 \quad (2.40)$$

Plates In a string diagram, a plate that is annotated $i \in A$ means the tensor product of the $|A|$ elements that appear inside the plate. A wire crossing from outside a plate boundary to the inside of a plate indicates an $|A|$ -fold copy map, which we indicate by placing a dot on the plate boundary. For our purposes, we do not define anything that allows wires to cross from the inside of a plate to the outside; wires must terminate within the plate.

Thus, given $\mathbb{K}_i : X \rightarrow Y$ for $i \in A$,

$$\bigotimes_{i \in A} \mathbb{K}_i := \boxed{\text{---} \boxed{\mathbb{K}_i} \text{---}}_{i \in A} \text{copy}_X^{|A|} \left(\bigotimes_{i \in A} \mathbb{K}_i \right) := \text{---} \bullet \boxed{\text{---} \boxed{\mathbb{K}_i} \text{---}}_{i \in A} \quad (2.41)$$

2.3.3 Commutative comonoid axioms

Diagrams in Markov categories satisfy the commutative comonoid axioms.

$$(2.42)$$

$$(2.43)$$

$$(2.44)$$

as well as compatibility with the monoidal structure

$$(2.45)$$

$$(2.46)$$

and the naturality of del , which means that

$$(2.47)$$

2.3.4 Manipulating String Diagrams

Planar deformations along with the applications of Equations 2.42 through to Equation 2.47 are almost the only rules we have for transforming one string diagram into an equivalent one. One further rule is given by Theorem 2.3.6.

Theorem 2.3.6 (Copy map commutes for deterministic kernels (Fong, 2013)).
For $\mathbb{K} : X \rightarrow Y$

$$(2.48)$$

holds iff \mathbb{K} is deterministic.

Examples

String diagrams can always be converted into definitions involving integrals and tensor products. A number of shortcuts can help to make the translations efficiently.

For arbitrary $\mathbb{K} : X \times Y \rightarrow Z$, $\mathbb{L} : W \rightarrow Y$

$$\begin{array}{c} \text{---} \\ \text{---} \end{array} \begin{array}{c} \boxed{\mathbb{K}} \\ \boxed{\mathbb{L}} \end{array} \text{---} = (\text{id}_X \otimes \mathbb{L})\mathbb{K} \quad (2.49)$$

$$[(\text{id}_X \otimes \mathbb{L})\mathbb{K}](A|x, w) = \int_Y \int_X \mathbb{K}(A|x', y') \mathbb{L}(dy'|w) \delta_x(dx') \quad (2.50)$$

$$= \int_Y \mathbb{K}(A|x, y') \mathbb{L}(dy'|w) \quad (2.51)$$

That is, an identity map “passes its input directly to the next kernel”.

For arbitrary $\mathbb{K} : X \times Y \times Y \rightarrow Z$:

$$\begin{array}{c} \text{---} \\ \text{---} \end{array} \bullet \begin{array}{c} \boxed{\mathbb{K}} \\ \bullet \end{array} \text{---} = (\text{id}_X \otimes \text{copy}_Y)\mathbb{K} \quad (2.52)$$

$$[(\text{id}_X \otimes \text{copy}_Y)\mathbb{K}](A|x, y) = \int_Y \int_Y \mathbb{K}(A|x, y', y'') \delta_y(dy') \delta_y(dy'') \quad (2.53)$$

$$= \mathbb{K}(A|x, y, y) \quad (2.54)$$

That is, the copy map “passes along two copies of its input” to the next kernel in the product.

For arbitrary $\mathbb{K} : X \times Y \rightarrow Z$

$$\begin{array}{c} \text{---} \\ \text{---} \end{array} \begin{array}{c} \boxed{\mathbb{K}} \\ \text{---} \end{array} = \text{swap}_{YX}\mathbb{K} \quad (2.55)$$

$$(\text{swap}_{YX}\mathbb{K})(A|y, x) = \int_{X \times Y} \mathbb{K}(A|x', y') \delta_y(dy') \delta_x(dx') \quad (2.56)$$

$$= \mathbb{K}(A|x, y) \quad (2.57)$$

The swap map before a kernel switches the input arguments.

For arbitrary $\mathbb{K} : X \rightarrow Y \times Z$

$$\text{---} \begin{array}{c} \boxed{\mathbb{K}} \\ \text{---} \end{array} \begin{array}{c} \text{---} \\ \text{---} \end{array} = \mathbb{K}\text{swap}_{YZ} \quad (2.58)$$

$$(\mathbb{K}\text{swap}_{YZ})(A \times B|x) = \int_{Y \times Z} \delta_y(B) \delta_z(A) \mathbb{K}(dy \times dz|x) \quad (2.59)$$

$$= \int_{B \times A} \mathbb{K}(dy \times dz|x) \quad (2.60)$$

$$= \mathbb{K}(B \times A|x) \quad (2.61)$$

Given $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : Y \rightarrow Z$:

$$(\mathbb{K} \odot \mathbb{L})(\text{id}_Y \otimes \text{del}_Z) = \begin{array}{c} X \text{ --- } \boxed{\mathbb{K}} \text{ --- } \bullet \begin{array}{l} \text{--- } Y \\ \text{--- } \boxed{\mathbb{L}} \text{ --- } * \end{array} \end{array} \quad (2.62)$$

$$= \begin{array}{c} X \text{ --- } \boxed{\mathbb{K}} \text{ --- } \bullet \begin{array}{l} \text{--- } Y \\ \text{--- } * \end{array} \end{array} \quad \text{by Eq. 2.47} \quad (2.63)$$

$$= X \text{ --- } \boxed{\mathbb{K}} \text{ --- } Y \quad \text{by Eq. 2.43} \quad (2.64)$$

Thus the action of the del map is to marginalise over the deleted wire. With integrals, we can write

$$(\mathbb{K} \odot \mathbb{L})(\text{id}_Y \otimes \text{del}_Z)(A \times \{*\}|x) = \int_Y \int_{\{*\}} \delta_y(A) \delta_*(\{*\}) \mathbb{L}(\text{d}z|y) \mathbb{K}(\text{d}y|x) \quad (2.65)$$

$$= \int_A \mathbb{K}(\text{d}y|x) \quad (2.66)$$

$$= \mathbb{K}(A|x) \quad (2.67)$$

2.4 Probability Sets

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

Definition 2.4.1 (Probability set). A probability set \mathbb{P}_C on (Ω, \mathcal{F}) is a collection of probability measures on (Ω, \mathcal{F}) . In other words it is a subset of $\mathcal{P}(\Delta(\Omega))$, where \mathcal{P} indicates the power set.

Given a probability set \mathbb{P}_C , we define marginal and conditional probabilities as probability measures and Markov kernels that satisfy Definitions 2.2.15 and 2.2.16 respectively for *all* base measures in \mathbb{P}_C . There are generally multiple Markov kernels that satisfy the properties of a conditional probability with respect to a probability set, and this definition ensures that marginal and conditional probabilities are “almost surely” unique (Definition 2.4.7) with respect to probability sets.

Definition 2.4.2 (Marginal probability with respect to a probability set). Given a sample space (Ω, \mathcal{F}) , a variable $X : \Omega \rightarrow X$ and a probability set \mathbb{P}_C , the marginal distribution $\mathbb{P}_C^X = \mathbb{P}_\alpha^X$ for any $\mathbb{P}_\alpha \in \mathbb{P}_C$ if a distribution satisfying this condition exists. Otherwise, it is undefined.

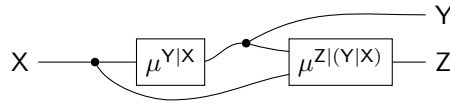
Definition 2.4.3 (Uniform conditional distribution). Given a sample space (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X$ and $Y : \Omega \rightarrow Y$ and a probability set \mathbb{P}_C , a uniform conditional distribution $\mathbb{P}_C^{Y|X}$ is any Markov kernel $X \rightarrow Y$ such that $\mathbb{P}_C^{Y|X}$ is an $Y|X$ conditional probability of \mathbb{P}_α for all $\mathbb{P}_\alpha \in \mathbb{P}_C$. If no such Markov kernel exists, $\mathbb{P}_C^{Y|X}$ is undefined.

Given a conditional distribution $\mu^{ZY|X}$ we can define a higher order conditional $\mu^{Z|(Y|X)}$, which is a version of $\mu^{Z|XY}$. This is useful because uniform conditionals don't always exist, but we can use higher order conditionals to show that if a probability set \mathbb{P}_C has a uniform conditional $\mathbb{P}_C^{ZY|X}$ then it also has a uniform conditional $\mathbb{P}_C^{Z|XY}$ (Theorems 2.4.26 and 2.4.27). Given $\mu^{XY|Z}$ and $X : \Omega \rightarrow X, Y : \Omega \rightarrow Y$ standard measurable, it has recently been proven that a higher order conditional $\mu^{Z|(Y|X)}$ exists Bogachev and Malofeev (2020), Theorem 3.5.

Definition 2.4.4 (Higher order conditionals). Given a probability space $(\mu, \Omega, \mathcal{F})$ and variables $X : \Omega \rightarrow X, Y : \Omega \rightarrow Y$ and $Z : \Omega \rightarrow Z$, a higher order conditional $\mu^{Z|(Y|X)} : X \times Y \rightarrow Z$ is any Markov kernel such that, for some $\mu^{Y|X}$,

$$\mu^{ZY|X}(B \times C|x) = \int_B \mu^{Z|(Y|X)}(C|x, y) \mu^{Y|X}(dy|x) \quad (2.68)$$

$$\iff \quad (2.69)$$



$$\mu^{ZY|X} = \quad (2.70)$$

Definition 2.4.5 (Uniform higher order conditional). Given a sample space (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X, Y : \Omega \rightarrow Y$ and $Z : \Omega \rightarrow Z$ and a probability set \mathbb{P}_C , if $\mathbb{P}_C^{ZY|X}$ exists then a uniform higher order conditional $\mathbb{P}_C^{Z|(Y|X)}$ is any Markov kernel $X \times Y \rightarrow Z$ that is a higher order conditional of some version of $\mathbb{P}_C^{ZY|X}$. If no $\mathbb{P}_C^{ZY|X}$ exists, $\mathbb{P}_C^{Z|(Y|X)}$ is undefined.

Definition 2.4.6 (Almost sure equality). Two Markov kernels $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are \mathbb{P}_C, X, Y -almost surely equal if for all $A \in \mathcal{X}, B \in \mathcal{Y}, \alpha \in C$

$$\int_A \mathbb{K}(B|x) \mathbb{P}_\alpha^X(dx) = \int_A \mathbb{L}(B|x) \mathbb{P}_\alpha^X(dx) \quad (2.71)$$

we write this as $\mathbb{K} \stackrel{\mathbb{P}_C}{\cong} \mathbb{L}$, as the variables X and Y are clear from the context.

Equivalently, \mathbb{K} and \mathbb{L} are almost surely equal if the set $C : \{x | \exists B \in \mathcal{Y} : \mathbb{K}(B|x) \neq \mathbb{L}(B|x)\}$ has measure 0 with respect to \mathbb{P}_α^X for all $\alpha \in C$.

2.4.1 Almost sure equality

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

Definition 2.4.7 (Almost sure equality). Two Markov kernels $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are almost surely equal $\stackrel{\mathbb{P}_C}{\cong}$ with respect to a probability set \mathbb{P}_C and variable $X : \Omega \rightarrow X$ if for all $\mathbb{P}_\alpha \in \mathbb{P}_C$,

$$\mathbb{P}_\alpha^X \odot \mathbb{K} = \mathbb{P}_\alpha^X \odot \mathbb{L} \quad (2.72)$$

Lemma 2.4.8 (Uniform conditional distributions are almost surely equal). *If $\mathbb{K} : X \rightarrow Y$ and $\mathbb{L} : X \rightarrow Y$ are both versions of $\mathbb{P}_C^{Y|X}$ then $\mathbb{K} \stackrel{\mathbb{P}_C}{\cong} \mathbb{L}$*

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

$$\mathbb{P}_\alpha^X \odot \mathbb{K} = \mathbb{P}_\alpha^{XY} \quad (2.73)$$

$$= \mathbb{P}_\alpha^X \odot \mathbb{L} \quad (2.74)$$

□

Lemma 2.4.9 (Substitution of almost surely equal Markov kernels). *Given \mathbb{P}_C , if $\mathbb{K} : X \times Y \rightarrow Z$ and $\mathbb{L} : X \times Y \rightarrow Z$ are almost surely equal $\mathbb{K} \stackrel{\mathbb{P}_C}{\cong} \mathbb{L}$, then for any $\mathbb{P}_\alpha \in \mathbb{P}_C$*

$$\mathbb{P}_\alpha^{Y|X} \odot \mathbb{K} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_\alpha^{Y|X} \odot \mathbb{L} \quad (2.75)$$

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

$$\mathbb{P}_\alpha^{XY} \odot \mathbb{K} \stackrel{\mathbb{P}_C}{\cong} (\mathbb{P}_\alpha^X \odot \mathbb{P}_C^{Y|X}) \odot \mathbb{K} \quad (2.76)$$

$$\stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_\alpha^X \odot (\mathbb{P}_C^{Y|X} \odot \mathbb{K}) \quad (2.77)$$

$$\stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_\alpha^X \odot (\mathbb{P}_C^{Y|X} \odot \mathbb{L}) \quad (2.78)$$

□

Theorem 2.4.10 (Semidirect product of uniform conditional distributions is a joint uniform conditional distribution). *Given a probability set \mathbb{P}_C on (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X$, $Y : \Omega \rightarrow Y$ and uniform conditional distributions $\mathbb{P}_C^{Y|X}$ and $\mathbb{P}_C^{Z|XY}$, then $\mathbb{P}_C^{YZ|X}$ exists and is equal to*

$$\mathbb{P}_C^{YZ|X} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{Y|X} \odot \mathbb{P}_C^{Z|XY} \quad (2.79)$$

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

$$\mathbb{P}_\alpha^{XYZ} = \mathbb{P}_\alpha^X \odot \mathbb{P}_\alpha^{YZ|X} \quad (2.80)$$

$$= \mathbb{P}_\alpha^X \odot (\mathbb{P}_\alpha^{Y|X} \odot \mathbb{P}_\alpha^{Z|YX}) \quad (2.81)$$

$$= \mathbb{P}_\alpha^X \odot (\mathbb{P}_C^{Y|X} \odot \mathbb{P}_C^{Z|YX}) \quad (2.82)$$

□

2.4.2 Extended conditional independence

Just like we defined uniform conditional probability as a version of “conditional probability” appropriate for probability sets, we need some version of “conditional independence” for probability sets. One such has already been given in some detail: it is the idea of *extended conditional independence* defined in Constantinou and Dawid (2017).

We will first define regular conditional independence. We define it in terms of a having a conditional that “ignores one of its inputs”, which, provided conditional probabilities exists, is equivalent to other common definitions (Theorem 2.4.12).

Definition 2.4.11 (Conditional independence). For a *probability model* \mathbb{P}_α and variables A, B, Z , we say B is conditionally independent of A given C , written $B \perp\!\!\!\perp_{\mathbb{P}_\alpha} A|C$, if

$$\mathbb{P}^{Y|WX} \stackrel{\mathbb{P}}{\cong} \begin{array}{c} W \text{ --- } \boxed{\mathbb{K}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (2.83)$$

$$\iff \mathbb{P}^{Y|WX}(A|w, x) \stackrel{\mathbb{P}}{\cong} \mathbb{K}(A|w) \quad \forall A \in \mathcal{Y} \quad (2.84)$$

Conditional independence can equivalently be stated in terms of the existence of a conditional probability that “ignores” one of its inputs.

Theorem 2.4.12. *Given standard measurable (Ω, \mathcal{F}) , a probability model \mathbb{P} and variables $W : \Omega \rightarrow W$, $X : \Omega \rightarrow X$, $Y : \Omega \rightarrow Y$, $Y \perp\!\!\!\perp_{\mathbb{P}} X|W$ if and only if there exists some version of $\mathbb{P}^{Y|WX}$ and $\mathbb{K} : W \rightarrow Y$ such that*

$$\mathbb{P}^{Y|WX} \stackrel{\mathbb{P}}{\cong} \begin{array}{c} W \text{ --- } \boxed{\mathbb{K}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (2.85)$$

$$\iff \mathbb{P}^{Y|WX}(A|w, x) \stackrel{\mathbb{P}}{\cong} \mathbb{K}(A|w) \quad \forall A \in \mathcal{Y} \quad (2.86)$$

Proof. See Cho and Jacobs (2019). □

Extended conditional independence as introduced by Constantinou and Dawid (2017) is defined in terms of “nonstochastic variables” on the choice set C . A nonstochastic variable is essentially a variable defined on C rather than on the sample space Ω

Definition 2.4.13 (Nonstochastic variable). Given a sample space (Ω, \mathcal{F}) , a choice set (C, \mathcal{C}) , a codomain (X, \mathcal{X}) and a probability set \mathbb{P}_C , a nonstochastic variable is a measurable function $\phi : C \rightarrow X$.

In particular, we want to consider *complementary* nonstochastic variable - that is, pairs of nonstochastic variables ϕ and ξ such that the sequence (ϕ, ξ) is invertible. For example, if $\phi := \text{idf}_C$, then

Definition 2.4.14 (Complementary nonstochastic variables). A pair of nonstochastic variables ϕ and ξ are complementary if (ϕ, ξ) is invertible.

Example 2.4.15. The letters ϕ and ξ are used to represent complementary nonstochastic variables.

Unlike Constantinou and Dawid (2017), we limit ourselves to a definition of extended conditional independence where regular uniform conditional probabilities exist. Our definition is otherwise identical.

Definition 2.4.16 (Extended conditional independence). Given a probability set \mathbb{P}_C , variables X, Y and Z and complementary nonstochastic variables ϕ and ξ , the extended conditional independence $Y \perp\!\!\!\perp_{\mathbb{P}_C}^e X\phi|Z\xi$ holds if for each $a \in \xi(C)$, $\mathbb{P}_{\xi^{-1}(a)}^{Y|XZ}$ and $\mathbb{P}_{\xi^{-1}(a)}^{Y|X}$ exist and

$$\mathbb{P}_{\xi^{-1}(a)}^{Y|XZ} \stackrel{\mathbb{P}_{\xi^{-1}(a)}}{\cong} \begin{array}{c} Z \text{ --- } \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (2.87)$$

$$\iff \quad (2.88)$$

$$\mathbb{P}_{\xi^{-1}(a)}^{Y|XZ}(A|x, z) \stackrel{\mathbb{P}_{\xi^{-1}(a)}}{\cong} \mathbb{P}_{\xi^{-1}(a)}^{Y|Z}(A|z) \quad \forall A \in \mathcal{Y}, (x, z) \in X \times Z \quad (2.89)$$

Very often, we consider a particular kind of extended conditional independence that does not explicitly make use of nonstochastic variables. We call this *uniform conditional independence*.

Definition 2.4.17 (Uniform conditional independence). Given a probability set \mathbb{P}_C and variables X, Y and Z , the uniform conditional independence $Y \perp\!\!\!\perp_{\mathbb{P}_C}^e XC|Z$ holds if $\mathbb{P}_C^{Y|XZ}$ and $\mathbb{P}_C^{Y|X}$ exist and

$$\mathbb{P}_C^{Y|XZ} \stackrel{\mathbb{P}_C}{\cong} \begin{array}{c} Z \text{ --- } \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (2.90)$$

$$\iff \quad (2.91)$$

$$\mathbb{P}_C^{Y|XZ}(A|x, z) \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{Y|Z}(A|z) \quad \forall A \in \mathcal{Y}, (x, z) \in X \times Z \quad (2.92)$$

For countable sets C (which, recall, is an assumption we generally accept), as shown by Constantinou and Dawid (2017) we can reason with collections of extended conditional independence statements as if they were regular conditional independence statements, with the provision that a complementary pair of nonstochastic variables must appear either side of the “|” symbol.

1. Symmetry: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|Z\xi$ iff $Y \perp\!\!\!\perp_{\mathbb{P}_C}^e X\phi|Z\xi$
2. $X \perp\!\!\!\perp_{\mathbb{P}_C}^e YC|YC$
3. Decomposition: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|W\xi$ and $Z \preceq Y$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Z\phi|W\xi$
4. Weak union:
 - (a) $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|W\xi$ and $Z \preceq Y$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|(Z, W)\xi$
 - (b) $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|W\xi$ and $\lambda \preceq \phi$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|(Z, W)(\xi, \lambda)$
5. Contraction: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Z\phi|W\xi$ and $X \perp\!\!\!\perp_{\mathbb{P}_C}^e Y\phi|(Z, W)\xi$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y, Z)\phi|W\xi$

The following forms of these properties are often used here:

1. Symmetry: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e YC|Z$ iff $Y \perp\!\!\!\perp_{\mathbb{P}}^e XC|Z$
2. Decomposition: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y, Z)C|W$ implies $X \perp\!\!\!\perp_{\mathbb{P}}^e YC|W$ and $X \perp\!\!\!\perp_{\mathbb{P}}^e ZC|W$
3. Weak union: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y, Z)C|W$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e YC|(Z, W)$
4. Contraction: $X \perp\!\!\!\perp_{\mathbb{P}_C}^e ZC|W$ and $X \perp\!\!\!\perp_{\mathbb{P}_C}^e YC|(Z, W)$ implies $X \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y, Z)C|W$

2.4.3 Examples

Example 2.4.18 (Choice variable). Suppose we have a decision procedure $\mathcal{S}_C := \{\mathcal{S}_\alpha | \alpha \in C\}$ that consists of a measurement procedure for each element of a denumerable set of choices C . Each measurement procedure \mathcal{S}_α is modeled by a probability distribution \mathbb{P}_α on a shared sample space (Ω, \mathcal{F}) such that we have an observable “choice” variable $(D, D \circ \mathcal{S}_\alpha)$ where $D \circ \mathcal{S}_\alpha$ always yields α .

Furthermore, Define $Y : \Omega \rightarrow \Omega$ as the identity function. Then, by supposition, for each $\alpha \in A$, \mathbb{P}_α^{YC} exists and for $A \in \mathcal{Y}$, $B \in \mathcal{C}$:

$$\mathbb{P}_\alpha^{YC}(A \times B) = \mathbb{P}_\alpha(A)\delta_\alpha(B) \quad (2.93)$$

This implies, for all $\alpha \in C$

$$\mathbb{P}_\alpha^{Y|D} = \mathbb{P}_\alpha^Y \quad (2.94)$$

Thus $\mathbb{P}_C^{Y|D}$ exists and

$$\mathbb{P}_C^{Y|D}(A|\alpha) = \mathbb{P}_\alpha^Y(A) \quad \forall A \in \mathcal{Y}, \alpha \in C \quad (2.95)$$

Because only deterministic marginals \mathbb{P}_α^D are available, for every $\alpha \in C$ we have $Y \perp\!\!\!\perp_{\mathbb{P}_\alpha} D$. This reflects the fact that *after we have selected a choice* α the value of C provides no further information about the distribution of Y , because D is deterministic given any α . It does not reflect the fact that “choosing different values of C has no effect on Y ”.

Theorem 2.4.19 (Uniform conditional independence representation). *Given a probability set \mathbb{P}_C with a uniform conditional probability $\mathbb{P}_C^{XY|Z}$,*

$$\mathbb{P}_C^{XY|Z} \stackrel{\mathbb{P}_C}{\cong} \begin{array}{c} Z \text{ --- } \bullet \begin{cases} \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ \boxed{\mathbb{P}_C^{X|Z}} \text{ --- } X \end{cases} \end{array} \quad (2.96)$$

$$\iff \quad (2.97)$$

$$\mathbb{P}_C^{XY|Z}(A \times B|z) \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{X|Z}(A|z) \mathbb{P}_C^{Y|Z}(B|z) \quad \forall A \in \mathcal{X}, B \in \mathcal{Y}, z \in Z \quad (2.98)$$

if and only if $Y \perp\!\!\!\perp_{\mathbb{P}_C}^e XC|Z$

Proof. If: By Theorem 2.4.27

$$\mathbb{P}_C^{XY|Z} = \begin{array}{c} Z \text{ --- } \bullet \begin{cases} \boxed{\mathbb{P}_C^{Y|ZX}} \text{ --- } Y \\ \boxed{\mathbb{P}_C^{X|Z}} \text{ --- } X \end{cases} \end{array} \quad (2.99)$$

$$\stackrel{\mathbb{P}_C}{\cong} \begin{array}{c} Z \text{ --- } \bullet \begin{cases} \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ \boxed{\mathbb{P}_C^{X|Z}} \text{ --- } X \end{cases} \end{array} \quad (2.100)$$

$$= \begin{array}{c} Z \text{ --- } \bullet \begin{cases} \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ \boxed{\mathbb{P}_C^{X|Z}} \text{ --- } X \end{cases} \end{array} \quad (2.101)$$

Only if: Suppose

$$\mathbb{P}_C^{XY|Z} \stackrel{\mathbb{P}_C}{\cong} \begin{array}{c} Z \text{ --- } \bullet \begin{cases} \boxed{\mathbb{P}_C^{Y|Z}} \text{ --- } Y \\ \boxed{\mathbb{P}_C^{X|Z}} \text{ --- } X \end{cases} \end{array} \quad (2.102)$$

and suppose for some $\alpha \in C$, $A \times C \in \mathcal{X} \otimes \mathcal{Z}$, $B \in \mathcal{Y}$ $\mathbb{P}_\alpha^{XZ}(A \times C) > 0$ and

$$\mathbb{P}_C^{Y|XZ}(B|x, z) > \mathbb{P}_C^{Y|Z}(B|z) \quad \forall (x, z) \in A \times C \quad (2.103)$$

then

$$\mathbb{P}_\alpha^{XYZZ}(A \times B \times C) = \int_{A \times C} \mathbb{P}_C^{Y|XZ}(B|x, z) \mathbb{P}_C^{X|Z}(dx|z) \mathbb{P}_\alpha^Z(dz) \quad (2.104)$$

$$> \int_{A \times C} \mathbb{P}_C^{Y|X}(B|z) \mathbb{P}_C^{X|Z}(dx|z) \mathbb{P}_\alpha^Z(dz) \quad (2.105)$$

$$= \int_C \mathbb{P}_C^{XY|X}(A \times B|z) \mathbb{P}_\alpha^Z(dz) \quad (2.106)$$

$$= \mathbb{P}_\alpha^{XYZZ}(A \times B \times C) \quad (2.107)$$

a contradiction. An analogous argument follows if we replace “>” with “<” in Eq. 2.103. \square

2.4.4 Maximal probability sets and valid conditionals

So far, we have been implicitly supposing that we first set up a probability set and from that set we may sometimes derive uniform conditional probabilities, extended conditional independences and so forth. However, sometimes we want to work backwards: start with a collection of uniform conditional probabilities, and work with the probability set implicitly defined by this collection. For example, when we have a Causal Bayesian Network, the collection of operations of the form “do($X = x$)” specify a probability set by a collection of uniform conditional probabilities on variables other than X , along with marginal probabilities of X . Specifically:

$$\mathbb{P}_{X=x}^{Y|Pa(Y)} = \begin{cases} \mathbb{P}_{\text{obs}}^{Y|Pa(Y)} & Y \text{ is a causal variable and not equal to } X \\ \delta_x & Y = X \end{cases} \quad (2.108)$$

The qualification “ Y is a causal variable” is usually not an explicit condition for causal Bayesian networks, but it is an important one. For example, $2X$ is not equal to X , but we cannot define a causal Bayesian network where both X and $2X$ are causal variables, see Example 2.4.23.

When working backwards like this, we can run into a couple of problems: we may end up with a probability set where some probabilities are non-unique, or we might inadvertently define an empty probability set. *Validity* is a condition that can ensure that we at least avoid the second problem.

Thus, if we start with a probability set, we know how to check if certain uniform conditional probabilities exist or not. However, there is a particular line of reasoning that comes up most often in the graphical models tradition of causal inference where we start with collections of conditional probabilities and assemble them into probability models as needed. A simple example of this is the causal Bayesian network given by the graph $X \longrightarrow Y$ and some observational probability distribution $\mathbb{P}^{XY} \in \Delta(X \times Y)$. Using the standard notion of “hard interventions on X ”, this model induces a probability set which we could informally describe as the set $\mathbb{P}_{\square} := \{\mathbb{P}_a^{XY} | a \in X \cup \{*\}\}$ where $*$ is a special element corresponding to the observational setting. The graph $X \longrightarrow Y$ implies the existence of the uniform conditional probability $\mathbb{P}_{\square}^{Y|X}$ under the nominated set of interventions, while the usual rules of hard interventions imply that $\mathbb{P}_a^X = \delta_a$ for $a \in X$.

Reasoning “backwards” like this – from uniform conditionals and marginals back to probability sets – must be done with care. The probability set associated with a collection of conditionals and marginals may be empty or nonunique. Uniqueness may not always be required, but an empty probability set is clearly not a useful model.

Consider, for example, $\Omega = \{0, 1\}$ with $X = (Z, Z)$ for $Z := \text{id}_{\Omega}$ and any measure $\kappa \in \Delta(\{0, 1\}^2)$ such that $\kappa(\{1\} \times \{0\}) > 0$. Note that $X^{-1}(\{1\} \times \{0\}) =$

$Z^{-1}(\{1\}) \cap Z^{-1}(\{0\}) = \emptyset$. Thus for any probability measure $\mu \in \Delta(\{0, 1\})$, $\mu^{\mathbf{X}}(\{1\} \times \{0\}) = \mu(\emptyset) = 0$ and so κ cannot be the marginal distribution of \mathbf{X} for any base measure at all.

We introduce the notion of *valid distributions* and *valid conditionals*. The key result here is: probability sets defined by collections of recursive valid conditionals and distributions are nonempty. While we suspect this condition is often satisfied by causal models in practice, we offer one example in the literature where it apparently is not. The problem of whether a probability set is valid is analogous to the problem of whether a probability distribution satisfying a collection of constraints exists discussed in Vorobev (1962). As that work shows, there are many questions of this nature that can be asked and that are not addressed by the criterion of validity.

There is also a connection between the notion of validity and the notion of *unique solvability* in Bongers et al. (2016). We ask “when can a set of conditional probabilities together with equations be jointly satisfied by a probability model?” while Bongers et. al. ask when a set of equations can be jointly satisfied by a probability model.

Definition 2.4.20 (Valid distribution). Given (Ω, \mathcal{F}) and a variable $\mathbf{X} : \Omega \rightarrow X$, an \mathbf{X} -valid probability distribution is any probability measure $\mathbb{K} \in \Delta(X)$ such that $\mathbf{X}^{-1}(A) = \emptyset \implies \mathbb{K}(A) = 0$ for all $A \in \mathcal{X}$.

Definition 2.4.21 (Valid conditional). Given (Ω, \mathcal{F}) , $\mathbf{X} : \Omega \rightarrow X$, $\mathbf{Y} : \Omega \rightarrow Y$ a $\mathbf{Y}|\mathbf{X}$ -valid conditional probability is a Markov kernel $\mathbb{L} : X \rightarrow Y$ that assigns probability 0 to impossible events, unless the argument itself corresponds to an impossible event:

$$\forall B \in \mathcal{Y}, x \in X : (\mathbf{X}, \mathbf{Y}) \bowtie \{x\} \times B = \emptyset \implies (\mathbb{L}(B|x) = 0) \vee (\mathbf{X} \bowtie \{x\} = \emptyset) \quad (2.109)$$

Definition 2.4.22 (Maximal probability set). Given (Ω, \mathcal{F}) , $\mathbf{X} : \Omega \rightarrow X$, $\mathbf{Y} : \Omega \rightarrow Y$ and a $\mathbf{Y}|\mathbf{X}$ -valid conditional probability $\mathbb{L} : X \rightarrow Y$ the maximal probability set $\mathbb{P}_C^{\mathbf{Y}|\mathbf{X}[M]}$ associated with \mathbb{L} is the probability set such that for all $\mathbb{P}_\alpha \in \mathbb{P}_C$, \mathbb{L} is a version of $\mathbb{P}_\alpha^{\mathbf{Y}|\mathbf{X}}$.

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

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

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

Note that validity is not a necessary condition for a conditional to define a non-empty probability set. The intuition for this is: if we have some $\mathbb{K} : X \rightarrow Y$,

\mathbb{K} might be an invalid $\mathbf{Y}|\mathbf{X}$ conditional on all of X , but might be valid on some subset of X , and so we might have some probability model \mathbb{P} that assigns measure 0 to the bad parts of X such that \mathbb{K} is a version of $\mathbb{P}^{\mathbf{Y}|\mathbf{X}}$. On the other hand, if we want to take the product of \mathbb{K} with arbitrary valid \mathbf{X} probabilities, then the validity of \mathbb{K} is necessary (Theorem ??).

Example 2.4.23. Body mass index is defined as a person's weight divided by the square of their height. Suppose we have a measurement process $\mathcal{S} = (\mathcal{W}, \mathcal{H})$ and $\mathcal{B} = \frac{\mathcal{W}}{\mathcal{H}^2}$ - i.e. we figure out someone's body mass index first by measuring both their height and weight, and then passing the result through a function that divides the second by the square of the first. Thus, given the random variables \mathbf{W}, \mathbf{H} modelling \mathcal{W}, \mathcal{H} , \mathcal{B} is the function given by $\mathbf{B} = \frac{\mathbf{W}}{\mathbf{H}^2}$.

With this background, suppose we postulate a decision model in which body mass index can be directly controlled by a variable \mathbf{C} , while height and weight are not. Specifically, we have a probability set \mathbb{P}_{\square} with

$$\mathbb{P}_{\square}^{\mathbf{B}|\mathbf{WHC}} = \begin{array}{c} \mathbf{H} \text{ ---} * \\ \mathbf{C} \text{ ---} \text{-----} \mathbf{B} \\ \mathbf{W} \text{ ---} * \end{array} \quad (2.110)$$

Then pick some $w, h, x \in \mathbb{R}$ such that $\frac{w}{h^2} \neq x$ and $(\mathbf{W}, \mathbf{H}) \bowtie (w, h) \neq \emptyset$ (which is to say, our measurement procedure could potentially yield (w, h) for a person's height and weight). We have $\mathbb{P}_{\square}^{\mathbf{B}|\mathbf{WHC}}(\{x\}|w, h, x) = 1$, but

$$(\mathbf{B}, \mathbf{W}, \mathbf{H}) \bowtie \{(x, w, h)\} = \{\omega | (\mathbf{W}, \mathbf{H})(\omega) = (w, h), \mathbf{B}(\omega) = \frac{w}{h^2}\} \quad (2.111)$$

$$= \emptyset \quad (2.112)$$

so $\mathbb{P}_{\square}^{\mathbf{B}|\mathbf{WHC}}$ is invalid. Thus there is some valid $\mu^{\mathbf{WHC}}$ such that the probability set $\mathbb{P}_{\square}^{\mathbf{B}|\mathbf{WHC}} = \mu^{\mathbf{WHC}} \odot \mathbb{P}_{\square}^{\mathbf{Y}|\mathbf{X}}$ is empty.

Validity rules out conditional probabilities like 2.110. We conjecture that in many cases this condition is implicitly taken into account – it is obviously silly to posit a model in which body mass index can be controlled independently of height and weight. We note, however, that presuming the authors intended their model to be interpreted according to the usual semantics of causal Bayesian networks, the invalid conditional probability 2.110 would be used to evaluate the causal effect of body mass index in the causal diagram found in Shahar (2009).

2.4.5 Existence of conditional probabilities

Lemma 2.4.24 (Conditional pushforward). *Suppose we have a sample space (Ω, \mathcal{F}) , variables $\mathbf{X} : \Omega \rightarrow X$ and $\mathbf{Y} : \Omega \rightarrow Y$, $\mathbf{Z} : \Omega \rightarrow Z$ and a probability set \mathbb{P}_C with conditional $\mathbb{P}_C^{\mathbf{X}|\mathbf{Y}}$ such that $\mathbf{Z} = f \circ \mathbf{Y}$ for some $f : Y \rightarrow Z$. Then there exists a conditional probability $\mathbb{P}_C^{\mathbf{Z}|\mathbf{X}} = \mathbb{P}_C^{\mathbf{Y}|\mathbf{X}} \mathbb{F}_f$.*

Proof. Note that $(X, Z) = (\text{id}_X \otimes f) \circ (X, Y)$. Thus, by Lemma 2.2.21, for any $\mathbb{P}_\alpha \in \mathbb{P}_C$

$$\mathbb{P}_\alpha^{XZ} = \mathbb{P}_\alpha^{XY} \mathbb{F}_{\text{id}_X \otimes f} \quad (2.113)$$

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

$$\mathbb{F}_{\text{id}_X \otimes f}(A \times B|x, y) = \delta_x(A) \delta_{f(y)}(B) \quad (2.114)$$

$$= \mathbb{F}_{\text{id}_X}(A|x) \otimes \mathbb{F}_f(B|y) \quad (2.115)$$

$$\implies \mathbb{F}_{\text{id}_X \otimes f} = \mathbb{F}_{\text{id}_X} \otimes \mathbb{F}_f \quad (2.116)$$

Thus

$$\mathbb{P}_\alpha^{XZ} = (\mathbb{P}_\alpha^X \odot \mathbb{P}_C^{Y|X}) \mathbb{F}_{\text{id}_X} \otimes \mathbb{F}_f \quad (2.117)$$

$$= \begin{array}{c} \text{--- } X \\ \diagup \\ \triangleleft \mathbb{P}_\alpha^X \text{---} \bullet \text{---} \boxed{\mathbb{P}_C^{Y|X}} \text{---} \boxed{\mathbb{F}_f} \text{---} Z \\ \diagdown \\ \text{--- } \end{array} \quad (2.118)$$

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

Theorem 2.4.25 (Existence of regular conditionals). *Suppose we have a sample space (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X$ and $Y : \Omega \rightarrow Y$ with Y standard measurable and a probability model \mathbb{P}_α on (Ω, \mathcal{F}) . Then there exists a conditional $\mathbb{P}_\alpha^{Y|X}$.*

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

Theorem 2.4.26 (Existence of higher order valid conditionals with respect to probability sets). *Suppose we have a sample space (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X$ and $Y : \Omega \rightarrow Y$, $Z : \Omega \rightarrow Z$ and a probability set \mathbb{P}_C with regular conditional $\mathbb{P}_C^{YZ|X}$ and Y and Z standard measurable. Then there exists a regular $\mathbb{P}_C^{Z|(Y|X)}$.*

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

$$n(f^{-1}(y)|x, y) = 1 \quad (2.119)$$

$$m(Y^{-1}(A) \cap B|x) = \int_A n(B|x, y) m \mathbb{F}_f(dy|x) \forall A \in \mathcal{Y}, B \in \mathcal{Y} \times \mathcal{Z} \quad (2.120)$$

In particular, $\mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{X}}$ is a Borel measurable map $X \rightarrow Y \times Z$. Thus equation 2.120 implies for all $A \in \mathcal{Y}, B \in \mathcal{Y} \times \mathcal{Z}$

$$\mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{X}}(\mathbf{Y}^{-1}(A) \cap B|x) = \int_A n(B|x, y) \mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{X}} \mathbb{F}_f(dy|x) \quad (2.121)$$

$$= \int_A n(B|x, y) \mathbb{P}_C^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (2.122)$$

Where Equation 2.122 follows from Lemma 2.4.24.

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

$$\mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{X}}(\mathbf{Y}^{-1}(A) \cap B|x) = \int_A n(B|x, y) \mathbb{P}_\alpha^{\mathbf{Y}|\mathbf{X}}(dy|x) \quad (2.123)$$

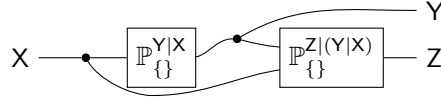
which implies n is a version of $\mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$. By Lemma 2.4.24, $n\mathbb{F}_f$ is a version of $\mathbb{P}_C^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$. \square

We might be motivated to ask whether the higher order conditionals in Theorem 2.4.26 can be chosen to be valid. Despite Lemma ?? showing that the existence of proper conditional probabilities implies the existence of valid ones, we cannot make use of this in the above theorem because Equation 2.119 makes n proper with respect to the “wrong” sample space $(Y \times Z, \mathcal{Y} \otimes \mathcal{Z})$ while what we would need is a proper conditional probability with respect to (Ω, \mathcal{F}) .

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

Theorem 2.4.27 (Higher order conditionals). *Suppose we have a sample space (Ω, \mathcal{F}) , variables $\mathbf{X} : \Omega \rightarrow X$ and $\mathbf{Y} : \Omega \rightarrow Y$, $\mathbf{Z} : \Omega \rightarrow Z$ and a probability set \mathbb{P}_C with conditional $\mathbb{P}_C^{\mathbf{Y}\mathbf{Z}|\mathbf{X}}$. Then $\mathbb{P}_C^{\mathbf{Z}|\mathbf{Y}|\mathbf{X}}$ is a version of $\mathbb{P}_C^{\mathbf{Z}|\mathbf{Y}\mathbf{X}}$*

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



$$\mathbb{P}_\alpha^{YZ|X} = \quad (2.124)$$

$$\Rightarrow \mathbb{P}_\alpha^{XYZ} = \triangleleft \mathbb{P}_\alpha^X \quad (2.125)$$

$$= \quad (2.126)$$

$$= \quad (2.127)$$

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

Theorem 2.4.28. *Given probability gap model \mathbb{P}_C , X, Y, Z such that $\mathbb{P}_C^{Z|YX}$ exists, $\mathbb{P}_C^{Z|Y}$ exists iff $Z \perp\!\!\!\perp_{\mathbb{P}_C} X|Y$.*

Proof. If: If $Z \perp\!\!\!\perp_{\mathbb{P}_C} X|Y$ then by Theorem 2.4.12, for each $\mathbb{P}_\alpha \in \mathbb{P}_C$ there exists $\mathbb{P}_\alpha^{Z|Y}$ such that

$$\mathbb{P}_\alpha^{Y|WX} = \begin{array}{c} W \text{ --- } \boxed{\mathbb{K}} \text{ --- } Y \\ X \text{ --- } * \end{array} \quad (2.128)$$

\square

Theorem 2.4.29 (Valid higher order conditionals). *Suppose we have a sample space (Ω, \mathcal{F}) , variables $X : \Omega \rightarrow X$ and $Y : \Omega \rightarrow Y$, $Z : \Omega \rightarrow Z$ and a probability set \mathbb{P}_C with regular conditional $\mathbb{P}_C^{YZ|X}$, Y discrete and Z standard measurable. Then there exists a valid regular $\mathbb{P}_C^{Z|XY}$.*

Proof. By Theorem 2.4.26, we have a higher order conditional $\mathbb{P}_C^{Z|(Y|X)}$ which, by Theorem 2.4.27 is also a version of $\mathbb{P}_C^{Z|XY}$.

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

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

$$\forall A \in \mathcal{Z}, x \in X, y \in Y : \mathbb{P}_C^{Y|X}(\{y\} \times A|x) = 0 \implies (\mathbb{Q}(A|x, y) = 0) \vee ((X, Y) \bowtie \{(x, y)\} = \emptyset) \quad (2.129)$$

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

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

Chapter 3

Models with choices and consequences

Probability sets, introduced in Chapter 2, will be used to model *decision problems*, which are problems that involve choices and consequences. In such problems, three things are given: a set of options, one of which must be chosen, a set of consequences and a means of judging which consequences are more desirable than others. Such a problem requires an understanding of how each choice corresponds to consequences, as far as this is able to be understood. The fundamental type of problem studied in this thesis is how to map choices to consequences.

In practice causal inference is concerned with a wider variety of problems than this. A great deal of empirical causal analysis is concerned with problems a step removed from this: the purpose is to advise other decision makers on a course of action rather than to recommend an action directly. Nevertheless, many theorists of causal inference signal an awareness that much of their analysis is ultimately motivated by problems involving a choice among options, even if such problems are not always directly addressed. Section 3.1 reviews briefly the role of decision problems according to theorists of causal inference and sets out the basic scheme by which probability sets are used to model problems of this type.

The reasons we provide for focusing on probability sets are not rigorous. The strongest motivation for this choice is *convention*: this chapter shows how many varieties of decision theory induce probability set models, and Chapter ?? shows how many causal inference frameworks also induce probability set models. Some decision theories examined in this chapter seek to justify their modelling choices by suggesting axioms for rational theories of decision under uncertainty, and deriving a particular type of model from these axioms. However, despite the various attempts at axiomatisation, the nature of such a theory is contested – there is no clear standard among the variety of theories surveyed here, nor among other proposed theories not considered in this work. This work is not

trying to resolve this dispute. However, modelling choices must still be made and it is still reasonable to ask what is being assumed by making these choices, so Section 3.2 provides an overview of several major decision theories along with their axiomatisations.

Section 3.2 describes in particular detail the connections between *statistical decision theory* (Wald, 1950) and probability set models of decision problems. There are two reasons for this: firstly, we are able to demonstrate a close connection between probability set models of decision problems and the classical statistical notion of *risk* of a decision rule, even though causal considerations are often not central to classical statistics. Secondly, the kind of probability set model – which we call a *see-do model* – shows up again in Chapter 4 where we consider the question of when a probability set model supports a certain notion of “the causal effect of a variable”, and again in Chapter ?? where we consider the kinds of probability set models induced by other causal reasoning frameworks.

The formal definition of a variable in a probabilistic model is straightforward (Definition 2.2.13). However, in practice the definitions of variables often includes informal content that enables the interpretation of a probabilistic model. In the field of causal models, one is likely to come across many different “kinds” of variables: for example, observed variables, unobserved variables, counterfactual variables and causal variables all play important roles in various causal inference frameworks. However, there is no formal distinction between these different kinds of variables – Definition 2.2.13 applies to them all. Section 3.3 is an attempt to clarify an understanding of the informal role of variables as “pointing to the parts of the world that the model is about”. In comparison to the wide variety of variable types encountered in the causal literature, it offers a very limited theory of the informal semantics of variables. In short, observed variables correspond to a measurement procedure (in a sense that will be made precise), and unobserved variables do not.

3.1 What is the point of causal inference?

Pearl and Mackenzie (2018) argues forcefully that causal reasoning frameworks should be understood by the questions that they answer. He also posits a “ladder” of types of causal question, where the n th level of question type also subsumes all lower levels. The question types are (Bareinboim et al., 2020):

1. *Associational*: informally, “questions about relationships and predictions”; formally defined as queries that can be answered by a single probability distribution
2. *Interventional*: informally, “questions about the consequences of interventions”; formally defined as queries that can be answered by a causal Bayesian network (CBN)
3. *Counterfactual*: informally, “questions concerning imagining alternate worlds”;

formally defined as queries that can be answered by a structural causal model (SCM)

As with theories of rational decision making under uncertainty, adjudicating whether this is indeed a sound basis for causal inference is not the purpose of this thesis. The approach to causal inference presented here is motivated by the informal version of rung 2, rather than the informal version of rung 3 (the formal theory is founded on probability sets, not CBNs or SCMs). It is easy to doubt that rung 3 is the right informal starting point for a theory of causal inference: the argument that the ladder is hierarchical depends on the formal characterisation of rungs in terms of questions that can be posed *with particular kinds of models*. This gets priorities backwards; the point is to find models that are fit for answering important questions, not to find questions that can be asked of important models.

Furthermore, the informal questions that characterise rung 3 can be posed in language more reminiscent of rungs 2. Rung 2 is characterised, informally, by questions like “if I take aspirin, will my headache go away?” while rung 3 is characterised by questions like “if I had taken the aspirin, would I still have a headache?”. However, the second question could be posed as “if I wave a magic wand and transport myself to a world where I have taken aspirin, would I still have a headache?”. Such an intervention is fanciful, but it nevertheless sounds like an intervention. Thus, perhaps, one can characterise rung 3 by questions concerning interventions that may or may not be fanciful.

Rung 2 questions are far more prominent in the causal inference literature than rung 3. Rubin (2005) argues that causal inference often informs a decision maker by providing “scientific knowledge”, but does not make recommendations by itself. (Imbens and Rubin, 2015) introduces causal inference as the study of “outcomes of manipulations” and (Spirtes et al., 2000) highlights the universal relevance of understanding how to control certain outcomes, while further arguing that clarifying commonsense ideas of causation is also an important aim of causal inference. Hernán and Robins (2020) present causal knowledge as critical for assessing the consequences of actions. Problems that involve comparing different choices on the basis of their consequences are an important class of problems. They are worth better understanding, even if this work does not shed any light on counterfactual questions.

3.1.1 Modelling decision problems

People who need to make a decisions might (and often do) make them with no mathematical reasoning at all. However, this work is concerned with making decisions assisted by mathematical reasoning. In order to reason mathematically about a decision to be made, we assume that somehow, we have access to two sets:

1. There is a set of choices C that need to be compared
2. There is a set of consequences Ω along with a utility function $u : \Omega \rightarrow \mathbb{R}$

Given some means of relating between C and Ω , the order on Ω will induce some order on C . There are a great number of different ways that of relating elements of C to Ω . For example, a binary relation between the two sets will, given a total order on Ω , induce a preorder on C . However, in this work the assumption is made that the relevant kinds of relations are either

- A Markov kernel $C \rightarrow \Omega$
- A Markov kernel $C \times H \rightarrow \Omega$ for some set of hypotheses H

That is, for each choice $c \in C$ we have either a probability distribution in $\Delta(\Omega)$ or a set of probability distributions indexed by $h \in H$. Sections 3.2.5 and 3.2.5 discuss each choice in more detail. Where it is needed, we also assume that a utility function $\Omega \rightarrow \mathbb{R}$ is available and that choices are evaluated using the principle of expected utility.

Usually, someone confronted with a decision problem will not know for certain the consequences that arise from any given choice, and yet they may have some views about which consequences are more likely than others. Probability has a long and successful history of representing uncertain knowledge of this type. There are many works that aim to show that any method for representing uncertain knowledge that adheres to certain principles must be a probability distribution de Finetti ([1937] 1992); Horvitz et al. (1986), along with criticism of these principles Halpern (1999). A notable alternative to representing uncertainty with a single probability distribution represents uncertainty with a set of probability distributions, which is a type of *vague probability* model (Walley, 1991).

More relevant to the question of modelling decision problems are a number of works that establish conditions under which “desirability” or “preference” relations over sets of choices or propositions must be represented by a probability distribution along with a utility function. These works are surveyed in Section 3.2. Ultimately, however, the question of whether probability is the right choice to represent uncertain knowledge in decision models is not a key focus of this work. It is a conventional choice, and one that is accepted here.

3.1.2 Formal definitions

We suppose that we are, at the outset, given a few basic ingredients: a set of choices C , a set of consequences Ω and a utility function $u : F \rightarrow \mathbb{R}$. We call these ingredients a “decision problem”.

Definition 3.1.1 (Decision problem). A decision problem is a triple (C, Ω, u) consisting of a measurable set (C, \mathcal{C}) of choices, (Ω, \mathcal{F}) consequences and a utility function $u : F \rightarrow \mathbb{R}$.

Our task is to find a *model* that relates C to Ω . We assume two forms of model – a *sharp model* associates each choice with a unique probability distribution, and a *vague model* associates each choice with a set of probability distributions.

Definition 3.1.2 (Sharp model). Given a decision problem (C, Ω, u) , a *sharp model* is a function $C \rightarrow \Omega$.

Definition 3.1.3 (Vague model). Given a decision problem (C, Ω, u) , a *vague model* is a function $C \times H \rightarrow \Omega$ for some hypothesis set H .

Both sharp models and vague models have probability sets induced by the range of the model.

Definition 3.1.4 (Induced probability set). Given a decision problem (C, Ω, u) and a model $\mathbb{P} : C \times H \rightarrow \Omega$, the induced probability set is $\mathbb{P}_{C \times H} := \{\mathbb{P}_\alpha | \alpha \in C \times H\}$.

3.2 Representation theorems for decision problems

We assume decision models are probabilistic functions $C \rightarrow \Delta(\Omega)$ for some sample space (Ω, \mathcal{F}) of “consequences”. Probability distributions, and the principle of expected utility in particular, are common choices for evaluation under uncertainty. Representation theorems offer a more formal justification for this choice; they propose a collection of axioms regulating the sets of evaluations we want some decision evaluation model to admit, and then show that this model can be represented with (for example) a probability distribution along with a utility function. The desirability of some of the axioms in these theorems is not obvious.

Here we review the representation theorems of Savage (1954) and Jeffrey (1965). We establish that both imply that choices are compared using a probabilistic function $C \rightarrow \Delta(\Omega)$ for a suitable selection of C and (Ω, \mathcal{F}) , along with a “desirability” function which differs in type between the two theorems.

Lewis’ *causal decision theory* is also briefly reviewed. While the particular considerations that motivated this theory are not examined, this theory introduces *dependency hypotheses*, which play a key role in the rest of this work.

The following discussion will often make reference to *complete preference relations*. A complete preference relation is a relation \succ, \prec, \sim on a set A such that for any a, b, c in A we have:

- Exactly one of $a \succ b$, $a \prec b$, $a \sim b$ holds
- $(a \succ b) \iff (b \prec a)$
- $a \succ b$ and $b \succ c$ implies $a \succ c$

In short, it is a total order without antisymmetry (a and b can be equally preferred even if they are not in fact equal).

This definition is meant to correspond to the common sense idea of having preferences over some set of things, where \succ can be read as “strictly better than”, \prec read as “strictly worse than” and \sim read as “as good as”. Given any

two things from the set, I can say which one I prefer, or if I prefer neither (and all of these are mutually exclusive). If I prefer a to a' then I think a' is worse than a . Furthermore, if I prefer a to a' and a' to a'' then I prefer a to a'' .

Define $a \preceq b$ to mean $a \prec b$ or $a \sim b$.

3.2.1 von Neumann-Morgenstern utility

Von Neumann and Morgenstern (1944) proved that when the *vNM axioms* hold (not defined here; see the original reference or Steele and Stefánsson (2020)), an agent's preferences between “lotteries” (probability distributions in $\Delta(\Omega)$ for some (Ω, \mathcal{F})) can be represented as the comparison of the expected value under each lottery of a utility function u unique up to affine transformation. That is, for lotteries \mathbb{P}_α and $\mathbb{P}_{\alpha'}$, there exists some $u : \Omega \rightarrow \mathbb{R}$ unique up to affine transformation such that $\mathbb{E}_{\mathbb{P}_\alpha}[u] > \mathbb{E}_{\mathbb{P}_{\alpha'}}[u]$ if and only if $\mathbb{P}_\alpha \succ \mathbb{P}_{\alpha'}$.

In vNM theory, the set of lotteries is the set of all probability measures on (Ω, \mathcal{F}) . Thus von Neumann-Morgenstern theorem gives conditions under which preferences *over distributions of consequences* can be represented using expected utility. It is silent on the question of whether each choice should be mapped to a unique probability distribution over consequences.

3.2.2 Savage decision theory

Savage's decision theory establishes conditions under which, given *acts* C , *consequences* Ω and *states* (S, \mathcal{S}) (which are “possible mappings from acts to consequences”), the preference relation over acts can be represented with a probability distribution over states and a utility function $\Omega \rightarrow \mathbb{R}$. This is much closer to the subject of this work than the theorem of von Neumann and Morgenstern.

Definition 3.2.1 (Elements of a Savage decision problem). A *Savage decision problem* features a measurable set of states (S, \mathcal{S}) , a set of consequences Ω and a set of acts C such that $|C| = \Omega^S$ and an evaluation function $T : S \times C \rightarrow F$ such that for any $f : S \rightarrow \Omega$ there exists $c \in C$ such that $T(\cdot, c) = f$.

Theorem 3.2.2. *Given any Savage decision problem (S, Ω, C, T) with a preference relation (\prec, \sim) on C that satisfies the Savage axioms 3.2.2, there exists a unique probability distribution $\mu \in \Delta(\mathcal{S})$ and a utility $u : \Omega \rightarrow \mathbb{R}$ unique up to affine transformation such that*

$$\alpha \preceq \alpha' \iff \int_S u(T(s, \alpha)) \mu(ds) \leq \int_S u(T(s, \alpha')) \mu(ds) \quad \forall \alpha, \alpha' \in C \quad (3.1)$$

Proof. Savage (1954) □

If we equip consequences with a measures (Ω, \mathcal{F}) , Savage's setup implies the existence of a unique probabilistic function $C \rightarrow \Delta(\Omega)$ representing the “probabilistic consequences” of each choice.

Theorem 3.2.3. *Given any Savage decision problem (S, Ω, C, T) with a preference relation (\prec, \sim) on C that satisfies the Savage axioms, and a σ -algebra \mathcal{F} on Ω such that T is measurable, there is a probabilistic function $\mathbb{P} : C \rightarrow \Delta(\Omega)$ and a utility $u : \Omega \rightarrow \mathbb{R}$ unique up to affine transformation such that*

$$\alpha \preceq \alpha' \iff \int_{\Omega} u(f) \mathbb{P}_{\alpha}(df) \leq \int_{\Omega} u(f) \mathbb{P}_{\alpha'}(df) \quad \forall \alpha, \alpha' \in C \quad (3.2)$$

Proof. Define $\mathbb{P} : C \rightarrow \Delta(\Omega)$ by

$$\mathbb{P}_{\alpha}(A) := \mu(T_{\alpha}^{-1}(A)) \quad \forall A \in \mathcal{F} \quad (3.3)$$

where $T_{\alpha} : S \rightarrow F$ is the function $s \mapsto T(s, \alpha)$. \mathbb{P}_{α} is the pushforward of T_{α} under μ .

Then

$$\int_{\Omega} u(f) \mathbb{P}_{\alpha}(df) = \int_S u \circ T_{\alpha}(s) \mu(ds) \quad (3.4)$$

$$= \int_S u(T(s, \alpha)) \mu(ds) \quad (3.5)$$

□

Savage axioms

Careful analysis of Savage's theorem is outside the scope of this work, but given the relevance of Savage's representation theorem we will reproduce the axioms from Savage (1954) with a small amount of commentary. Keep in mind that Savage's theorem establishes that the following are sufficient for representation with a probability set, not necessary, and furthermore the probability set representation of preferences satisfying these axioms is unique.

Given acts C , states (S, \mathcal{S}) and consequences F and a map $T : S \times C \rightarrow F$, let all greek letters α, β etc. be elements of C . Savage's axioms are:

P1: There is a complete preference relation \preceq on C

D1: $\alpha \preceq \beta$ given $B \in \mathcal{S}$ if and only if $\alpha' \preceq \beta'$ for every α' and β' such that $T(\alpha, s) = T(\alpha', s)$ for $s \in B$ and $T(\alpha', r) = T(\beta', r)$ for $r \notin B$, and $\beta' \preceq \alpha'$ either for every such pair or for none.

P2: For every α, β and $B \in \mathcal{S}$, $\alpha \preceq \beta$ given B or $\beta \preceq \alpha$ given B

D2: for $q, q' \in F$, $q \preceq q'$ if and only if $\alpha \preceq \alpha'$ where $T(\alpha, s) = q$ and $T(\alpha', s) = q'$ for all $s \in S$

D2: $B \in \mathcal{S}$ is null if and only if $\alpha \preceq \beta$ given B for every $\alpha, \beta \in C$

P3: If $T(\alpha, s) = q$ and $T(\alpha', s) = q'$ for every $s \in B$, $B \in \mathcal{S}$ non-null, then $\alpha \preceq \alpha'$ given B if and only if $q \preceq q'$

- D4: For $A, B \in \mathcal{S}$, $A \leq B$ if and only if $\alpha_A \preceq \alpha_B$ or $q \preceq q'$ for all $\alpha_A, \alpha_B \in C$, $q, q' \in F$ such that $T(\alpha_A, s) = q$ for $s \in A$, $T(\alpha_A, s') = q'$ for $s' \notin A$, $T(\alpha_B, s) = q$ for $s \in B$, $T(\alpha_B, s') = q'$ for $s' \notin B$. Read \leq as “is less probable than”
- P4: For every $A, B \in \mathcal{S}$, $A \leq B$ or $B \leq A$
- P5: For some α, β , $\alpha \prec \beta$
- P6: Suppose $\alpha \not\preceq \beta$. Then for every γ there is a finite partition of S such that if α' agrees with α and β' agrees with β except on some element B of the partition, α' and β' being equal to γ on B , then $\alpha \not\preceq \beta'$ and $\alpha' \not\preceq \beta$
- D5: $\alpha \preceq q$ for $q \in F$ given B if and only if $\alpha \preceq \beta$ given B where $T(\beta, s) = q$ for all $s \in S$
- P7: If $\alpha \preceq T(\beta, s)$ given B for every $s \in B$, then $\alpha \preceq \beta$ given B
- P7': The proposition given by inverting every expression in D5 and P7

Our initial view of decision problems was that the consequences Ω are a set of things we know how to rank and choices C are the things we want to rank. This is not exactly Savage’s setup – he assumes a preference relation ranking “acts” C to begin with. Furthermore, Savage also introduces a set of states S and assumes that the set of acts corresponds to the set of all function $S \rightarrow \Omega$. Many decision problems might be able to be extended with states and the set of acts enriched so as to satisfy these requirements, but it is not obvious that this is always possible.

D1 formalises the idea of one act α being not preferred to another β given the knowledge that the true state lies in the set B (in short: “given B ” or “conditional on B ”). P2 is sometimes called the “sure thing principle”, as it implies the following: for any α, β if α is better than β on some states and no worse on any other, then $\alpha \succ \beta$. In Savage’s model, the “likelihood” that of any state cannot depend on the act chosen.

D4 + P4 defines the “probability preorder” \leq on (S, \mathcal{S}) and assumes it is complete.

P5 is the requirement that the preference relation is non-trivial; not everything is equally desirable. This doesn’t seem like it should be a practical requirement to me; we might hope that a model can distinguish between some of our options, but that doesn’t mean we should assume it can. Savage claims that this requirement is “innocuous” because any exception must be trivial, but I’m not sure I agree.

P6 is a requirement of continuity; for any $\alpha \preceq \beta$, we can divide S finely enough to squeeze a “small slice” of any third outcome γ into the gap between the two.

P7 in combination with the other axioms forces preferences to be bounded.

3.2.3 Jeffrey's decision theory

Jeffrey's decision theory is an alternative to Savage's that starts from a different set of assumptions. One of the key differences is in what is assumed at the outset: where Savage assumes a set of states S , acts C and consequences Ω , Jeffrey's theory only considers a single space \mathcal{F} , which is a complete atomless boolean algebra. Elements of \mathcal{F} are said to be propositions, although the structure of \mathcal{F} means we can't understand it as, for example, a set of propositions regarding the result of a particular measurement procedure (Section 3.3). The theory is set out in Jeffrey (1965), and the key representation theorem proved in Bolker (1966).

Recall that our fundamental problem is relating a set C of things we can choose to a set F of things we can compare. Jeffrey's theory uses a different strategy to accomplish this than Savage's; where identifies a set of acts C with all functions $S \rightarrow F$ and proposes axioms that constrain a preference relation on C , Jeffrey assumes that choices are elements of the algebra \mathcal{F} , along with propositions that do not correspond to choices. Jeffrey's axioms pertain to a preference relation on \mathcal{F} . The ultimate result is, for our purposes, very similar.

Theorem 3.2.4. *Suppose there is a complete atomless Boolean algebra \mathcal{F} with a preference relation \preceq . If \preceq satisfies the Bolker axioms (Section 3.2.3) then there exists a desirability function $\text{des} : \mathcal{F} \rightarrow \mathbb{R}$ and a probability distribution $\mu \in \Delta(\mathcal{F})$ such that for $A, B \in \mathcal{F}$ and finite partition $D_1, \dots, D_n \in \mathcal{F}$:*

$$(A \preceq B) \iff \sum_i^n \text{des}(D_i) \mu(D_i|A) \leq \sum_i^n \text{des}(D_i) \mu(D_i|B) \quad (3.6)$$

where $\mu(D_i|A) := \frac{\mu(A \cap D_i)}{\mu(A)}$ for $\mu(A) > 0$, undefined otherwise.

Proof. Bolker (1966) □

As mentioned, in Jeffrey's theory the *choices* C are a subset of \mathcal{F} . Thus we can deduce from a Jeffrey model a function $C \rightarrow \Delta(\mathcal{F})$ that “represents the consequences of choices” in the sense of Theorem 3.2.5.

Theorem 3.2.5. *Suppose there is a complete atomless Boolean algebra \mathcal{F} with a preference relation \preceq that satisfies the Bolker axioms, and a set of choices C over which a preference relation is sought with $\mu(\alpha) > 0$ for all $\alpha \in C$. Then there is a function $\mathbb{P} : C \rightarrow \Delta(\mathcal{F})$ such that for any $\alpha, \alpha' \in C$ and finite partition $D_1, \dots, D_n \in \mathcal{F}$:*

$$\alpha \preceq \alpha' \iff \sum_i^n \text{des}(D_i) \mathbb{P}_\alpha(D_i) \leq \sum_i^n \text{des}(D_i) \mathbb{P}_{\alpha'}(D_i) \quad (3.7)$$

Where μ and des are as in Theorem 3.2.4

Proof. Define \mathbb{P} by $\alpha \mapsto \mu(\cdot|\alpha)$. Then Equation 3.7 follows from Equation 3.6. □

Bolker axioms

$\underline{\mathcal{F}}$ a complete, atomless Boolean algebra with the impossible proposition. An example of such a set is constructed from the set of Lebesgue measurable sets on $[0, 1]$ identifying any two sets that differ by a set of measure zero identified Bolker (1967). This is not a σ -algebra.

A1: \preceq is a complete preference relation

B2: $\underline{\mathcal{F}}$ is a complete, atomless Boolean algebra with the impossible proposition removed

C3: For $A, B \in \underline{\mathcal{F}}$, if $A \cap B = \emptyset$, then

- a) If $A \succ B$ then $A \succ A \cup B \succ B$
- b) If $A \sim B$ then $A \sim A \cup B \sim B$

D4: Given $A \cap B = \emptyset$ and $A \sim B$, if $A \cup G \sim B \cup G$ for some G where $A \cap G = B \cap G = \emptyset$ and $G \not\succeq A$, then $A \cup G \sim B \cup G$ for every such G

D1: The supremum (infimum) of a subset $W \subset \underline{\mathcal{F}}$ is a set G (D) such that for all $A \in W$, $G \subset A$ ($A \subset D$), and for any E that also has this property, $G \subset E$ ($E \subset D$)

E5: Given $W := \{W_i\}_{i \in M \subset \mathbb{N}}$ with $i < j \implies W_j \subset W_i$ and $W \subset \underline{\mathcal{F}}$ with supremum G (infimum D), whenever $A \prec G \prec B$ ($A \prec D \prec B$) then there exists some $k \in M$ such that $i \geq k$ ($i \leq k$) implies $A \prec W_i \prec B$.

Like Savage's theory, A1 requires the preference relation to be complete.

A3 is the assumption that the desirability of disjunctions of events lies between the desirability of each event; it is sometimes called "averaging". It notably rules out the following: if $A \succ B$ we cannot have $A \cup B \sim A$. In the Jeffrey-Bolker theory, propositions all have positive probabilities.

A4 allows a probability order to be defined on $\underline{\mathcal{F}}$. The conditions $A \cap B = \emptyset$, $A \sim B$, $A \cup G \sim B \cup G$ for some G where $A \cap G = B \cap G = \emptyset$ and $G \not\succeq A$ can be seen as a test for A and B being "equally probable". A4 requires that if A and B are rated as equally probable by one such test, then they are rated as equally probable by all such tests.

A5 is an axiom of continuity.

3.2.4 Causal decision theory

Causal decision theory was developed after both Jeffrey's and Savage's theory. A number of authors Lewis (1981); Skyrms (1982) felt that Jeffrey's theory erred by treating the consequences of a choice as an "ordinary conditional probability". Lewis (1981) suggested that causal decision theory can be used to evaluate choices when we are given a set Ω of consequences over which preferences are known, a set C of choices and a set H of dependency hypotheses (the letters

have been changed to match usage in this work; in the original the consequences were called S , the choices A and the dependency hypotheses H). Choices are then evaluated according to the causal decision rule. We have taken the liberty to state Lewis' rule in the language of the present work.

Definition 3.2.6 (Causal decision rule). Given a set C of choices, sample space (Ω, \mathcal{F}) , variables $H : \Omega \rightarrow H$ (the *dependency hypothesis*) and $S : \Omega \rightarrow S$ (the *consequence*) and a utility $u : \Omega \rightarrow \mathbb{R}$, the *causal utility* of a choice $\alpha \in C$ is given by

$$U(\alpha) := \int_S \int_H u(s) \mathbb{P}_\alpha^{S|H}(ds|h) \mathbb{P}_C^H(dh) \quad (3.8)$$

For some probabilistic function $\mathbb{P} : C \rightarrow \Delta(\Omega)$.

The reasons why Lewis wanted to introduce dependency hypothesis and modify Jeffrey's rule to Equation 3.8 are controversial and do not come up in this work. However, causal decision theory is still relevant to this work in two ways: firstly, once again is a probabilistic function $\mathbb{P} : C \rightarrow \Delta(\Omega)$. Secondly, causal decision theory introduces the notion of the dependency hypothesis H . The dependency hypothesis is similar to the state in Savage's theory, however Lewis does not require a deterministic map from dependency hypotheses to consequences, nor does he require a choice to correspond to every possible function from dependency hypotheses to states.

Dependency hypotheses are quite an important idea in causal reasoning. Together Lewis' decision rule connect the theory of probability sets with *statistical decision theory*, as Section 3.2.5 will show. Chapter 4 goes into considerable detail concerning the question of when probability sets support certain types of dependency hypothesis. While they are typically not explicitly represented in common frameworks for causal inference, Chapter ?? discusses how dependency hypotheses are often implicit in these approaches, and shows how they can be made explicit.

3.2.5 Statistical decision theory

Statistical decision theory (SDT), created by Wald (1950), predates all of the decision theories discussed above. Savage's theory appears to have developed in part to explain some features of SDT Savage (1951), and Jeffrey's theory and subsequent causal decision theories were in turn influenced by Savage's decision theory. While the later decision theories were concerned with articulating why their theory fit the role of a theory for rational decision under uncertainty, Wald focused much more on the mathematical formalism and solutions to statistical problems. Statistical decision theory introduced many fundamental ideas that have since entered the "water supply" of machine learning theory, such as *decision rules* and *risk* as a measure of the quality of a decision rule.

In contrast to the later decision theories, SDT has no explicit representation of the "consequences" of a decision. Rather, it is assumed that a loss function

is given that maps decisions and hypotheses directly to a loss, which is a kind of desirability score similar to a utility (although it is minimised rather than maximised).

Definition 3.2.7 (Statistical decision problem). A statistical decision problem is a tuple (X, H, D, l, \mathbb{P}) where (X, \mathcal{X}) is a set of outcomes, (H, \mathcal{H}) is a set of hypotheses, (D, \mathcal{D}) is a set of decisions, $l : D \times H \rightarrow \mathbb{R}$ is a loss function and $\mathbb{P} : H \rightarrow \mathcal{X}$ is a Markov kernel from hypotheses to outcomes.

Statistical decision theory is concerned with the selection of *decision rules*, rather than the selection of decisions directly. A decision rule maps observations to decisions, and may be deterministic or stochastic.

Definition 3.2.8 (Decision rule). Given a statistical decision problem (X, H, D, l, \mathbb{P}) , a decision rule is a Markov kernel $\mathbb{D}_\alpha : \Omega \rightarrow D$.

Because decision rules in SDT play the role of what we call *choices*, we denote the set of all available decision rules by C . A further feature of SDT that is unlike the later decision theories is that SDT does not offer a single rule for assessing the desirability of any choice in C . Instead, it offers a rule for assessing the desirability of a choice *relative to a particular hypothesis*, and considers a number of different rules for turning this “intermediate assessment” into a final assessment of the available choices.

Definition 3.2.9 (Risk). Given a statistical decision problem (X, H, D, l, \mathbb{P}) and decision functions C , the *risk functional* $R : C \times H \rightarrow \mathbb{R}$ is defined by

$$R(\mathbb{D}_\alpha, h) := \int_X \int_D l(d, h) \mathbb{D}_\alpha(\mathrm{d}d | f) \mathbb{P}_h(\mathrm{d}f) \quad (3.9)$$

SDT offers two decision rules that can be used to determine a best decision function when the hypothesis is unknown. The first is to choose the decision function minimising the *Bayes risk*, which requires a prior over hypotheses. The second decision rule is to choose the decision rule minimising the maximum risk, and this rule does not require a prior.

Definition 3.2.10 (Bayes risk). Given a statistical decision problem (X, H, D, l, \mathbb{P}) , decision functions C and a prior $\mu \in \Delta(\mathcal{H})$, the *Bayes risk* is

$$R_\mu(\mathbb{D}_\alpha) := \int_H R(\mathbb{D}_\alpha, h) \mu(\mathrm{d}h) \quad (3.10)$$

$$= \int_X \int_D l(d, h) \mathbb{D}_\alpha(\mathrm{d}d | f) \mathbb{P}_h(\mathrm{d}f) \mu(\mathrm{d}h) \quad (3.11)$$

Definition 3.2.11 (Minimax rule). Given a statistical decision problem (X, H, D, l, \mathbb{P}) and decision functions C , a *minimax decision function* is any decision function \mathbb{D}_α satisfying

$$\sup_{h \in H} R(\mathbb{D}_\alpha, h) = \inf_{\alpha' \text{ in } C} \sup_{h \in H} R(\mathbb{D}_{\alpha'}, h) \quad (3.12)$$

From consequences to statistical decision problems

To create a model *with consequences* that induces the SDT risk, we need to make a few assumptions. Recall that the basic model for a decision problem is either a Markov kernel $\mathbb{P} : C \rightarrow \Omega$ or $\mathbb{P}_{\cdot} : C \times H \rightarrow \Omega$ for some set of hypotheses H , and that a utility function u is also available. Here, Ω represents both observations (if any are taken) and consequences of decisions. Statistical decision problems correspond to a particular instance of this general scheme. In particular, statistical decision problems are associated with probability functions in which:

- The set of “choices” is identified with the set of decision rules
- There is some $X : \Omega \rightarrow X$ representing the observations made before any decision rule is enacted
 - Given any hypothesis, observations are independent of the choice of decision rule
- There are consequences $Y : \Omega \rightarrow Y$ and decisions $D : \Omega \rightarrow D$
 - Given a hypothesis and a decision, the consequences Y are independent of both the choice of decision rule and the observations
 - Given a decision function, the decision D is independent of the hypothesis
- The utility $u : Y \rightarrow \mathbb{R}$ is a function of the consequences

We call models satisfying these requirements *see-do models*

Definition 3.2.12 (See-do model). A probability set model of a statistical decision problem, or a *see-do model* for short, is a tuple $(\mathbb{P}_{C \times H}, X, Y, D)$ where $\mathbb{P}_{C \times H}$ is a probability set indexed by elements of $C \times H$ on (Ω, \mathcal{F}) , $X : \Omega \rightarrow X$ are the observations, $Y : \Omega \rightarrow Y$ are the consequences and $D : \Omega \rightarrow D$ are the decisions. $\mathbb{P}_{C \times H}$ must observe the following conditional independences:

$$X \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e C | H \quad (3.13)$$

$$Y \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e (X, C) | (D, H) \quad (3.14)$$

$$D \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e H | C \quad (3.15)$$

where $C : C \times H \rightarrow C$ and $H : C \times H \rightarrow H$ are the respective projections (see Definition 2.4.16 for the definition of extended conditional independence).

In order to produce a statistical decision problem, we also need a loss. We identify the required loss with the negative expected utility, conditional on a particular decision.

Definition 3.2.13 (Induced loss). Given a see-do model $(\mathbb{P}_{C \times H}, \mathbf{X}, \mathbf{Y}, \mathbf{D})$ and a utility $u : Y \rightarrow \mathbb{R}$, the induced loss $l : D \times H \rightarrow \mathbb{R}$ is defined as

$$l(d, h) := - \int_Y u(y) \mathbb{P}_{C \times \{h\}}^{\mathbf{Y}|\mathbf{D}}(dy|d) \quad (3.16)$$

where the uniform conditional $\mathbb{P}_{C \times \{h\}}^{\mathbf{Y}|\mathbf{D}}$'s existence is guaranteed by $\mathbf{Y} \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e (\mathbf{X}, \mathbf{C}) | (\mathbf{D}, \mathbf{H})$.

The probability set model induces a set of decision functions: for each $\alpha \in C$, there is a probability distribution $\mathbb{P}_\alpha^{\mathbf{D}|\mathbf{X}}$.

Theorem 3.2.14 (Induced risk). *Given a see-do model $(\mathbb{P}_{C \times H}, \mathbf{X}, \mathbf{Y}, \mathbf{D})$ along with a utility $u : Y \rightarrow \mathbb{R}$, the expected utility for each choice of decision rule $\alpha \in C$ and hypothesis $h \in H$ pair is equal to the negative risk of the decision rule and hypothesis pair.*

Proof. The expected utility given α and h is

$$\int_Y u(y) \mathbb{P}_{\alpha, h}^{\mathbf{Y}}(dy) = \int_Y \int_D \int_X u(y) \mathbb{P}_{\alpha, h}^{\mathbf{Y}|\mathbf{D}\mathbf{X}}(dy|d, x) \mathbb{P}_{\alpha, h}^{\mathbf{D}|\mathbf{X}}(dd|x) \mathbb{P}_{\alpha, h}^{\mathbf{X}}(dx) \quad (3.17)$$

$$= \int_X \int_D \int_Y u(y) \mathbb{P}_{\alpha, h}^{\mathbf{Y}|\mathbf{D}}(dy|d) \mathbb{P}_{\alpha, h}^{\mathbf{D}|\mathbf{X}}(dd|x) \mathbb{P}_{\alpha, h}^{\mathbf{X}}(dx) \quad (3.18)$$

$$= \int_X \int_D \int_Y u(y) \mathbb{P}_{C \times \{h\}}^{\mathbf{Y}|\mathbf{D}}(dy|d) \mathbb{P}_{\{\alpha\} \times H}^{\mathbf{D}|\mathbf{X}}(dd|x) \mathbb{P}_{C \times \{h\}}^{\mathbf{X}}(dx) \quad (3.19)$$

$$= - \int_D \int_X l(d, h) \mathbb{P}_{\{\alpha\} \times H}^{\mathbf{D}|\mathbf{X}}(dd|x) \mathbb{P}_{C \times \{h\}}^{\mathbf{X}}(dx) \quad (3.20)$$

$$= -R(\mathbb{P}_{\{\alpha\} \times H}^{\mathbf{D}|\mathbf{X}}, h) \quad (3.21)$$

where Equation 3.27 follows from $\mathbf{Y} \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e (\mathbf{X}, \mathbf{C}) | (\mathbf{D}, \mathbf{H})$, the uniform conditional $\mathbb{P}_{\{\alpha\} \times H}^{\mathbf{D}|\mathbf{X}}$ exists due to $\mathbf{D} \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e \mathbf{H} | \mathbf{C}$ and the uniform conditional $\mathbb{P}_{C \times \{h\}}^{\mathbf{X}}$ exists due to $\mathbf{X} \perp\!\!\!\perp_{\mathbb{P}_{C \times H}}^e \mathbf{C} | \mathbf{H}$. \square

We can also define a “Bayesian probability set model”, which is a probability set model \mathbb{P}_C along with a prior $\mathbb{P}_C^{\mathbf{H}}$ over hypotheses. In this case, the expected utility of a decision rule is equal to the Bayes risk.

Definition 3.2.15 (Bayesian see-do model). A Bayesian see-do model is a tuple $(\mathbb{P}_C, \mathbf{X}, \mathbf{Y}, \mathbf{D}, \mathbf{H})$ where \mathbb{P}_C is a probability set on (Ω, \mathcal{F}) , $\mathbf{X} : \Omega \rightarrow X$ are the observations, $\mathbf{Y} : \Omega \rightarrow Y$ are the consequences, $\mathbf{D} : \Omega \rightarrow D$ are the decisions and $\mathbf{H} : \Omega \rightarrow H$ is the hypothesis. \mathbb{P}_C must observe the following conditional independences:

$$\mathbf{X} \perp\!\!\!\perp_{\mathbb{P}_C}^e \mathbf{C} | \mathbf{H} \quad (3.22)$$

$$\mathbf{Y} \perp\!\!\!\perp_{\mathbb{P}_C}^e (\mathbf{X}, \mathbf{C}) | (\mathbf{D}, \mathbf{H}) \quad (3.23)$$

$$\mathbf{D} \perp\!\!\!\perp_{\mathbb{P}_C}^e \mathbf{H} | \mathbf{C} \quad (3.24)$$

$$\mathbf{H} \perp\!\!\!\perp_{\mathbb{P}_C}^e \mathbf{C} \quad (3.25)$$

Theorem 3.2.16 (Induced Bayes risk). *Given a Bayesian see-do model $(\mathbb{P}_C, \mathbf{X}, \mathbf{Y}, \mathbf{D}, \mathbf{H})$ of a statistical decision problem along with a utility $u : Y \rightarrow \mathbb{R}$, the expected utility for each choice of decision rule $\alpha \in C$ is equal to the negative Bayes risk of that decision rule relative to the prior \mathbb{P}_C^H .*

Proof. The expected utility of $\alpha \in C$ is

$$\int_Y u(y) \mathbb{P}_\alpha^Y(dy) = \int_Y \int_D \int_X \int_H u(y) \mathbb{P}_\alpha^{Y|D \times H}(dy|d, x, h) \mathbb{P}_\alpha^{D|X \times H}(dd|x, h) \mathbb{P}_\alpha^{X|H}(dx|h) \mathbb{P}_\alpha^H(dh) \quad (3.26)$$

$$= \int_X \int_D \int_Y \int_H u(y) \mathbb{P}_\alpha^{Y|D \times H}(dy|d, h) \mathbb{P}_\alpha^{D|X}(dd|x) \mathbb{P}_\alpha^{X|H}(dx|h) \mathbb{P}_\alpha^H(dh) \quad (3.27)$$

$$= \int_X \int_D \int_Y \int_H u(y) \mathbb{P}_C^{Y|D \times H}(dy|d, h) \mathbb{P}_\alpha^{D|X}(dd|x) \mathbb{P}_C^{X|H}(dx|h) \mathbb{P}_C^H(dh) \quad (3.28)$$

$$= - \int_D \int_X \int_H l(d, h) \mathbb{P}_\alpha^{D|X}(dd|x) \mathbb{P}_C^{X|H}(dx|h) \mathbb{P}_C^H(dh) \quad (3.29)$$

$$= - \int_H R(\mathbb{P}_\alpha^{D|X}, h) \mathbb{P}_C^H(dh) \quad (3.30) \quad = -R_{\mathbb{P}_C^H}(\mathbb{P}_\alpha^{D|X})$$

□

Theorem 3.2.17 (Extension to a Bayesian see-do model). *Given a see-do model $(\mathbb{P}_{C \times H}, \mathbf{X}, \mathbf{Y}, \mathbf{D})$*

Complete class theorem

A Bayesian see-do model is a probability set \mathbb{P}_C and a see-do model is a probability set $\mathbb{P}_{C \times H}$. We can obtain the former by combining the latter with a prior $\mu \in \Delta(H)$; call this an induced Bayesian model with respect to μ . The *complete class theorem* establishes that any admissible decision rule for a see-do model $\mathbb{P}_{C \times H}$ must minimise the Bayes risk for the Bayesian model induced by some $\mu \in \Delta(H)$.

3.3 Variables

In probability theory, it is standard to assume the existence of a probability space $(\mu, \Omega, \mathcal{F})$ and to define *random variables* as measurable functions from (Ω, \mathcal{F}) to $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. However, variables aren't *just* functions – they're also typically understood to correspond to some measured aspect of the real world. For example, Pearl (2009) offers the following two, purportedly equivalent, definitions of variables:

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

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

However, these quotes are describing different things – in fact, they’re not even describing the same *kind* of thing. The first is talking about a *measurement*, which is something we can do in the real world that produces as a result an element of a mathematical set. The second is talking about a *function*, which is a purely mathematical object with a domain and a codomain and a mapping from the former into the latter.

The way we address this distinction is: a procedure that takes place in the real world and yields as results elements of a mathematical set is called a *measurement procedure*. We suppose for a given problem that there is a “complete measurement procedure” \mathcal{S} , and the result of every specific measurement procedure of interest can be reconstructed from the result of the complete measurement procedure by applying a function to the latter result. The function X that yields the result of a specific measurement procedure \mathcal{X} given the result of the complete measurement procedure \mathcal{S} is the *variable* associated with the measurement procedure \mathcal{X} .

In this way, the variable X – which is by itself just a mathematical function – is made relevant to the real-world by combining it with a total measurement procedure \mathcal{S} .

We can even use this scheme to address situations where it is possible to make different choices. We can simply posit a sub-procedure \mathcal{C} that yields the choice α that we eventually make. However, modelling this can be tricky. If we want to use the consequence model to help make the decision, then it seems that the model of the decision procedure \mathcal{C} will need to be self-referential. Furthermore, even if we have a model of \mathcal{C} that says we will certainly decide on a particular element α^* , we still need to map every element of C to a consequence because this is what enables the comparison of elements of C . Thus, modelling \mathcal{C} makes the model more complicated, and it’s not obvious that this is answering the question that we are interested in. This complication is not obviously intractable, but we do not address it here. Instead, we simply assume that we have a collection \mathcal{S}_α of total measurement procedures that all yield elements of the same set, and each of which are executed if the choice made is α .

3.3.1 Variables and measurement procedures

We illustrate this approach with the example of Newton’s second law in the form $F = MA$. This model relates “variables” F , M and A . As Feynman (1979) noted, in order to understand this law, we must bring some pre-existing understanding of force, mass and acceleration independent of the law itself. Furthermore, we contend, this knowledge cannot be expressed in any purely mathematical statement. In order to say what the net force on a given object is, even a highly knowledgeable physicist will have to go and do some measurements, which is a procedure that they carry out involving interacting with the real world somehow and obtaining as a result a vector representing the net forces on that object.

That is, the variables F , M and A are referring to the *results of measurement procedures*. We will introduce a separate notation to refer to these measurement procedures – \mathcal{F} is the procedure for measuring force, \mathcal{M} and \mathcal{A} for mass and acceleration respectively. A measurement procedure \mathcal{F} is akin to Menger (2003)’s notion of variables as “consistent classes of quantities” that consist of pairing between real-world objects and quantities of some type. Force \mathcal{F} itself is not a well-defined mathematical thing, as measurement procedures are not mathematically well-defined. At the same time, the set of values it may yield *are* well-defined mathematical things. No actual procedure can be guaranteed to return elements of a mathematical set known in advance – anything can fail – but we assume that we can study procedures reliable enough that we don’t lose much by making this assumption.

Note that, because \mathcal{F} is not a purely mathematical thing, we cannot perform mathematical reasoning with \mathcal{F} directly. Rather, we introduce a variable F which, as we will see, is a well-defined mathematical object, assert that it corresponds to \mathcal{F} and conduct our reasoning using F .

3.3.2 Measurement procedures

Definition 3.3.1 (Measurement procedure). A *measurement procedure* \mathcal{B} is a procedure that involves interacting with the real world somehow and delivering an element of a mathematical set X as a result. A procedure \mathcal{B} is said to takes values in a set B .

We adopt the convention that the procedure name \mathcal{B} and the set of values B share the same letter.

Definition 3.3.2 (Values yielded by procedures). $\mathcal{B} \bowtie x$ is the proposition that the procedure \mathcal{B} will yield the value $x \in X$. $\mathcal{B} \bowtie A$ for $A \subset X$ is the proposition $\bigvee_{x \in A} \mathcal{B} \bowtie x$.

Definition 3.3.3 (Equivalence of procedures). Two procedures \mathcal{B} and \mathcal{C} are equal if they both take values in X and $\mathcal{B} \bowtie x \iff \mathcal{C} \bowtie x$ for all $x \in X$.

If two involve different measurement actions in the real world but necessarily yield the same result, we say they are equivalent.

It is worth noting that this notion of equivalence identifies procedures with different real-world actions. For example, “measure the force” and “measure everything, then discard everything but the force” are often different – in particular, it might be possible to measure the force only before one has measured everything else. Thus the result yielded by the first procedure could be available before the result of the second. However, if the first is carried out in the course of carrying out the second, they both yield the same result in the end and so we treat them as equivalent.

Measurement procedures are like functions without well-defined domains. Just like we can compose functions with other functions to create new functions, we can compose measurement procedures with functions to produce new measurement procedures.

Definition 3.3.4 (Composition of functions with procedures). Given a procedure \mathcal{B} that takes values in some set B , and a function $f : B \rightarrow C$, define the “composition” $f \circ \mathcal{B}$ to be any procedure \mathcal{C} that yields $f(x)$ whenever \mathcal{B} yields x . We can construct such a procedure by describing the steps: first, do \mathcal{B} and secondly, apply f to the value yielded by \mathcal{B} .

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

$$(g \circ f) \circ \mathcal{B} \text{ yields } x \iff \mathcal{B} \text{ yields } (g \circ f)^{-1}(x) \quad (3.31)$$

$$\iff \mathcal{B} \text{ yields } f^{-1}(g^{-1}(x)) \quad (3.32)$$

$$\iff f \circ \mathcal{B} \text{ yields } g^{-1}(x) \quad (3.33)$$

$$\iff g \circ (f \circ \mathcal{B}) \text{ yields } x \quad (3.34)$$

One might wonder whether there is also some kind of “tensor product” operation that takes a standalone \mathcal{M} and a standalone \mathcal{A} and returns a procedure $(\mathcal{M}, \mathcal{A})$. Unlike function composition, this would be an operation that acts on two procedures rather than a procedure and a function. Thus this “append” combines real-world operations somehow, which might introduce additional requirements (we can’t just measure mass and acceleration; we need to measure the mass and acceleration of the same object at the same time), and may be under-specified. For example, measuring a subatomic particle’s position and momentum can be done separately, but if we wish to combine the two procedures then we can get different results depending on the order in which we combine them.

Our approach here is to suppose that there is some complete measurement procedure \mathcal{S} to be modeled, which takes values in the observable sample space (Ψ, \mathcal{E}) and for all measurement procedures of interest there is some f such that the procedure is equivalent to $f \circ \mathcal{S}$ for some f . In this manner, we assume that any problems that arise from a need to combine real world actions have already been solved in the course of defining \mathcal{S} .

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

One could also consider measurement procedures to produce results in $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ (i.e. the reals with the Borel sigma-algebra) or a set isomorphic to it. This choice is often made in practice, and following standard practice we also often consider variables to take values in sets isomorphic to $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. However, for measurement in particular this seems to be a choice of convenience rather than necessity – for any measurement with finite precision and range, it is possible to specify a finite set of possible results.

3.3.3 Observable variables

Our *complete* procedure \mathcal{S} represents a large collection of subprocedures of interest, each of which can be obtained by composition of some function with \mathcal{S} . We call the pair consisting of a subprocedure of interest \mathcal{X} along with the variable X used to obtain it from \mathcal{S} an *observable variable*.

Definition 3.3.5 (Observable variable). Given a measurement procedure \mathcal{S} taking values in (Ψ, \mathcal{E}) , an observable variable is a pair $(X \circ \mathcal{S}, X)$ where $X : (\Psi, \mathcal{E}) \rightarrow (X, \mathcal{X})$ is a measurable function and $\mathcal{X} := X \circ \mathcal{S}$ is the measurement procedure induced by X and \mathcal{S} .

For the model $F = MA$, for example, suppose we have a complete measurement procedure \mathcal{S} that yields a triple (force, mass, acceleration) taking values in the sets X, Y, Z respectively. Then we can define the “force” variable (\mathcal{F}, F) where $\mathcal{F} := F \circ \mathcal{S}$ and $F : X \times Y \times Z \rightarrow X$ is the projection function onto X .

A measurement procedure yields a particular value when it is completed. We will call a proposition of the form “ \mathcal{X} yields x ” an *observation*. Note that \mathcal{X} need not be a complete procedure here. Given the complete procedure \mathcal{S} , a variable $X : \Psi \rightarrow X$ and the corresponding procedure $\mathcal{X} = X \circ \mathcal{S}$, the proposition “ \mathcal{X} yields x ” is equivalent to the proposition “ \mathcal{S} yields a value in $X^{-1}(x)$ ”. Because of this, we define the *event* $X \bowtie x$ to be the set $X^{-1}(x)$.

Definition 3.3.6 (Event). Given the complete procedure \mathcal{S} taking values in Ψ and an observable variable $(X \circ \mathcal{S}, X)$ for $X : \Psi \rightarrow X$, the *event* $X \bowtie x$ is the set $X^{-1}(x)$ for any $x \in X$.

If we are given an observation “ \mathcal{X} yields x ”, then the corresponding event $X \bowtie x$ is *compatible with this observation*.

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

An *impossible event* is the empty set. If $X \bowtie x = \emptyset$ this means that we have identified no possible outcomes of the measurement process \mathcal{S} compatible with the observation “ X yields x ”.

3.3.4 Model variables

Observable variables are special in the sense that they are tied to a particular measurement procedure \mathcal{S} . However, the measurement procedure \mathcal{S} does not enter into our mathematical reasoning; it guides our construction of a mathematical model, but once this is done mathematical reasoning proceeds entirely with mathematical objects like sets and functions, with no further reference to the measurement procedure.

A *model variable* is simply a measurable function with domain (Ψ, \mathcal{E}) .

Model variables do not have to be derived from observable variables. We may instead choose a sample space for our model (Ω, \mathcal{F}) that does not correspond to the possible values that \mathcal{S} might yield. In that case, we require a surjective model variable $S : \Omega \rightarrow \Psi$ called the complete observable variable, and every observable variable $(X' \circ \mathcal{S}, X')$ is associated with the model variable $X := X' \circ S$.

An *unobserved variable* is a variable whose set of possible values is not constrained by the results of the measurement procedure.

Definition 3.3.7 (Unobserved variable). Given a sample space (Ω, \mathcal{F}) and a complete observable variable $S : \Omega \rightarrow \Psi$, a model variable $Y : \Omega \rightarrow Y$ is *unobserved* if $Y(S \bowtie s) = Y$ for all $s \in \Psi$.

3.3.5 Variable sequences and partial order

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

Define the partial order on variables \preceq where $X \preceq Y$ can be read “ X is completely determined by Y ”.

Definition 3.3.8 (Variables determined by another variable). Given a sample space (Ω, \mathcal{F}) and variables $X : \Omega \rightarrow X$, $Y : \Omega \rightarrow Y$, $X \preceq Y$ if there is some $f : Y \rightarrow X$ such that $X = f \circ Y$.

Clearly, $X \preceq (X, Y)$ for any X and Y .

3.3.6 Decision procedures

The kind of problem we want to solve requires us to compare the consequences of different choices from a set of possibilities C . We take the *consequences of* $\alpha \in C$ to refer to the values obtained by some measurement procedure \mathcal{S}_α associated with the choice α .

As we have said, what exactly a “measurement procedure” is is a bit vague – it’s “what we actually do to get the numbers we associate with variables”. It seems we could describe the above in terms of a single measurement procedure \mathcal{S} , which involves:

1. Choose α
2. Proceed according to \mathcal{S}_α

However, \mathcal{S} is problematic to model. The model is often part of the process of choosing α , and so a model of \mathcal{S} that involves the step “choose α ” will be self-referential. Because of this, we don’t try to model \mathcal{S} , and whether this changes anything is an open question.

Definition 3.3.9 (Decision procedure). A decision procedure is a collection $\{\mathcal{S}_\alpha\}_{\alpha \in C}$ of measurement procedures.

Like measurement procedures, a decision procedure $\{\mathcal{S}_\alpha\}_{\alpha \in A}$ isn’t a well-defined mathematical object; it’s not really a “set”, because the contents are real-world actions.

Chapter 4

Repeatable Response Functions

4.1 When do response functions exist?

We model decision problems with probability sets \mathbb{P}_C for some set of choices C . If we have a pair of variables X and Y such that the uniform conditional $\mathbb{P}_C^{Y|X}$ exists (Definition 2.4.3), then the joint outcome \mathbb{P}_α^{XY} of any choice $\alpha \in C$ can be computed from the marginal distribution \mathbb{P}_α^X alone.

We're interested in models that feature a particular kind of “causal effect” that we call a *conditionally independent and identical response functions*, or just “response functions” for short. They are a causal analogue of conditionally independent and identical sequences of random variables. Concretely, a model with response functions is a probability set \mathbb{P}_C with variables $Y := (Y_i)_{i \in M}$, $(X_i)_{i \in M}$ for some index set M and some H such that $Y \perp\!\!\!\perp_{\mathbb{P}_C}^e C|X_i H$ and $H \perp\!\!\!\perp_{\mathbb{P}_C}^e X_i C$ (see Section 2.4.2 for extended conditional independence) and $\mathbb{P}_C^{Y_i|X_i H} = \mathbb{P}_C^{Y_j|X_j H}$ for all $i, j \in M$ (the identical response conditional requirement).

is that OK?

We will focus on the case where $\mathbb{P}_\alpha^{H|Y_A X_A}$ approaches a deterministic distribution as $|A| \rightarrow \infty$, for appropriate $\alpha \in C$. We could say this is the case of “identifiable” response conditionals.

Put together, these conditions say: in the limit of infinite samples under an appropriate sampling regime $\alpha \in C$, the model converges to a probabilistic function $X \rightarrow Y$ that represents “the probability of Y_i given X_i ” for any unobserved (X_i, Y_i) . For arbitrary $\alpha' \in C$, we may instead converge to a set containing this limiting distribution. We think – although it usually isn't stated in these terms – that given a causal Bayesian network, if the function “ $x \mapsto \mathbb{P}(Y|\text{do}(X = x))$ ” is identifiable, then it is an instance of the kind of function that we described in the previous sentences.

We prove our result under the simplifying assumption that $X_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{<i} C|X_{<i}$. This is a limiting assumption – for example, it excludes cases where X_i depends

on $(X_{<i}, Y_{<i})$. In the more general case where this does not hold, the conditions we provide must still hold for any $C' \subset C$ such that $X_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{<i>C'} | X_{<i}$, but providing sufficient conditions in this case is the topic of further work.

Under this assumption, we show that *causal contractibility* is necessary and sufficient for the existence of response conditionals. Causal contractibility can be broken down into two sub-assumptions: *exchange commutativity* and *consequence locality*. The first is the assumption that the uniform conditional probability $\mathbb{P}_C^{Y|X}$ “commutes” with the permutation operation, and the second is the assumption that X_i “has no effect” on any Y_j for $j \neq i$.

4.1.1 Relevance to previous work

Both sub-assumptions have precedent in existing literature, but these precedents tend to have been stated at somewhat informally.

Post-treatment exchangeability found in Dawid (2020) is implied by exchange commutativity, but not the reverse. “Causal exchangeability” notions are also found in Greenland and Robins (1986) and Banerjee et al. (2017); a subtle difference between these notions and exchange commutativity is that these latter notions are given as symmetries of *decision procedures* – they involve actually swapping actions taken or individuals in an experiment – while exchange commutativity is a symmetry of probability sets.

Consequence locality is similar to the stable unit treatment distribution assumption (SUTDA) in Dawid (2020), although consequence locality is distinguished by being a concrete extended conditional independence (Definition 4.1.1) while SUTDA is given as the assumption that Y_i “depends only on” X_i . Consequence locality is also similar to stable unit treatment value assumption (SUTVA). The stable unit treatment value assumption (SUTVA) is given as (Rubin, 2005):

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

String diagram statements of both sub-assumptions can give some intuition about what they mean. For clarity, these diagrams illustrate the assumptions with exactly two inputs and two outputs, while the general definitions are for any countable number of inputs and outputs.

Exchange commutativity for two inputs and outputs is given by the following equality:

$$\begin{array}{c} D_1 \\ D_2 \end{array} \begin{array}{c} \diagup \\ \diagdown \end{array} \boxed{\mathbb{P}_C^{Y_{\{1,2\}} | D_{\{1,2\}}}} \begin{array}{c} \diagdown \\ \diagup \end{array} \begin{array}{c} Y_1 \\ Y_2 \end{array} = \begin{array}{c} D_1 \\ D_2 \end{array} \boxed{\mathbb{P}_C^{Y_{\{1,2\}} | D_{\{1,2\}}}} \begin{array}{c} \diagdown \\ \diagup \end{array} \begin{array}{c} Y_1 \\ Y_2 \end{array} \quad (4.1)$$

While consequence locality for two inputs and outputs is given by the following pair of equalities:

$$\begin{array}{c} X_1 \\ X_2 \end{array} \begin{array}{|c} \hline \mathbb{P}_C^{Y_{1,2}|X_{1,2}} \\ \hline \end{array} \begin{array}{c} Y_1 \\ * \end{array} = \begin{array}{c} X_1 \\ X_2 \end{array} \begin{array}{|c} \hline \mathbb{P}_C^{Y_1|X_1} \\ \hline \end{array} \begin{array}{c} Y_1 \\ * \end{array} \quad (4.2)$$

$$\begin{array}{c} X_1 \\ X_2 \end{array} \begin{array}{|c} \hline \mathbb{P}_C^{Y_{1,2}|X_{1,2}} \\ \hline \end{array} \begin{array}{c} * \\ Y_2 \end{array} = \begin{array}{c} X_1 \\ X_2 \end{array} \begin{array}{|c} \hline \mathbb{P}_C^{Y_2|X_2} \\ \hline \end{array} \begin{array}{c} * \\ Y_2 \end{array} \quad (4.3)$$

4.1.2 Causal contractibility

Here we set out formal definitions of exchange commutativity and locality of consequences, as well as “consequence contractibility”, which is the conjunction of both conditions.

Definition 4.1.1 (Locality of consequences). Suppose we have a sample space (Ω, \mathcal{F}) and a probability set \mathbb{P}_C where $\mathbf{Y} := \mathbf{Y} := (Y_i)_M$, $\mathbf{D} := \mathbf{D}_M := (D_i)_M$, $M \subseteq \mathbb{N}$. If for any $A \subset M$, $\mathbf{Y}_A \perp\!\!\!\perp_{\mathbb{P}_C}^e \mathbf{D}_{A^c} C | \mathbf{D}_A$ then \mathbb{P}_C exhibits $(\mathbf{D}; \mathbf{Y})$ -local consequences.

If \mathbb{P}_C exhibits $(\mathbf{D}; \mathbf{Y})$ -local consequences then, given two different choices α and α' such that $\mathbb{P}_\alpha^{\mathbf{D}_A} = \mathbb{P}_{\alpha'}^{\mathbf{D}_A}$ then $\mathbb{P}_\alpha^{\mathbf{Y}_A} = \mathbb{P}_{\alpha'}^{\mathbf{Y}_A}$. However, \mathbb{P}_C may exhibit consequence locality even if no such pair of choices exists.

Note that consequence locality implies $\mathbf{Y}_M \perp\!\!\!\perp_{\mathbb{P}_C}^e C | \mathbf{D}_M$, and hence we have the uniform conditional $\mathbb{P}_C^{\mathbf{Y}_M | \mathbf{D}_M}$. We assume the existence of such a conditional for the next definition.

Definition 4.1.2 (Swap map). Given $M \subset \mathbb{N}$ a finite permutation $\rho : M \rightarrow M$ and a variable $\mathbf{X} : \Omega \rightarrow X^M$ such that $\mathbf{X} = (X_i)_{i \in M}$, define the Markov kernel $\text{swap}_{\rho(\mathbf{X})} : X^M \rightarrow X^M$ by $(d_i)_{i \in \mathbb{N}} \mapsto \delta_{(d_{\rho(i)})_{i \in \mathbb{N}}}$.

Definition 4.1.3 (Exchange commutativity). Suppose we have a sample space (Ω, \mathcal{F}) and a probability set \mathbb{P}_C with uniform conditional probability $\mathbb{P}_C^{\mathbf{Y} | \mathbf{D}}$ where $\mathbf{Y} := \mathbf{Y} := (Y_i)_M$, $\mathbf{D} := \mathbf{D}_M := (D_i)_M$, $M \subseteq \mathbb{N}$. If for any finite permutation $\rho : M \rightarrow M$

$$\text{swap}_{\rho(\mathbf{D})} \mathbb{P}_C^{\mathbf{Y} | \mathbf{D}} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{\mathbf{Y} | \mathbf{D}} \text{swap}_{\rho(\mathbf{Y})} \quad (4.4)$$

Then $\mathbb{P}_C^{\mathbf{Y} | \mathbf{D}}$ is $(\mathbf{D}; \mathbf{Y})$ -exchange commutative.

If \mathbb{P}_C is $(\mathbf{D}; \mathbf{Y})$ -exchange commutative and we have $\alpha, \alpha' \in C$ such that $\mathbb{P}_\alpha^C = \mathbb{P}_{\alpha'}^C \text{swap}_{\rho(\mathbf{D})}$, then $\mathbb{P}_\alpha^{\mathbf{Y}} = \mathbb{P}_{\alpha'}^{\mathbf{Y}} \text{swap}_{\rho(\mathbf{Y})}$. However, \mathbb{P}_C may commute with exchange even if there are no such α and $\alpha' \in C$.

Theorem 4.1.4 shows that neither condition implies the other.

Theorem 4.1.4. *Exchange commutativity does not imply locality of consequences or vice versa.*

Proof. Appendix ??.

□

If we are modelling the treatment of several patients whom who have already been examined, we might assume consequence locality – patient B’s treatment does not affect patient A – but not exchange commutativity – we don’t expect the same results from giving patient A’s treatment to patient B as we would from giving patient A’s treatment to patient A.

A model of stimulus payments might exhibit exchange commutativity but not consequence locality. If exactly n payments of \$10 000 are made, we might suppose that it doesn’t matter much exactly who receives the payments, but the amount of inflation induced depends on the number of payments made; making 100 such payments will have a negligible effect on inflation, while making payments to everyone in the country will have a substantial effect. Dawid (2000) offers the example of herd immunity in vaccination campaigns as a situation where post-treatment exchangeability holds but locality of consequences does not.

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

$$\mathbb{P}_S^{Y_1|D}(y_1|d_1, d_2) = 0.5 \quad (4.5)$$

$$\mathbb{P}_S^{Y_2|D}(y_2|d_1, d_2) = 0.5 \quad (4.6)$$

- The marginal distribution of both experiments in isolation is Bernoulli(0.5) no matter what choices I make, so a model of this experiment would satisfies Definition 4.1.1
- Nevertheless, the choice for the first experiment affects the result of the second experiment

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

Definition 4.1.5 (Causal contractibility). A probability set \mathbb{P}_C is $(D; Y)$ -*causally contractible* if it is both exchange commutative and exhibits consequence locality.

Theorem 4.1.6 (Equality of reduced conditionals). *A probability set \mathbb{P}_C that is $(D; Y)$ -causally contractible has, for any $A, B \subset M$ with $|A| = |B|$*

$$\mathbb{P}_C^{Y_A|D_A} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{Y_B|D_B} \quad (4.7)$$

Proof. Only if: For any $A, B \subset M$, let $s_{BA} : D^M \rightarrow D^M$ be the swap map that sends the B indices to A indices and $s_{AB} : Y^M \rightarrow Y^M$ be the swap map that sends A indices to B indices.

$$\begin{array}{c} D_A \text{---} \boxed{\mathbb{P}_C^{Y_A|D_A}} \text{---} Y_A \\ D_{M \setminus A} \text{---} * \end{array} = D_{M \setminus A} \text{---} \boxed{\mathbb{P}_C^{Y_A Y_{M \setminus A} | D_A D_{M \setminus A}}} \text{---} Y_A \quad (4.8)$$

$$= D_{M \setminus A} \text{---} \boxed{s_{BA}} \text{---} \boxed{\mathbb{P}_C^{Y_A Y_{M \setminus A} | D_A D_{M \setminus A}}} \text{---} \boxed{s_{AB}} \text{---} Y_A \quad (4.9)$$

$$= D_{M \setminus B} \text{---} \boxed{\mathbb{P}_C^{Y_B Y_{M \setminus B} | D_A D_{M \setminus B}}} \text{---} Y_B \quad (4.10)$$

Thus

$$\mathbb{P}_C^{Y_A|D_A D_{M \setminus A}} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{Y_B|D_B D_{M \setminus B}} \quad (4.11)$$

$$\stackrel{\mathbb{P}_C}{\cong} \begin{array}{c} D_A \text{---} \boxed{\mathbb{P}_C^{Y_A|D_A}} \text{---} Y_A \\ D_{M \setminus A} \text{---} * \end{array} \quad (4.12)$$

$$\implies \mathbb{P}_C^{Y_A|D_A} \stackrel{\mathbb{P}_C}{\cong} \mathbb{P}_C^{Y_B|D_B} \quad (4.13)$$

□

4.1.3 Existence of response conditionals

The main result in this section is Theorem 4.1.8 which shows that a probability set \mathbb{P}_C is causally contractible if and only if it can be represented as the product of a distribution over hypotheses \mathbb{P}_H and a collection of identical uniform conditionals $\mathbb{P}_C^{Y_1|D_1 H}$. Note the hypothesis H that appears in this conditional; it can be given the interpretation of a random variable that expresses the “true but initially unknown” $Y_1|D_1$ conditional probability.

Theorem 4.1.7. *Given a probability set \mathbb{P}_C and variables $D := (D_i)_{i \in \mathbb{N}}$ and $Y := (Y_i)_{i \in \mathbb{N}}$, \mathbb{P}_C is $(D; Y)$ -causally contractible if and only if there exists a column exchangeable probability distribution $\mu^{Y^D} \in \Delta(Y^{D \times \mathbb{N}})$ such that*

$$\mathbb{P}_C^{Y|D} = \begin{array}{c} \triangle \mu^{Y^D} \\ D \text{---} \boxed{\mathbb{F}_{\text{ev}}} \text{---} Y \end{array} \quad (4.14)$$

$$\iff \quad (4.15)$$

$$\mathbb{P}_C^{Y|D}(y|(d_i)_{i \in \mathbb{N}}) = \mu^{Y^D} \Pi_{(d_i)_{i \in \mathbb{N}}}(y) \quad (4.16)$$

Where $\Pi_{(d_i i)_{i \in \mathbb{N}}} : Y^{|D| \times \mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is the function that projects the (d_i, i) indices for all $i \in \mathbb{N}$ and \mathbb{F}_{ev} is the Markov kernel associated with the evaluation map

$$ev : D^{\mathbb{N}} \times Y^{D \times \mathbb{N}} \rightarrow Y \quad (4.17)$$

$$((d_i)_{i \in \mathbb{N}}, (y_{ij})_{i \in D, j \in \mathbb{N}}) \mapsto (y_{d_i i})_{i \in \mathbb{N}} \quad (4.18)$$

Proof. Appendix ??.

□

We would prefer to talk about Y^D as a latent variable, rather than needing to refer to the factorisation of a model in terms of μ^{Y^D} in Equation 4.14. This motivates the definition of an *augmented* causally contractible model.

An augmented causally contractible model looks in some respects similar to a potential outcomes model - both have a distribution over an unobserved “tabular” variable Y^D , and the value of Y_i given D is deterministically equal to the Y_i^D (abusing notation). However, the Y^D in an augmented causally contractible model usually can’t be interpreted as potential outcomes. For example, consider a series of bets on fair coin flips. Model the consequence Y_i as uniform on $\{0, 1\}$ for any decision D_i , for all i . Specifically, $D = Y = \{0, 1\}$ and $\mathbb{P}_\alpha^{Y^n}(y) = \prod_{i \in [n]} 0.5$ for all n , $y \in Y^n$, $\alpha \in R$. Then the construction of \mathbb{P}^{Y^D} following the method in Lemma 4.1.7 yields $\mathbb{P}^{Y_i^D}(y_i^D) = \prod_{j \in D} 0.5$ for all $y_i^D \in Y^D$. In this model Y_i^0 and Y_i^1 are independent and uniformly distributed. However, if we wanted Y_i^0 to be interpretable as “what would happen if I bet on outcome 0 on turn i ” and Y_i^1 to represent “what would happen if I bet on outcome 1 on turn i ”, then we ought to have $Y_i^0 = 1 - Y_i^1$.

The following is the main theorem of this section, that establishes the equivalence between causal contractibility and the existence of response conditionals. The argument in outline is: because $\mathbb{P}_C^{Y^D}$ is a column exchangeable probability distribution we can apply De Finetti’s theorem to show $\mathbb{P}_C^{Y^D}$ is representable as a product of identical parallel copies of $\mathbb{P}_C^{Y_1^D | H}$ and a common prior \mathbb{P}_C^H . This in turn can be used to show that $\mathbb{P}_C^{Y^D}$ can be represented as a product of identical parallel copies of $\mathbb{P}_C^{Y_1 | D_1 H}$ and the same common prior \mathbb{P}_C^H .

Theorem 4.1.8. *Suppose we have a sample space (Ω, \mathcal{F}) and a probability set \mathbb{P}_C and variables $D := (D_i)_{i \in \mathbb{N}}$ and $Y := (Y_i)_{i \in \mathbb{N}}$. Suppose also that \mathbb{P}_C is $(D; Y)$ -causally contractible if and only if there exists some $H : \Omega \rightarrow H$ such that*

\mathbb{P}_C^H and $\mathbb{P}_C^{Y_i|HD_i}$ exist for all $i \in \mathbb{N}$ and

$$\mathbb{P}_C^{Y|D} = \begin{array}{c} \triangle \mathbb{P}_C^H \\ \swarrow \\ \bullet \text{---} \boxed{\mathbb{P}_C^{Y_0|HD_0} \text{---} Y_i} \\ \uparrow D \\ \text{---} \end{array} \quad i \in \mathbb{N} \quad (4.19)$$

$$\iff \quad (4.20)$$

$$Y_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{N \setminus i}, D_{N \setminus i} C | HD_i \quad \forall i \in \mathbb{N} \quad (4.21)$$

$$\wedge H \perp\!\!\!\perp_{\mathbb{P}_C}^e DC \quad (4.22)$$

$$\wedge D_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{<i>} C | D_{<i>} \quad (4.23)$$

$$\wedge \mathbb{P}_C^{Y_i|HD_i} = \mathbb{P}^{Y_0|HD_0} \quad \forall i \in \mathbb{N} \quad (4.24)$$

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

Proof. Appendix ??.

□

4.1.4 Elaborations and examples

Theorem 4.1.8 requires an infinite sequence of causally contractible pairs. In practice we only want to model finite sequences of variables, but this theorem applies as long as it is possible to extend the finite model to an infinite model maintaining causal contractibility.

Theorem 4.1.8 applies whatever procedure we use to obtain the (D_i, Y_i) pairs – the D_i s may be randomised, passive observations or active choices. Purely passive observations can be modeled with a probability set of size 1, and in this case an exchangeable sequence of (D_i, Y_i) will also be causally contractible.

If we are modelling M passive observations followed by N active choices, then we will have a model \mathbb{P}_C with $D_{[M]} \perp\!\!\!\perp_{\mathbb{P}_C}^e C$ (because these are passive observations). If this model is $(D; Y)$ -causally contractible, then one consequence of this is an “observational imitation” condition: any choice α that makes $\mathbb{P}_\alpha^{D_{[M+N]}}$ exchangeable also makes $\mathbb{P}_\alpha^{Y_{[M+N]}}$ exchangeable. That is, if for some permutation swap_ρ

$$\mathbb{P}_\alpha^{D_{[M+N]}} \text{swap}_\rho = \mathbb{P}_\alpha^{D_{[M+N]}} \quad (4.25)$$

then by commutativity of exchange

$$\mathbb{P}_\alpha^{Y_{[M+N]}} = \mathbb{P}_\alpha^{D_{[M+N]}} \mathbb{P}_C^{Y_{[M+N]} | D_{[M+N]}} \quad (4.26)$$

$$= \mathbb{P}_\alpha^{D_{[M+N]}} \text{swap}_\rho \mathbb{P}_C^{Y_{[M+N]} | D_{[M+N]}} \quad (4.27)$$

$$= \mathbb{P}_\alpha^{D_{[M+N]}} \mathbb{P}_C^{Y_{[M+N]} | D_{[M+N]}} \text{swap}_\rho \quad (4.28)$$

$$= \mathbb{P}_\alpha^{Y_{[M+N]}} \text{swap}_\rho \quad (4.29)$$

If we assume a probability set \mathbb{P}_C is $(D, X; Y)$ -causally contractible and $X_i \perp\!\!\!\perp_{\mathbb{P}_C}^e D_i C | H$ – that is, D_i is must be independent of X_i conditional on H – then we get a version of the “backdoor adjustment” formula. Specifically

$$\mathbb{P}_\alpha^{Y_i | D_i H}(A | d, h) = \int_X \mathbb{P}_\alpha^{Y_i | X_i D_i H}(A | d, x, h) \mathbb{P}_\alpha^{X_i | D_i H}(dx | d, h) \quad (4.30)$$

$$= \int_X \mathbb{P}_C^{Y_i | X_i D_i H}(A | d, x, h) \mathbb{P}_C^{X_i | H}(dx | h) \quad (4.31)$$

If we additionally assume $\mathbb{P}_C^{X_i | H} \cong \mathbb{P}_C^{X_1 | H}$ then

$$\mathbb{P}_\alpha^{Y_i | D_i H}(A | d, h) = \int_X \mathbb{P}_C^{Y_1 | X_1 D_1 H}(A | d, x, h) \mathbb{P}_C^{X_1 | H}(dx | h) \quad (4.32)$$

Equation 4.32 is identical to the backdoor adjustment formula for an intervention on D_1 targeting Y_1 where X_1 is a common cause of both.

While it is formally possible to use a causally contractible model for a decision procedure that involves both passive observations and active choices, causal contractibility is a very strong assumption. Suppose we have a decision procedure in which M passive observations are made (D_M, Y_M) , followed by M active choices $(D_{(M, 2M]}, Y_{(M, 2M]})$. If a model \mathbb{P}_C of this procedure is (D_{2M}, Y_{2M}) -causally contractible model then the following holds (see corollary 4.3.5):

$$\mathbb{P}_C^{Y_{[2, M+1]} | D_{[2, M+1]}} = \mathbb{P}^{Y_{(M, 2M]} | D_{(M, 2M]}} \quad (4.33)$$

$$\implies \mathbb{P}_C^{Y_{M+1} | D_{[2, M+1]} Y_{[2, M]}} = \mathbb{P}^{Y_{M+1} | D_{(M, 2M]} Y_{(M+1, 2M]}} \quad (4.34)$$

That is, causal contractibility implies that there is no difference between conditioning on observational results or on the results of active choices; active choices are as good for predicting observations as vice-versa. Normally one might consider randomised experimental results to be “better” than passive observations, but this is not compatible with the assumption of causal contractibility.

4.1.5 Assessing causal contractibility

Assessing when a particular sequence of experiments should be modeled with a causally contractible model can be difficult. One way to justify the assumption is in two steps: first, all the repetitions of the experiment that yield the values of each of the (D_i, Y_i) pairs are indistinguishable “at the time of model construction”, and they are still indistinguishable after learning the value of D – because, for example, D is deterministic for each choice $\alpha \in C$.

Two step justifications of this form are common in literature on causal identifiability. For example, Greenland and Robins (1986) explain

Equivalence of response type may be thought of in terms of exchangeability of individuals: if the exposure states of the two individuals had been exchanged, the same data distribution would have resulted.

Note that exchanging individuals involved in an experiment and exchanging the individuals' exposure states are two different things, and the former doesn't imply the latter. We may consider a model that is symmetric to permutations of individual identifiers but is not symmetric to permuting individual identifiers and leaving exposure states fixed.

Dawid (2020) suggests (with many qualifications) that “post-treatment exchangeability” for a decision problem regarding taking aspirin to treat a headache may be acceptable if the data are from

A group of individuals whom I can regard, in an intuitive sense, as similar to myself, with headaches similar to my own.

Dawid points to the “first step” in our two step justification for causal contractibility: that the people involved are “similar” in an appropriate sense. However, under Dawid's approach there is a background assumption here that whether or not I take the aspirin is deterministic given the choice I end up making, which is the second step in our two step justification.

Finally, Rubin (2005) explicitly discusses two separate assumptions to justify causal identifiability:

indexing of the units is, by definition, a random permutation of $1, \dots, N$, and thus any distribution on the science must be row-exchangeable [...] The second critical fact is that if the treatment assignment mechanism is ignorable (e.g., randomized), then when the expression for the assignment mechanism (2) is evaluated at the observed data, it is free of dependence on Y_{mis}

Here we have a more abstract statement about the row-exchangeability of “the science”, rather than individual people involved in an experiment, but we regard it as similar in spirit to assumptions that people involved in the experiment are “similar”. Rubin explicitly mentions a second condition: that the treatment assignment is randomized.

Theorem 4.1.12 formalises these ideas. As an example of its application, consider an experiment where N patients, each with an individual identifier I_i , receive treatment D_i and experience outcome Y_i . We assume a $((D, I); Y)$ -causally contractible model \mathbb{P}_C is appropriate. This reflects two judgments; firstly, that treatment D_i and identifiers I_i screen off all other variables from Y_i (Definition 4.1.1), and secondly that the order in which the individuals appear and the treatments are received does not alter the consequences we expect to see (Definition 4.1.3). The fact that we need a preliminary assumption of causal contractibility is similar to how, in the potential outcomes framework, a preliminary assumption of SUTVA is required in order to justify the use of potential outcomes.

Next, we assume that, no matter which choice $\alpha \in C$ is decided on, all identifiers can be swapped without altering the distribution over consequences, and finally that for each choice $\alpha \in C$ the treatment vector D is deterministic. Then, according to Theorem 4.1.12, \mathbb{P}_C is also $(D; Y)$ -causally contractible. This can be extended to the case where D is a function of a “random signal” R .

Proof. For arbitrary $\nu \in \Delta(I^{\mathbb{N}})$, by assumption of $((\mathbf{D}, \mathbf{I}); \mathbf{Y})$ -causal contractibility and Theorem 4.1.8

$$\Rightarrow \mathbb{P}_C^{\mathbf{Y}|\mathbf{D}} = \begin{array}{c} \begin{array}{c} \text{Diagram of } \mathbb{P}_C^{\mathbf{Y}|\mathbf{D}} \text{ block} \end{array} \end{array} \quad (4.38)$$

Lemma 4.1.11. *Suppose we have a probability set \mathbb{P}_C where $Y \perp_{\mathbb{P}_C}^e C|(D, I)$ and I is an index variable. If for each permutation $\rho: \mathbb{N} \rightarrow \mathbb{N}$*

Proof. Taking $A \subset \mathbb{N}^{\mathbb{N}}$ to be the set of permutations of \mathbb{N} , note that for every $i \in A$, $B \in \mathcal{Y}^{\mathbb{N}}$, $d \in D^{\mathbb{N}}$ we can take $\rho_i : \mathbb{N} \rightarrow \mathbb{N}$ such that the image of i under

ρ is \mathbb{N} ; that is, $\rho_i(i) = \text{id}_{\mathbb{N}}$. Then

$$\mathbb{P}_C^{\mathbf{Y}|\mathbf{D}}(B|i, d) = (\text{swap}_{\rho_i(I)} \otimes \text{Id}_X) \mathbb{P}_C^{\mathbf{Y}|\mathbf{D}}(B|i, d) \quad (4.40)$$

$$= \mathbb{P}_C^{\mathbf{Y}|\mathbf{D}}(B|\text{id}_{\mathbb{N}}, d) \quad (4.41)$$

Therefore

$$\mathbb{P}_C^{\mathbf{Y}|\mathbf{D}} \stackrel{\mathbb{P}_C}{\cong} \text{erase}_{\mathbb{N}^{\mathbb{N}}} \otimes \mathbb{K} \quad (4.42)$$

where $\mathbb{K} : D^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is the kernel

$$(B|d) \mapsto \mathbb{P}_C^{\mathbf{Y}|\mathbf{D}}(B|\text{id}_{\mathbb{N}}, d) \quad (4.43)$$

□

Theorem 4.1.12. *Suppose we have a probability set \mathbb{P}_C , C countable, that $((D, \mathbf{l}); \mathbf{Y})$ -causally contractible for variables $\mathbf{Y} : \Omega \rightarrow Y^{\mathbb{N}}$, $D : \Omega \rightarrow D^{\mathbb{N}}$ and index variable $\mathbf{l} : \Omega \rightarrow \mathbb{N}^{\mathbb{N}}$. If for each $\alpha \in C$, $\rho : \mathbb{N} \rightarrow \mathbb{N}$*

$$\mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{l}} = \text{swap}_{\rho(I)} \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{l}} \quad (4.44)$$

and there is an invertible function $f : C \rightarrow D$ such that

$$\mathbb{P}_{\alpha}^{\mathbf{D}} = \mathbb{F}_f \quad (4.45)$$

then \mathbb{P}_C is $(D; \mathbf{Y})$ -causally contractible.

Proof. The map $\mathbb{Q} : (B|i, \alpha) \mapsto \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{l}}(B|i)$ is itself a Markov kernel. By assumption, for any α ,

$$\mathbb{Q} \stackrel{\mathbb{P}_{\alpha}}{\cong} (\text{id}_{\mathbb{N}^{\mathbb{N}}} \otimes \mathbb{F}_f) \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{D}} \quad (4.46)$$

Furthermore, by assumption

$$(\text{swap}_{\rho(I)} \otimes \text{Id}_C) \mathbb{Q} = \mathbb{Q} \quad (4.47)$$

Therefore

$$(\text{swap}_{\rho(I)} \otimes \text{Id}_C) \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{D}} = (\text{id}_{\mathbb{N}^{\mathbb{N}}} \otimes \mathbb{F}_f)(\text{id}_{\mathbb{N}^{\mathbb{N}}} \otimes \mathbb{F}_{f^{-1}})(\text{swap}_{\rho(I)} \otimes \text{Id}_C) \mathbb{Q} \quad (4.48)$$

$$= (\text{id}_{\mathbb{N}^{\mathbb{N}}} \otimes \mathbb{F}_{f^{-1}}) \mathbb{Q} \quad (4.49)$$

$$= \mathbb{P}_{\alpha}^{\mathbf{Y}|\mathbf{D}} \quad (4.50)$$

Then by Lemma 4.1.11 we have $\mathbf{Y} \perp_{\mathbb{P}_C}^e \mathbf{l} | \mathbf{D}$ and by Lemma 4.1.9 we have $(D; \mathbf{Y})$ -causal contractibility. □

If we suppose Theorem 4.1.12 applies to $C' \subset C$ such that D is deterministic for all $\alpha \in C'$, while C consists of C' and all “random choices between elements of C' ”; that is for all $\beta \in C$

$$\mathbb{P}_\beta = \sum_{c \in C} k_c \mathbb{P}_c \quad (4.51)$$

Then it follows that

$$\mathbb{P}_\beta^{Y|D} = \sum_{c \in C} k_c \mathbb{P}_c^{Y|D} \quad (4.52)$$

and hence \mathbb{P}_C is still $(D; Y)$ -causally contractible. Note that in order to actually implemented a random choice, we would typically consult a known random source R and set D deterministically on the basis of the value of R .

4.1.6 Body mass index revisited

If we have a probability set \mathbb{P}_C with $B := (B_i)_{i \in M}$ representing body mass index and $Y := (Y_i)_{i \in M}$ representing health outcomes of interest, the previous considerations don't support a judgement of causal contractibility for $(B; Y)$, because the choices we imagine we might have available do not allow B to be a deterministic invertible function of the choice. Note that we haven't established that causal contractibility cannot be appropriate, merely that we have no reason to accept it on the basis of arguments so far.

Causal contractibility is the a priori assumption that there is a response function relating each pair from a sequence of pairs of variables. However, we could also consider the possibility that we conclude that there is such a response function after reviewing the data.

Suppose we are in possession of a $(D; (B, Y))$ -causally contractible probability set \mathbb{P}_C , such that each (D_i, B_i, Y_i) is related by the response conditional $\mathbb{P}_C^{Y_1 B_1 | D_1 H}$. Suppose that we also have an “oracle” available that performs an infinite number of samples under appropriate conditions and reveals the value $h \in H$ yielded by the variable H . Then we can consider the new probability set $\mathbb{P}_{C,h}$ where for arbitrary $Z : \Omega \rightarrow Z$, $A \in \mathcal{Z}$

$$\mathbb{P}_{C,h}^Z(A) = \mathbb{P}_C^{Z|H}(A|h) \quad (4.53)$$

Note that $\mathbb{P}_{C,h}$ remains $(D; (B, Y))$ -causally contractible with response conditionals $\mathbb{P}_{C,h}^{Y_1 B_1 | D_1}$. Furthermore that by Theorem 2.4.27, we have $((D, B); Y)$ -causal contractibility with response conditionals $\mathbb{P}_{C,h}^{Y_1 | D_1 B_1}$. In this case, we could find that $Y_i \perp\!\!\!\perp_{\mathbb{P}_{C,h}}^e D_i | B_i$, and so by Lemma 4.1.9 $\mathbb{P}_{C,h}$ is also $(B; Y)$ -causally contractible.

In this case, it seems much more reasonable to describe the fact that B has a causal effect on Y as a *finding*, rather than an *assumption*.

4.2 Allowing dependence on observations

We want to remove the assumption $D_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{<i}C | D_{<i}$. It turns out that this is fairly straightforward to do – essentially, we require that causal contractibility holds for some C' such that $D_i \perp\!\!\!\perp_{\mathbb{P}_{C'}}^e Y_{<i}C' | D_{<i}$, and then we require that \mathbb{P}_C extends $\mathbb{P}_{C'}$ in the appropriate way. Before we do this, however, we will introduce the idea of a *comb*, which is an important generalisation of the idea of a conditional probability. The result we will eventually arrive at is that a model features response conditionals if it has a causally contractible comb of the appropriate type.

To begin with an example, consider a probability set \mathbb{P}_C , variables $D_{\mathbb{N}}$ and $Y_{\mathbb{N}}$ and a subsequence $(D_i, Y_i)_{i \in [2]}$ of length 2. Lifting the assumption $D_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{<i}C | D_{<i}$ means that only the following holds:

$$Y_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{[\mathbb{N}] \setminus i}, D_{\mathbb{N} \setminus i}C' | HD_i \quad \forall i \in \mathbb{N} \quad (4.54)$$

$$\wedge H \perp\!\!\!\perp_{\mathbb{P}_C}^e DC \quad (4.55)$$

$$\wedge \mathbb{P}_C^{Y_i | HD_i} = \mathbb{P}^{Y_i | HD_i} \quad \forall i \in \mathbb{N} \quad (4.56)$$

In this case, for arbitrary $\alpha \in C$

$$\mathbb{P}_\alpha^{Y_{[2]}} = \quad (4.57)$$

we still have response conditionals $\mathbb{P}_C^{Y_i | D_i H}$ but now the D_i 's can depend on the previous Y_i 's.

Given $\mathbb{P}_C^{Y_1 | D_1 H}$ and $\mathbb{P}_C^{Y_2 | D_2 H}$, we can define

$$\mathbb{P}_C^{Y_{[2]} | D_{[2]}} := \quad (4.58)$$

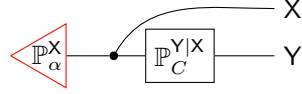
and $\mathbb{P}_C^{Y_{[2]} | D_{[2]}}$ is still “(D; Y)-causally contractible”, but it is no longer a uniform conditional probability. It is a *uniform 2-comb*.

4.2.1 Combs

Where uniform conditional probabilities map probability distributions to probability distributions via the semidirect product, 2-combs map conditional probabilities to conditional probabilities via an “insert” operation. Similarly, higher

order combs map Graphically, the semidirect product looks like

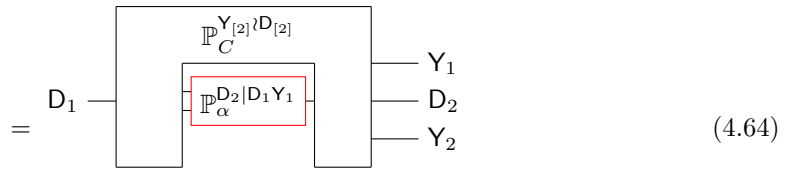
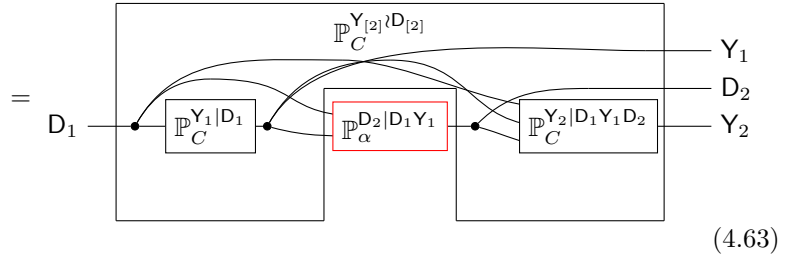
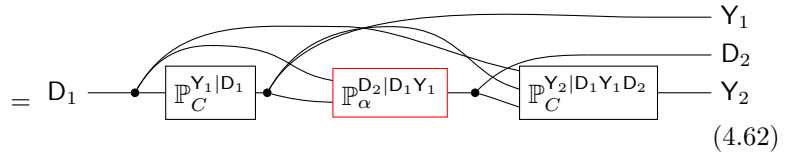
$$\mathbb{P}_\alpha^{XY} = \mathbb{P}_\alpha^X \odot \mathbb{P}_C^{Y|X} \quad (4.59)$$



$$= \quad (4.60)$$

and the insert operation looks like

$$\mathbb{P}_\alpha^{Y_1 D_2 Y_2 | D_1} = \text{insert}(\mathbb{P}_\alpha^{D_2 | D_1 Y_1}, \mathbb{P}_C^{Y_{[2]} | D_{[2]}}) \quad (4.61)$$



While Equation 4.62 is a well-formed string diagram in the category of Markov kernels, Equation 4.64 is not. In the case that all the underlying sets are discrete, Equation 4.64 can be defined using an extended string diagram notation appropriate for the category of real-valued matrices (Jacobs et al., 2019). We're not going to go to introduce the extension here, but we will give an algebraic definition of the insert operation for discrete sets.

Definition 4.2.1 (Uniform n -Comb). Given a probability set \mathbb{P}_C with variables $Y_i : \Omega \rightarrow Y$ and $D_i : \Omega \rightarrow D$ for $i \in [n]$ an uniform conditional probabilities $\{\mathbb{P}_C^{Y_i | D_{[i]} Y_{[i-1]}} | i \in [n]\}$, the uniform n -comb $\mathbb{P}_C^{Y_{[n]} | D_{[n]}} : D^n \rightarrow Y^n$ is the Markov

kernel given by the recursive definition

$$\mathbb{P}_C^{Y_1 \wr D_1} = \mathbb{P}_C^{Y_1 | D_1} \quad (4.65)$$

$$\mathbb{P}_C^{Y_{[m]} \wr D_{[m]}} = \quad (4.66)$$

Definition 4.2.2 (Uniform \mathbb{N} -comb). Given a probability set \mathbb{P}_C with variables $Y_i : \Omega \rightarrow Y$ and $D_i : \Omega \rightarrow D$ for $i \in \mathbb{N}$ and uniform conditional probabilities $\{\mathbb{P}_C^{Y_i | D_{[i]} Y_{[i-1]}} | i \in \mathbb{N}\}$, the uniform \mathbb{N} -comb $\mathbb{P}_C^{Y_{\mathbb{N}} \wr D_{\mathbb{N}}} : D^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is the Markov kernel such that for all $n \in \mathbb{N}$

$$\mathbb{P}_C^{Y_{\mathbb{N}} \wr D_{\mathbb{N}}}(\text{id}_{Y^n} \otimes \text{del}_{Y^{\mathbb{N}}}) = \mathbb{P}_C^{Y_{[n]} \wr D_{[n]}} \otimes \text{del}_{Y^{\mathbb{N}}} \quad (4.67)$$

Theorem 4.2.3 (Existence of \mathbb{N} -combs). Given a probability set \mathbb{P}_C with variables $Y_i : \Omega \rightarrow Y$ and $D_i : \Omega \rightarrow D$ for $i \in \mathbb{N}$ and uniform conditional probabilities $\{\mathbb{P}_C^{Y_i | D_{[i]} Y_{[i-1]}} | i \in \mathbb{N}\}$, a uniform \mathbb{N} -comb $\mathbb{P}_C^{Y_{\mathbb{N}} \wr D_{\mathbb{N}}} : D^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ exists.

Proof. For each $n \in \mathbb{N}$ $m < n$, we have

$$\mathbb{P}_C^{Y_{[n]} \wr D_{[n]}}(\text{id}_{Y^{n-m}} \otimes \text{del}_{Y^m}) = \mathbb{P}_C^{Y_{[n-m]} \wr D_{[n-m]}} \otimes \text{del}_{Y^m} \quad (4.68)$$

Therefore the existence of $\mathbb{P}_C^{Y_{\mathbb{N}} \wr D_{\mathbb{N}}}$ is a consequence of Lemma 4.3.8. \square

We now define the insert operation for discrete sets only

Definition 4.2.4 (Comb insert). Given an n -comb $\mathbb{P}_C^{Y_{[n]} \wr D_{[n]}}$ and an $n-1$ comb $\mathbb{P}_\alpha^{D_{[1,n]} | Y_{[n-1]}}$, (D, \mathcal{D}) and (Y, \mathcal{Y}) discrete, for all $y_i \in Y$ and $d_i \in D$

$$\text{insert}(\mathbb{P}_\alpha^{D_{[1,n]} | Y_{[n-1]}}, \mathbb{P}_C^{Y_{[n]} \wr D_{[n]}})(y_{[n]}, d_{[n-1]} | d_1) = \mathbb{P}_C^{Y_{[n]} \wr D_{[n]}}(y_n | d_{[1,n]}, d_1) \mathbb{P}_\alpha^{D_{[1,n]} | Y_{[n-1]}}(d_{[1,n]} | y_{[n-1]}) \quad (4.69)$$

4.2.2 Response conditionals in models with history dependence

The main theorem of this section, Theorem 4.2.5 is an almost trivial extension of Theorem 4.1.8. Instead of starting with a probability set \mathbb{P}_C with a uniform conditional $\mathbb{P}_C^{Y | D}$, we assume instead a collection of uniform conditionals $\mathbb{P}_C^{Y_i | D \leq_i Y_{<i}}$ (this is a strictly weaker condition, as we can derive these as higher order conditionals from the assumption of $\mathbb{P}_C^{Y | D}$). Theorem 4.2.5 then says if the \mathbb{N} -comb $\mathbb{P}_C^{Y \wr D}$ is causally contractible, then we have response conditionals

$\mathbb{P}_C^{Y_i|D_iH}$. Because $\mathbb{P}_C^{Y_iD}$ is a Markov kernel with the same signature as $\mathbb{P}_C^{Y|D}$, the reasoning from Theorem 4.2.5 carries over almost directly..

Unlike causal contractibility in the history independent case, we are not aware of an equivalent condition to the causal contractibility of $\mathbb{P}_C^{Y_iD}$ in terms of relabeled random variables. The difficulty is, because D_i can depend on $Y_{<i}$, we usually don't have symmetry with respect to relabeling variables.

Theorem 4.2.5. *Suppose we have a sample space (Ω, \mathcal{F}) , a probability set \mathbb{P}_C and variables $D := (D_i)_{i \in \mathbb{N}}$ and $Y := (Y_i)_{i \in \mathbb{N}}$ with uniform \mathbb{N} -comb $\mathbb{P}_C^{Y_iD}$. $\mathbb{P}_C^{Y_iD}$ is causally contractible if and only if Ω can be extended with some $H : \Omega \times H \rightarrow H$ such that \mathbb{P}_C^H and $\mathbb{P}_C^{Y_i|HD_i}$ exist for all $i \in \mathbb{N}$ and*

$$\mathbb{P}_C^{Y_iD} = \begin{array}{c} \triangle \mathbb{P}_C^H \\ \text{---} \bullet \text{---} \boxed{\mathbb{P}_C^{Y_0|HD_0} \text{---} Y_i} \\ \text{---} D \end{array} \quad i \in \mathbb{N} \quad (4.70)$$

$$\iff \quad (4.71)$$

$$Y_i \perp\!\!\!\perp_{\mathbb{P}_C}^e Y_{N \setminus i}, D_{N \setminus i} C | HD_i \quad \forall i \in \mathbb{N} \quad (4.72)$$

$$\wedge H \perp\!\!\!\perp_{\mathbb{P}_C}^e DC \quad (4.73)$$

$$\wedge \mathbb{P}_C^{Y_i|HD_i} = \mathbb{P}_C^{Y_0|HD_0} \quad \forall i \in \mathbb{N} \quad (4.74)$$

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

lemma I've yet to prove

Proof. Apply Lemma to $\mathbb{P}_C^{Y_iD}$. □

4.2.3 Validity

4.2.4 Combs are the output of the “fix” operation

There is a relationship between combs and the “fix” operation defined in Richardson et al. (2017). In particular, suppose we have a probability \mathbb{P}_α and a comb $\mathbb{P}_\alpha^{Y_{[2]}|D_{[2]}}$. Then (assuming discrete sets)

$$\mathbb{P}_\alpha^{Y_{[2]}|D_{[2]}}(y_1, y_2 | d_1, d_2) = \mathbb{P}_\alpha^{Y_1|D_1}(y_1 | d_1) \mathbb{P}_\alpha^{Y_2|D_2}(y_2 | d_2) \quad (4.75)$$

$$= \frac{\mathbb{P}_\alpha^{Y_1|D_1}(y_1 | d_1) \mathbb{P}_\alpha^{D_2|Y_1D_1}(d_2 | y_1, d_1) \mathbb{P}_\alpha^{Y_2|D_2}(y_2 | d_2)}{\mathbb{P}_\alpha^{D_2|Y_1D_1}(d_2 | y_1, d_1)} \quad (4.76)$$

$$= \frac{\mathbb{P}_\alpha^{Y_{[2]}|D_2|D_1}(y_1, y_2, d_2 | d_1)}{\mathbb{P}_\alpha^{D_2|Y_1D_1}(d_2 | y_1, d_1)} \quad (4.77)$$

This is precisely the “division by a conditional probability” used in the fix operation. We speculate that the fix operation is precisely an alternative definition of an n -comb, but we have not proven this.

4.3 Weaker assumptions than causal contractibility

Definition 4.3.1 (Local). A Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is *local* if for all $n \in \mathbb{N}$, $A_i \in \mathcal{Y}$, $(x_{[n]}, x_{[n]^c}) \in \mathbb{N}$ there exists $\mathbb{L} : X^n \rightarrow Y^n$ such that

$$\boxed{\dots} = \boxed{\dots} \quad (4.78)$$

$$\iff \quad (4.79)$$

$$\mathbb{K}(\bigtimes_{i \in [n]} A_i \times Y^{\mathbb{N}} | x_{[n]}, x_{[n]^c}) = \mathbb{L}(\bigtimes_{i \in [n]} A_i | x_{[n]}) \quad (4.80)$$

Definition 4.3.2 (Exchange commutativity). A Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ *commutes with exchange* if for all finite permutations $\rho : \mathbb{N} \rightarrow \mathbb{N}$, $A_i \in \mathcal{Y}$, $(x_{[n]}, x_{[n]^c}) \in \mathbb{N}$

$$\mathbb{K} \text{swap}_{\rho_Y} = \text{swap}_{\rho_X} \mathbb{K} \quad (4.81)$$

$$\iff \quad (4.82)$$

$$\mathbb{K}(\bigtimes_{i \in \mathbb{N}} A_{\rho(i)} | \bigotimes_{i \in \mathbb{N}} x_i) = \mathbb{K}(\bigtimes_{i \in \mathbb{N}} A_i | \bigotimes_{i \in \mathbb{N}} x_{\rho(i)}) \quad (4.83)$$

Definition 4.3.3 (Causal contractibility). A Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is *causally contractible* if it is local and commutes with exchange.

Definition 4.3.4 (Finite set swap). Given two equally sized sequences $A = (a_i)_{i \in [n]}$, $B = (b_i)_{i \in [n]}$, $A \leftrightarrow B : \mathbb{N} \rightarrow \mathbb{N}$ is the permutation that sends the i th element of A to the i th element of B and vice versa. Note that $A \leftrightarrow B$ is its own inverse.

Theorem 4.3.5 (Equality of equally sized contractions). *A Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is causally contractible if and only if for every $n \in \mathbb{N}$ there exists $\mathbb{L}_n : X^n \rightarrow Y^n$ such that for all $A \subset \mathbb{N}$ where $|A| = n$*

$$\text{swap}_{[n] \leftrightarrow A} \mathbb{K} \text{swap}_{[n] \leftrightarrow A} (\text{id}_{[n]} \otimes \text{del}_{[n]^c}) = \mathbb{L}_n \otimes \text{del}_{[n]^c} \quad (4.84)$$

Proof. Only if: By exchange commutativity

$$\text{swap}_{[n] \leftrightarrow A} \mathbb{K} = \mathbb{K} \text{swap}_{A \leftrightarrow n} \quad (4.85)$$

multiply both sides by $\text{swap}_{A \leftrightarrow n}$ on the right and we have (because $\text{swap}_{[n] \leftrightarrow A}$ is its own inverse)

$$\text{swap}_{[n] \leftrightarrow A} \mathbb{K} \text{swap}_{A \leftrightarrow n} = \mathbb{K} \quad (4.86)$$

Then by locality, there exists some $\mathbb{L} : X^n \rightarrow Y^n$ such that

$$\text{swap}_{[n] \leftrightarrow A} \mathbb{K} \text{swap}_{A \leftrightarrow n} (\text{id}_{[n]} \otimes \text{del}_{[n]^c}) = \mathbb{K} (\text{id}_{[n]} \otimes \text{del}_{[n]^c}) \quad (4.87)$$

$$= \mathbb{L} \otimes \text{del}_{[n]^c} \quad (4.88)$$

If: Taking $A = [n]$ establishes locality immediately.

For exchange commutativity, note that for all $x \in X^{\mathbb{N}}$, $n \in \mathbb{N}$, we have

$$\text{swap}_{A \leftrightarrow [n]} \mathbb{K} \text{swap}_{A \leftrightarrow [n]} (\text{id}_{[n]} \otimes \text{del}_{[n]^c}) = \mathbb{K} (\text{id}_{[n]} \otimes \text{del}_{[n]^c}) \quad (4.89)$$

Then by Lemma 4.3.8

$$\text{swap}_{A \leftrightarrow [n]} \mathbb{K} \text{swap}_{A \leftrightarrow [n]} = \mathbb{K} \quad (4.90)$$

Consider an arbitrary finite permutation $\rho : \mathbb{N} \rightarrow \mathbb{N}$. ρ can be decomposed into a finite set of cyclic permutations on disjoint orbits. Each cyclic permutation is simply the composition of a sequence of 1-cycles of the form $A \leftrightarrow [n]$, and so ρ itself can be written as a composition of a sequence of 1-cycles. Thus for any finite $\rho : \mathbb{N} \rightarrow \mathbb{N}$

$$\text{swap}_{\rho} \mathbb{K} \text{swap}_{\rho} = \mathbb{K} \quad (4.91)$$

□

Theorem 4.3.6. *A Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is causally contractible if and only if there exists a column exchangeable probability distribution $\mu \Delta(Y^{|X| \times \mathbb{N}})$ such that*

$$\mathbb{K} = \begin{array}{c} \triangle \mu \\ \text{---} \\ X \text{---} \boxed{\mathbb{F}_{\text{ev}}} \text{---} Y \end{array} \quad (4.92)$$

$$\iff \quad (4.93)$$

$$\mathbb{K}(A | (x_i)_{i \in \mathbb{N}}) = \mu \Pi_{(x_i)_{i \in \mathbb{N}}}(A) \forall A \in \mathcal{Y}^{\mathbb{N}} \quad (4.94)$$

Where $\Pi_{(d_i i)_{i \in \mathbb{N}}} : Y^{|X| \times \mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is the function

$$(y_{ji})_{j, i \in X \times \mathbb{N}} \mapsto (y_{di})_{i \in \mathbb{N}} \quad (4.95)$$

that projects the (x_i, i) indices of y for all $i \in \mathbb{N}$, and \mathbb{F}_{ev} is the Markov kernel associated with the evaluation map

$$\text{ev} : X^{\mathbb{N}} \times Y^{X \times \mathbb{N}} \rightarrow Y \quad (4.96)$$

$$((x_i)_{i \in \mathbb{N}}, (y_{ji})_{j, i \in X \times \mathbb{N}}) \mapsto (y_{x_i i})_{i \in \mathbb{N}} \quad (4.97)$$

Proof. Only if: Choose $e := (e_i)_{i \in \mathbb{N}}$ such that $e_{i+|X|j}$ is the i th element of X for all $i, j \in \mathbb{N}$.

Define

$$\mu \left(\bigtimes_{(i,j) \in X \times \mathbb{N}} A_{ij} \right) := \mathbb{K} \left(\bigtimes_{(i,j) \in X \times \mathbb{N}} A_{ij} | e \right) \forall A_{ij} \in \mathcal{Y} \quad (4.98)$$

Now consider any $x := (x_i)_{i \in \mathbb{N}} \in X^{\mathbb{N}}$. By definition of e , $e_{x_i i} = x_i$ for any $i, j \in \mathbb{N}$.

Define

$$\mathbb{Q} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}} \quad (4.99)$$

$$\mathbb{Q} := \begin{array}{c} \triangle \mu \\ \text{---} X \text{---} \boxed{\mathbb{F}_{\text{ev}}} \text{---} Y \end{array} \quad (4.100)$$

and consider some $A \subset \mathbb{N}$, $|A| = n$ and $B := (x_i, i)_{i \in A}$. Note that the subsequence of e indexed by B , $e_B := (e_{x_i i})_{i \in A} = x_A$. Thus given the swap map $\text{swap}_{A \leftrightarrow B} : \mathbb{N} \rightarrow \mathbb{N}$ that sends the first element of A to the first element of B and so forth, $\text{swap}_{A \leftrightarrow B}(e_B) = x_A$. For arbitrary $\{C_i \in \mathcal{Y} \mid i \in A\}$, define $C_A := \text{swap}_{[n] \leftrightarrow A}(\times_{i \in [n]} C_i \times Y^{\mathbb{N}})$. Then, for arbitrary $x \in X^{\mathbb{N}}$

$$\mathbb{Q}(C_A | x) = \mu(\text{ev}_x^{-1}(C_A)) \quad (4.101)$$

The argument of μ is

$$\text{ev}_x^{-1}(C_A) = \{(y_{ji})_{j, i \in X \times \mathbb{N}} \mid (y_{x_i i})_{i \in \mathbb{N}} \in C_A\} \quad (4.102)$$

$$= \bigtimes_{i \in \mathbb{N}} \bigtimes_{j \in X} D_{ji} \quad (4.103)$$

where

$$D_{ji} = \begin{cases} C_i & (j, i) \in B \\ Y & \text{otherwise} \end{cases} \quad (4.104)$$

and so

$$\text{swap}_{A \leftrightarrow B}(\text{ev}_x^{-1}(C_A)) = C_A \quad (4.105)$$

Substituting Equation 4.105 into 4.101

$$\mathbb{Q}(C_A | x) = \mu \text{swap}_{A \leftrightarrow B}(C_A) \quad (4.106)$$

$$= \mathbb{K} \text{swap}_{A \leftrightarrow B}(C_A | e) \quad (4.107)$$

$$= \mathbb{K} \text{swap}_{A \leftrightarrow B}(C_A | e_B, \text{swap}_{B \leftrightarrow A}(x)_B^C) \quad \text{by locality} \quad (4.108)$$

$$= \mathbb{K} \text{swap}_{A \leftrightarrow B}(C_A | \text{swap}_{B \leftrightarrow A}(x)) \quad (4.109)$$

$$= \text{swap}_{B \leftrightarrow A} \mathbb{K} \text{swap}_{A \leftrightarrow B}(C_A | x) \quad (4.110)$$

$$= \mathbb{K}(C_A | x) \quad \text{by commutativity of exchange} \quad (4.111)$$

Because this holds for all x , $A \subset \mathbb{N}$, by Lemma 4.3.8

$$\mathbb{Q} = \mathbb{K} \quad (4.112)$$

Next we will show μ is column exchangeable. Consider any column swap $\text{swap}_c : X \times \mathbb{N} \rightarrow X \times \mathbb{N}$ that acts as the identity on the X component and a finite permutation on the \mathbb{N} component. From the definition of e , $\text{swap}_c(e) = e$. Thus by commutativity of exchange, for any $A \in \mathcal{Y}^{\mathbb{N}}$

$$\mathbb{K}(A|e) = \text{swap}_c \mathbb{K} \text{swap}_c(A|e) \quad (4.113)$$

$$= \mathbb{K} \text{swap}_c(A|\text{swap}_c(e)) \quad (4.114)$$

$$= \mathbb{K} \text{swap}_c(A|e) \quad (4.115)$$

If: Suppose

$$\mathbb{K} = \begin{array}{c} \triangleleft \mu \\ \text{---} X \text{---} \boxed{\mathbb{F}_{\text{ev}}} \text{---} Y \end{array} \quad (4.116)$$

where μ is column exchangeable, and consider any two $x, x' \in X^{\mathbb{N}}$ such that some subsequences are equal $x_S = x'_T$ with $S, T \subset \mathbb{N}$ and $|S| = |T| = [n]$.

For any $\{A_i \in \mathcal{Y} | i \in S\}$, let $A_S = \text{swap}_{[n] \leftrightarrow S} \times_{i \in [n]} A_i \times Y^{\mathbb{N}}$, $A_T = \text{swap}_{S \leftrightarrow T}(A_S)$, $B = (x_i i)_{i \in S}$ and $C = (x_i i)_{i \in T} = (x_{\text{swap}_{S \leftrightarrow T}}(i) i)_{i \in S}$. By Equations 4.101 and 4.105

$$\mathbb{K}(A_S|x) = \mu \text{swap}_{S \leftrightarrow B}(A_S) \quad (4.117)$$

$$= \mu \text{swap}_{T \leftrightarrow C}(A_T) \quad \text{by column exchangeability of } \mu \quad (4.118)$$

$$= \mathbb{K}(A_T|\text{swap}_{S \leftrightarrow T}(x)) \quad (4.119)$$

$$= \text{swap}_{S \leftrightarrow T} \mathbb{K}(A_T|x) \quad (4.120)$$

$$= \text{swap}_{S \leftrightarrow T} \mathbb{K} \text{swap}_{S \leftrightarrow T}(A_S|x) \quad (4.121)$$

so \mathbb{K} is causally contractible by Theorem 4.3.5. \square

Theorem 4.3.7. A kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ is causally contractible if and only if there exists some set H , $\mu \in \Delta(H)$ and $\mathbb{L} : H \times X \rightarrow Y$ such that

$$\mathbb{K} = \begin{array}{c} \triangleleft \mu \\ \text{---} \bullet \text{---} \boxed{\mathbb{L}} \text{---} Y \\ \text{---} X \text{---} \end{array} \quad i \in \mathbb{N} \quad (4.122)$$

$$\iff \quad (4.123)$$

$$\mathbb{K}(\bigotimes_{i \in \mathbb{N}} A_i | (x_i)_{i \in \mathbb{N}}) = \int_H \prod_{i \in \mathbb{N}} \mathbb{L}(A_i | h, x_i) \mu(dh) \quad (4.124)$$

Proof. By Theorem 4.3.6, we can represent the conditional probability \mathbb{K} as

$$\mathbb{K} = \begin{array}{c} \triangleleft \mu \\ \text{---} X \text{---} \boxed{\mathbb{F}_{\text{ev}}} \text{---} Y \end{array} \quad (4.125)$$

$$\mathbb{K} = \left[\begin{array}{c} \text{Diagram: A triangle labeled } \mathbb{P}_C^H \text{ is connected to a black dot. This dot is connected to a box labeled } \mathbb{P}_C^{Y_0|HD_0}. \text{ The box also receives input } D. \text{ The output of the box is } Y_i. \text{ The entire structure is indexed by } i \in \mathbb{N}. \end{array} \right] \quad (4.136)$$

Where we can connect the outputs of μ to the inputs of \mathbb{F}_{evs} “inside the plate” as the plates in Equations 4.126 and 4.133 are equal in number and each connected wire represents a single copy of Y^D .

If: By assumption, for any $\{A_i \in \mathcal{Y} | i \in \mathbb{N}\}$, $x := (x_i)_{i \in \mathbb{N}} \in X^{\mathbb{N}}$

$$\mathbb{K}(\bigotimes_{i \in \mathbb{N}} A_i | x) = \int_H \prod_{i \in \mathbb{N}} \mathbb{L}(A_i | h, x_i) \mu(dh) \quad (4.137)$$

Consider any $S, T \subset \mathbb{N}$ with $|S| = |T|$, and define $A_S := \times_{i \in \mathbb{N}} B_i$ where $B_i = Y$ if $i \notin S$, otherwise A_i is an arbitrary element of \mathcal{Y} . Define $A_T := \times_{i \in \mathbb{N}} B_{\text{swap}_{S \leftrightarrow T}(i)}$.

$$\mathbb{K}(A_S | x) = \int_H \prod_{i \in S} \mathbb{L}(A_i | h, x_i) \mu(dh) \quad (4.138)$$

$$= \int_H \prod_{i \in T} \mathbb{L}(A_i | h, x_{\text{swap}_{S \leftrightarrow T}(i)}) \mu(dh) \quad (4.139)$$

$$= \text{swap}_{S \leftrightarrow T} \mathbb{K}(A_T | x) \quad (4.140)$$

$$= \text{swap}_{S \leftrightarrow T} \mathbb{K} \text{swap}_{S \leftrightarrow T}(A_S | x) \quad (4.141)$$

So by Theorem 4.3.5, \mathbb{K} is causally contractible. \square

Lemma 4.3.8 (Infinitely extended kernels). *Given a collection of Markov kernels $\mathbb{K}_i : X^i \rightarrow Y^i$ for all $i \in \mathbb{N}$, if we have for every $j > i$*

$$\mathbb{K}_j(\text{id}_{X_i} \otimes \text{del}_{X_{j-i}}) = \mathbb{K}_i \otimes \text{del}_{X_{j-i}} \quad (4.142)$$

then there is a unique Markov kernel $\mathbb{K} : X^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$ such that for all $i, j \in \mathbb{N}, j > i$

$$\mathbb{K}(\text{id}_{X_i} \otimes \text{del}_{X_{j-i}}) = \mathbb{K}_i \otimes \text{del}_{X_{j-i}} \quad (4.143)$$

Proof. Take any $x \in X^{\mathbb{N}}$ and let $x|_m \in X^m$ be the first m elements of x . By Equation 4.142, for any $A_i \in \mathcal{Y}$, $i \in [m]$

$$\mathbb{K}_n(\bigotimes_{i \in [m]} A_i \times Y^{n-m} | x|_n) = \mathbb{K}_m(\bigotimes_{i \in [m]} A_i | x|_m) \quad (4.144)$$

Furthermore, by the definition of the swap map for any permutation $\rho : [n] \rightarrow [n]$

$$\mathbb{K}_n \text{swap}_{\rho}(\bigotimes_{i \in [m]} A_{\rho(i)} \times Y^{n-m} | x|_n) = \mathbb{K}_n(\bigotimes_{i \in [m]} A_i \times Y^{n-m} | x|_n) \quad (4.145)$$

Thus by the Kolmogorov Extension Theorem (Çinlar, 2011), for each $x \in X^{\mathbb{N}}$ there is a unique probability measure $\mathbb{Q}_x \in \Delta(Y^{\mathbb{N}})$ satisfying

$$\mathbb{Q}_d(\bigotimes_{i \in [n]} A_i \times Y^{\mathbb{N}}) = \mathbb{K}_n(\bigotimes_{i \in [n]} A_{\rho(i)} | d|_n) \quad (4.146)$$

Furthermore, for each $\{A_i \in \mathcal{Y} | i \in \mathbb{N}\}$, $n \in \mathbb{N}$ note that for $p > n$

$$\mathbb{Q}_d(\bigtimes_{i \in [n]} A_i \times Y^{\mathbb{N}}) \geq \mathbb{Q}_d(\bigtimes_{i \in [p]} A_i \times Y^{\mathbb{N}}) \quad (4.147)$$

$$\geq \mathbb{Q}_d(\bigtimes_{i \in \mathbb{N}} A_i) \quad (4.148)$$

so by the Monotone convergence theorem, the sequence $\mathbb{Q}_d(\bigtimes_{i \in [n]} A_i)$ converges as $n \rightarrow \infty$ to $\mathbb{Q}_d(\bigtimes_{i \in \mathbb{N}} A_i)$. $d \mapsto \mathbb{Q}_d^{Z^n}(\bigtimes_{i \in [n]} A_i)$ is measurable for all n , $\{A_i \in \mathcal{Y} | i \in \mathbb{N}\}$ by Equation 4.146, and so $d \mapsto Q_d$ is also measurable.

Thus $d \mapsto Q_d$ is the desired $\mathbb{P}_C^{Y^{\mathbb{N}} | D^{\mathbb{N}}} : D^{\mathbb{N}} \rightarrow Y^{\mathbb{N}}$. \square

Theorem 4.3.9 (Conditional independence in augmented causally contractible model). *Suppose we have a probability set \mathbb{P}_C on Ω with causally contractible \mathbb{N} -comb $\mathbb{P}_C^{Y^{\mathbb{N}} | D^{\mathbb{N}}}$, $Y = (Y_i)_{i \in \mathbb{N}}$, $D = (D_i)_{i \in \mathbb{N}}$, augmented with the hypothesis variable H as in Lemma 4.3.11. Then for all n , $Y_n \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y_{<n}, D_{<n}) | (H, D_n)$.*

Proof. For all $\alpha \in C$, $n \in \mathbb{N}$

$$\mathbb{P}_\alpha^{Y_{[n]} D_{[n]} | H} = \quad (4.149)$$

Thus for all α , by Theorem 2.4.27

$$\mathbb{P}_C^{Y_n | D_n H} \otimes \text{del}_{D^{n-1} \times Y^{n-1}} \stackrel{\mathbb{P}_\alpha}{\cong} \mathbb{P}_C^{Y_n | D_{[n]} Y_{<n} H} \quad (4.150)$$

which implies $Y_n \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y_{<n}, D_{<n}) | (H, D_n)$. \square

Theorem 4.3.10 (Conditional independence in non-data-dependant causally contractible model). *Suppose we have a probability set \mathbb{P}_C on Ω with causally contractible \mathbb{N} -comb $\mathbb{P}_C^{Y^{\mathbb{N}} | D^{\mathbb{N}}}$, $Y = (Y_i)_{i \in \mathbb{N}}$, $D = (D_i)_{i \in \mathbb{N}}$, augmented with the hypothesis variable H as in Lemma 4.3.11. Then for all n , $Y_n \perp\!\!\!\perp_{\mathbb{P}_C}^e (Y_{<n}, D_{<n}) | (H, D_n)$.*

Lemma 4.3.11. *Suppose we have a probability set \mathbb{P}'_C on (Ω', \mathcal{F}') with causally contractible \mathbb{N} -comb $\mathbb{P}'_C^{Y'^{\mathbb{N}} | D'^{\mathbb{N}}}$, $Y' = (Y'_i)_{i \in \mathbb{N}}$, $D' = (D'_i)_{i \in \mathbb{N}}$. Then there exists H and an augmented model \mathbb{P}_C on $(\Omega, \mathcal{F}) := ((\Omega' \times H, \mathcal{F}' \otimes \mathcal{H}))$ such that $\mathbb{P}_C \Pi_{\Omega'} = \mathbb{P}'_C$ and, defining $H : \Omega' \times H \rightarrow H$ as the projection onto H ,*

$$\mathbb{P}_C^{Y^{\mathbb{N}} | D^{\mathbb{N}}} = \quad (4.151)$$

Proof. Let H be the hypothesis space from Theorem 4.3.7 and \mathbb{P}_C^H be some directing random measure for $\mathbb{P}_C^{Y^D}$. Then by Theorem 4.151, Equation 4.151 holds. Furthermore, by equality of combs, we have for all $\alpha \in C$

$$\mathbb{P}_\alpha^{YD} = \mathbb{P}_\alpha^{Y'D'} \quad (4.152)$$

Define $W' : \Omega' \rightarrow \Omega'$ as the identity function on Ω' , $W : \Omega' \times H \rightarrow \Omega'$ as the projection to Ω' . For each $\alpha \in C$, define \mathbb{P}_α by

$$\mathbb{P}_\alpha^W = \mathbb{P}_\alpha^{YD} \odot \mathbb{P}_\alpha^{W'|Y'D'}(\text{del}_{Y^{ND^N}H} \otimes \text{id}_{\Omega'}) \quad (4.153)$$

Then

$$\mathbb{P}_\alpha^W = \mathbb{P}_\alpha \Pi_{\Omega'} \quad (4.154)$$

$$= \mathbb{P}_\alpha^{Y'D'} \odot \mathbb{P}_\alpha^{W'|Y'D'}(\text{del}_{Y^{ND^N}H} \otimes \text{id}_{\Omega'}) \quad (4.155)$$

$$= \mathbb{P}_\alpha^{W'} \quad (4.156)$$

□

Chapter 5

Statistical Decision Theory

Can non-comb version of probability set represent statistical decision theory?

5.1 Summary

Statistical models are ubiquitous in the analysis of inference problems. A statistical model features a set of *states*, and each state is mapped to a probability distribution over *outcomes*. If we want to model problems involving *decisions* and *consequences*, we need to consider different kinds of statistical models. We introduce two types of model for this purpose: *two player statistical models* which differ from classical statistical model in that the state is assumed to consist of a decision and a *hypothesis* (the two players are the decision maker “player D” and the hypothesis selector “player H”). They model the consequences of decisions under various hypotheses. A *see-do model* is a special case of two player statistical model that can be used in situations where some *observations* are available for review prior to selecting a decision. See-do models are the main focus of work here and problems involving observations, decisions and consequences will be discussed at length in Chapter 4.

A common simplifying assumption made when using classical statistical models is that they are *conditionally independent and identically distributed* (conditionally IID); this means that the model maps each state to an independent and identically distributed (IID) sequence of observations. This just a common choice, it is not a strict requirement. A similar assumption is likely to be useful for see-do models. If consequences depend on choices, then it does not make sense to assert that observations and consequences together form a single IID sequence of random variables, so we need to consider alternatives. We propose that models where observations are an IID sequence and choices and consequences together are *independent and functionally identical* (IFI; defined later in this chapter) are similar to conditionally IID statistical models.

Instead of directly assuming that a conditionally IID model is appropriate, *De Finetti's representation theorem* shows that probability models where the sequence of observations is *exchangeable* induce conditionally IID statistical models. The assumption of exchangeable observations is preferable to the assumption of IID observations if a probability model is being used to represent subjective uncertainty. We investigate whether there is an analogous result relating “exchangeability-like” assumptions for see-do models to “IID-like” assumptions. We show that there is: in particular, a see-do “forecast” with exchangeable observations and *functionally exchangeable* decision to consequence maps induces a see-do model with IID observations and IFI consequences.

The assumption of functional exchangeability will appear again in Chapter 6 as part of the definition of *counterfactual models*, and the joint assumptions of exchangeable observations and functionally exchangeable consequences to motivate the assumption of *imitability* in Chapter ??, an assumption that in combination with a number of other assumptions allows for inference of consequences from data.

5.2 Modelling observations, choices and consequences

5.2.1 Modelling observations with statistical models

Statistical models are a ubiquitous type of model in statistics and machine learning. They consist of a set of *states* (S, \mathcal{S}) , and for each state the model prescribes a single probability distribution on a given measurable set of *outcomes* (O, \mathcal{O}) .

Definition 5.2.1 (Statistical model). A statistical model is a set of states (S, \mathcal{S}) , a set of outcomes (O, \mathcal{O}) and a stochastic map $\mathbb{T} : S \rightarrow \Delta(\mathcal{O})$.

Definition 5.2.2 (State and outcome variables). Given a statistical model $(\mathbb{T}, (O, \mathcal{O}), (S, \mathcal{S}))$, define the *state variable* $\mathbf{S} : S \times O \rightarrow S$ as the projection from $S \times O \rightarrow S$ and define the *outcome variable* $\mathbf{O} : S \times O \rightarrow O$ as the projection onto O .

The common example of a potentially biased coin is modelled with a statistical model. We suppose our coin has some rate of heads $\theta \in [0, 1]$, and we furthermore suppose that for each θ the result of flipping the coin can be modeled (in some sense) by the probability distribution $\text{Bernoulli}(\theta)$. The statistical model here is the set of states $S = [0, 1]$ (corresponding to *rates of heads*), the observation space $O = \{0, 1\}^n$ with the discrete sigma-algebra (where n is the number of flips observed) and the stochastic map $\mathbb{B} : [0, 1] \rightarrow \Delta(\mathcal{P}(0, 1))$ which is given by $\mathbb{B} : \theta \rightarrow \text{Bernoulli}(\theta)$.

Almost any theoretical treatment of statistics or machine learning will at some point make use of statistical models to describe the problem they are addressing – for a collection of examples from the last 100 years, see Goodfellow

et al. (2016); Vapnik (2013); Bishop (2006); Le Cam (1996); Freedman (1963); Wald (1950); de Finetti ([1937] 1992); Fisher ([1925] 1992). They are often simply assumed without a great deal of discussion of why this type of model is chosen, or what role they play.

If we want to reason about how well some learning algorithm performs in some context, we typically require a reasonable model of the context in which the learning algorithm operates. The algorithms themselves may not give us such a model. Because learning almost always operates in a context with noise and uncertainty, we need models that can handle noise and uncertainty. Probability models are a very common choice for this. In addition, it is often assumed that we do not know with certainty the exact probability model that should be used to model a context. A statistical model assumes a certain number of states may prevail – reflecting uncertainty in the “mechanics of the world” – and given any state it gives us a probability distribution – reflecting uncertainty remaining after the mechanics of the world are well-understood.

Learning algorithms don’t necessarily implement reasonable models of the world. For example, consider a linear regressor that takes a set of predictors $x \in X$ and targets $y \in Y$ and returns some $\beta \in B$ such that $(y - x^T \beta)^2$ is as small as possible. It is possible to interpret B as a set of states, and consider the learner to be implementing the statistical model $(\mathbb{T}, B, \mathcal{L}_{X \rightarrow Y})$ where $\mathcal{L}_{X \rightarrow Y}$ is the set of linear function $X \rightarrow Y$ given by $\{x \mapsto x^T \beta | \beta \in B\}$ and \mathbb{T} maps a state $\beta \in B$ deterministically to the function $x \mapsto x^T \beta$. This is, formally, a statistical model, but it is not one that would typically be considered a good model of the world in the kinds of problems that a linear regressor is used to solve. One problem with this model is that it is deterministic - the outcome for any $\beta \in B$ will be a particular function $X \rightarrow Y$. However, it will almost never be the case that some set of targets y will be an exact function of some set of predictors x , and insisting on an exact functional relationship will typically give very poor generalisation results if this demand can be satisfied at all.

Suppose we want to ask whether the function f given by a linear regressor is useful for some purpose. In order to address this question, we want to consider a more appropriate model of the world than the crude statistical model given above, and consider what behaviour we will see from the regressor under different assumptions imposed on this model. Statistical models typically are used to serve the purpose of a “more appropriate model of the world”. In this example, we might consider a statistical model (\mathbb{T}, H, O) where for each $h \in H$, $\mathbb{T}_h \in \Delta(\mathcal{Y} \otimes \mathcal{X})$ such that $\mathbb{T}_h^{\mathcal{Y}|\mathcal{X}} = \text{Normal}(\mathbf{X}^T \beta_h, \sigma_h)$. If we assume the data generating process is described by such a probability distribution for some h , we can ask questions like “does the linear regressor output a β such that $\mathbf{Y} - \mathbf{X}^T \beta < \epsilon$ with high probability for all $h \in H$?”

5.2.2 Modelling choices and consequences with two-player statistical models

The states in a statistical model are usually considered to be “under the control of nature”. In the possibly biased coin example above, if we were to consider some “player D” acutally flipping the coin and trying to infer the bias, we would typically assume that their opinion about the coin’s bias does not affect the coin’s actual bias; they could decide it is biased towards heads when in fact it is completely unbiased. In some cases player D can make choices that affect the outcomes. Suppose player D has the option to choose how high to toss the coin – perhaps they can aim for a toss height anywhere from 10 to 50cm. This plausibly affects the outcomes of their coin toss and, unlike the coin’s bias, they gets to choose the height they intend to toss it. If they decide to toss it to a height of 15cm then 15cm is the height they have chosen to toss it to. Unlike the state, which can differ from whatever player D ultimately decides on, the choice made by player D is the same thing as whatever they ultimately decide on. We call features of the state that are not under player D’s control *hypotheses* and features that are under player D’s control *decisions*, and statistical models in which the state is the Cartesian product of a set of hypotheses and a set of decisions “two player statistical models” (the two players being nature or “player H” and the decision maker or “player D”).

Definition 5.2.3 (Two player statistical model). A *two-player statistical model* is a tuple $(\mathbb{T}, \mathcal{H}, \mathcal{D}, \mathcal{O})$ where $(\mathbb{T}, (H \times D, \mathcal{H} \otimes \mathcal{D}), (\mathcal{O}, \mathcal{O}))$ is a statistical model and $H : H \times D \times O \rightarrow H$, $D : H \times D \times O \rightarrow D$ and $O : H \times D \times O \rightarrow O$ are measurable functions that project elements of $H \times D \times O$ to their respective codomains. H is called the *hypothesis*, D the *decision* and O the *outcome*.

Whenever we propose a two player statistical model, we will also assume for any random variables $X : H \times D \times O \rightarrow X$ and $Y : H \times D \times O \rightarrow Y$, a disintegration $\mathbb{K}^{Y|XDH} : X \times D \times H \rightarrow \Delta(\mathcal{Y})$ exists (see Theorem ?? for a sufficient condition).

The problems that we will mostly study in this work, in addition to having a second player (“player D”), will often involve some data X that is observed before the second player is able to make a choice. Two player statistical models with *observations* are called *see-do models*.

Definition 5.2.4 (See-Do model). A *see-do model* $(\mathbb{T}, \mathcal{H}, \mathcal{D}, X, Y)$ is a two-player statistical model along with two additional random variables: the *observation* $X : H \times D \times O \rightarrow X$ and the *consequence* $Y : H \times D \times O \rightarrow Y$. The outcome variable is defined to be the coupled product of the observation and the consequence $O = (X, Y)$, and we will leave this implicit when specifying a see-do model. A see-do model must observe the conditional independence:

$$X \perp\!\!\!\perp_{\mathbb{T}} D | H \quad (5.1)$$

We can informally read the independence requirement as saying “the observations are independent of the decision given the hypothesis”. This does not

imply that probability models we construct from \mathbb{T} will necessarily have the property that D and X will be independent conditional on H , and in fact this will often not be the case. In Chapter 6 we will argue that this requirement captures the intuition that observations are not “affected” by decisions. For now, we will observe that this independence requirement means that \mathbb{T} can be drawn with no path from D to X .

Explicitly, the independence on line 5.1 implies that the kernel \mathbb{T} can be drawn as follows:

$$\mathbb{T} := \begin{array}{c} \begin{array}{ccc} H & \bullet & \boxed{\mathbb{T}^{X|H}} - X \\ & \searrow & \\ D & \longrightarrow & \boxed{\mathbb{T}^{Y|HD}} - Y \end{array} \end{array} \quad (5.2)$$

In this picture, again informally, Y takes input from D but X does not.

Chapter 6

See-do models, interventions and counterfactuals

6.1 How do see-do models relate to other approaches to causal inference?

- Review of approaches: CBN, CBN soft intervention, CBN fat-hand intervention, CBN noise intervention, SEM (Pearl/Heckman), PO unit model, PO population model, SWIG, Dawid decision theoretic model, Heckerman decision theoretic model, Rohde/Lattimore Bayesian model
- Focus on CBN, PO unit model, PO population model

6.2 Interpretations of the choice set

- Decisions or actions we could actually make - decision problem
- Idealised/hypothetical choices constrained by a set of causal relationships - interventions
- Suppositions - counterfactuals
- Further possibility - intervention \rightarrow decisions might be actuator randomisation

6.3 Causal Bayesian Networks as see-do models

- Definition of CBN, intervention set (recall: existence of disintegrations, decomposability)

- How interventions differ from decisions: no effect strength uncertainty, side effects, may be more interventions than what we actually know how to do
- Example: sets of CBNs and d-separation

6.4 Unit Potential Outcomes models

- Counterfactual random variables Y_x answer a question: "what would Y be supposing X was x ?"
- Proposed formalisation of suppositions: (....)
- Implies existence of counterfactual random variables
- Difference between suppositions and decisions: determinism, other conditions
- "3-player models": hypotheses, suppositions and interventions/decisions
- Error in key theorem of Ruben, Imbens (ignorability does not imply functional exchangeability)
- What can be represented by a 3 player model?
 - "1 of 2 counterfactuals": anything
 - "3 of 2 counterfactuals": very restrictive
 - "2 of 3 counterfactuals": Bell's theorem, counterfactual definiteness

This chapter is currently a disorganised cut and paste

The field of causal inference is additionally concerned with types of questions called "counterfactual" by Pearl. There is substantial theoretical interest in counterfactual questions, but counterfactual questions are much more rarely found in applications than interventional questions. Even though see-do models are motivated by the need to answer interventional questions, the theory developed here is surprisingly applicable to counterfactuals as well. In particular, the theory of see-do models offers explanations for three key features of counterfactual models:

- **Apparent absence of choices:** *Potential outcomes* models, which purportedly answer counterfactual questions, are standard statistical models *without choices* (Rubin, 2005)
- **Deterministic dependence on unobserved variables:** Counterfactual models involve *deterministic* dependence on unobserved variables (Pearl, 2009; Rubin, 2005; Richardson and Robins, 2013)

- **Residual dependence on observations:** Counterfactual questions depend on the given data *even if the joint distribution of this data is known*. For example, Pearl (2009) introduces a particular method for conditioning a known joint distribution on observations that he calls *abduction*

Potential outcomes models lack a notion of “choices” because there is a generic method to “add choices” to a potential outcomes model, which is implicitly used whenever potential outcomes models are used. Furthermore, we show that a see-do model induces a potential outcome model if and only if it is a model of *parallel choices*, and in this case the observed consequences depend deterministically on the unobserved potential outcomes in precisely the manner as given in Rubin (2005). Parallel choices can be roughly understood as models of sequences of experiments where an action can be chosen for each experiment, and with the special properties that repeating the same action deterministically yields the same consequence, and the consequences of a sequence of actions doesn’t depend on the order in which the actions are taken. That is, we show that the fundamental property of any “counterfactual” model is *deterministic reproducibility* and *action exchangeability*, and while these models may admit a “counterfactual” interpretation, they are fundamentally just a special class of see-do models.

But the proof is still in my notebook

Interestingly, it seems to be possible to construct a see-do model where the “hypothesis” is a quantum state, and quantum mechanics + locality seems to rule out parallel choices in such models in a manner similar to Bell’s theorem. “Seems to” because I haven’t actually proven any of these things.

The residual dependence on observations exhibited by counterfactual questions is a generic property of see-do models, and it is a particular property of *decision problems* are notable in that it is often

Where to discuss the connections to statistical decision theory?

See-do models are closely related to *statistical decision theory* introduced by Wald (1950) and elaborated by Savage (1954) after Wald’s death. See-do models equipped with a *utility function* induce a slightly generalised form of statistical decision problems, and the complete class theorem is applicable to these models.

A stylistic difference between see-do models and most other causal models is that see-do models explicitly represent both the observation model and the consequence model and their coupling, making them “two picture” causal models. Causal Bayesian Networks and Single World Intervention Graphs (Richardson and Robins, 2013) use “one picture” to represent the observation model and the consequence model. However, both of these approaches employ “graph mutilation”, so one picture on the page actually corresponds to many pictures when combined with the mutilation rules. For more on how these different types of models relate, see Section ?? . Lattimore and Rohde (2019)’s Bayesian causal inference employs two-picture causal models, as do “twin networks” (Pearl, 2009).

Sometimes we are interested in modelling situations where we can also make some choices that also affect the eventual consequences. For example, I might hypothesise H_1 : the switch on the wall controls my light, H_2 : the switch on the wall does not control my light. Then, given H_1 I can choose to toggle the switch, and I will see my light turn on, or I can choose not to toggle the switch and I will not see my light turn on. Given H_2 , neither choice will result in a light turned on. Choices are clearly different to hypotheses: the choice I make depends on what I want to happen, while whether or not a hypothesis is true has no regard for my ambitions.

A “statistical model with choices” is simply a map $\mathbb{T} : D \times H \rightarrow \Delta(\mathcal{E})$ for some set of choices D , hypotheses H and outcome space (E, \mathcal{E}) . We can also distinguish two types of outcomes: *observations* which are given prior to a choice being made and *consequences* which happen after a choice is made. Observations cannot be affected by the choices made, while consequences are not subject to this restriction. That is, observations are what we might *see* before making a choice, which depends on the hypothesis alone, and if we are lucky we may be able to invert this dependence to learn something about the hypothesis from observations. On the other hand, the consequences of what we *do* depends jointly on the hypothesis and the choice we make and we judge which choices are more desirable on the basis of which consequences we expect them to produce.

What we are studying is a family of models that generalises of statistical models to include hypotheses, choices, observations and consequences. These models are referred to as *see-do models*. Hypotheses, observations, consequences and choices are not individually new ideas. *Statistical decision problems* (Wald, 1950; ?) extend statistical models with decisions and *losses*. Like consequences, losses depend on which choices are made. However, unlike consequences, losses must be ordered and reflect the preferences of a decision maker. *Influence diagrams* are directed graphs created to represent decision problems that feature “choice nodes”, “chance nodes” and “utility nodes”. An influence diagram may be associated with a particular probability distribution Nilsson and Lauritzen (2013) or with a set of probability distributions Dawid (2002).

See-do models have deep roots in decision theory. Decision theory asks, out of a set of available acts, which ones ought to be chosen. See-do models answer an intermediate question: out of a set of available acts, what are the consequences of each? This question is described by Pearl (2009) as an “interventional” question.

See-do models depend crucially on a set of choices D . While these models can obviously answer questions like “what is likely to happen if I choose $d \in D$?”, this construction appears to rule out “causal” questions like “Does rain cause wet roads?”. We define a restricted idea of causation called *D-causation*. Roughly, if the roads get wet when it rains regardless of my choice of $d \in D$, then rain “*D*-causes” wet roads. *D-causation* is closely related to the idea *limited invariance* put forward by Heckerman and Shachter (1995).

6.4.1 D-causation

The choice set D is a primitive element of a see-do model. However, while we claim that see-do models are the basic objects studied in causal inference, so far we have no notion of “causation”. What we call *D-causation* is one such notion. It is called *D-causation* because it is a notion of causation that depends on the set of choices available. A similar idea, called *limited unresponsiveness*, is discussed extensively in the decision theoretic account of causation found in Heckerman and Shachter (1995). The main difference is that see-do maps are fundamentally stochastic while Heckerman and Shachter work with “states” (approximately hypotheses in our terminology) that map decisions deterministically to consequences. In addition, while we define *D-causation* relative to a see-do map \mathbb{T} , Heckerman and Shachter define limited unresponsiveness with respect to *sets* of states.

Section ?? explores the difficulty of defining “objective causation” without reference to a set of choices. D need not be interpreted as the set of choices available to an agent, but however we want to interpret it, all existing examples of causal models seem to require this set.

See Section ?? for the definition of random variables in Kernel spaces.

One way to motivate the notion of *D-causation* is to observe that for many decision problems, I may wish to include a very large set of choices D . Suppose I aim to have my light switched on, and there is a switch that controls the light. Often, the relevant choices for such a problem would appear to be $D_0 = \{\text{flip the switch, don't flip the switch}\}$. However, this doesn't come close to exhausting the set of things I might choose to do, and I might wish to consider a larger set of possibilities. For simplicity's sake, suppose I have instead the following set of options:

$$D_1 := \{ \text{“walk to the switch and press it with my thumb”}, \\ \text{“trip over the lego on the floor, hop to the light switch and stab my finger at it”}, \\ \text{“stay in bed”} \}$$

If having the light turned on is all that matters, I could consider any acts in D_1 to be equivalent if, in the end, the light switch ends up in the same position. In this case, I could say that the light switch position D_1 -causes the state of the light. Subject to the assumption that the light switch position D_1 -causes the state of the light, I can reduce my problem to one of choosing from D_0 (noting that some choices correspond to mixtures of elements of D_0).

If I consider an even larger set of possible acts D_2 , I might not accept that the switch position D_2 -causes the state of the light. Let D_2 be the following acts:

$D_2 := \{$ “walk to the switch and press it with my thumb”,
“trip over the lego on the floor, hop to the light switch and stab my finger at it”,
“stay in bed”,
“toggle the mains power, then flip the light switch” $\}$

In this case, it would be unreasonable to suppose that all acts that left the light switch in the “on” position would also result in the light being “on”. Thus the switch does not D_2 -cause the light to be on.

Formally, D -causation is defined in terms of conditional independence. Given a see-do model $\mathbb{T} : H \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$, define the *consequence model* $\mathbb{C} : H \times D \rightarrow \Delta(\mathcal{Y})$ as $\mathbb{C} := \mathbb{T}^{\mathcal{Y}|\mathcal{H}\mathcal{D}}$.

Definition 6.4.1 (D -causation). Given a hypothesis $h \in H$ and a consequence model $\mathbb{C} : H \times D \rightarrow \Delta(\mathcal{Y})$, random variables $\mathbf{Y}_1 : Y \times D \rightarrow Y_1$, $\mathbf{Y}_2 : Y \times D \rightarrow Y_2$ and $\mathbf{D} : Y \times D \rightarrow D$ (defined the usual way), \mathbf{Y}_1 D -causes \mathbf{Y}_2 iff $\mathbf{Y}_2 \perp\!\!\!\perp_{\mathbb{C}} \mathbf{D} | \mathbf{Y}_1 H$.

6.4.2 D-causation vs Limited Unresponsiveness

Heckerman and Shachter study deterministic “consequence models”. Furthermore, what we call hypotheses $h \in H$, Heckerman and Schachter call states $s \in S$. Heckerman and Shachter’s notion of causation is defined by *limited unresponsiveness* rather than *conditional independence*, which depends on a partition of states rather than a particular hypothesis.

Definition 6.4.2 (Limited unresponsiveness). Given states S , deterministic consequence models $\mathbb{C}_s : D \rightarrow \Delta(F)$ for each $s \in A$ and a random variables $\mathbf{Y}_1 : F \rightarrow Y_1$, $\mathbf{Y}_2 : F \rightarrow Y_2$, \mathbf{Y}_1 is unresponsive to \mathbf{D} in states limited by \mathbf{Y}_2 if $\mathbb{C}_{(s,d)}^{\mathbf{Y}_2|\mathbf{SD}} = \mathbb{C}_{(s,d')}^{\mathbf{Y}_2|\mathbf{SD}} \implies \mathbb{C}_{(s,d)}^{\mathbf{Y}_1|\mathbf{SD}} = \mathbb{C}_{(s,d')}^{\mathbf{Y}_1|\mathbf{SD}}$ for all $d, d' \in D$, $s \in S$. Write $\mathbf{Y}_1 \not\prec_{\mathbf{Y}_2} \mathbf{D}$

Lemma 6.4.3 (Limited unresponsiveness implies D -causation). *For deterministic consequence models, $\mathbf{Y}_1 \not\prec_{\mathbf{Y}_2} \mathbf{D}$ implies \mathbf{Y}_2 D -causes \mathbf{Y}_1 .*

Proof. By the assumption of determinism, for each $s \in S$ and $d \in D$ there exists $y_1(s, d)$ and $y_2(s, d)$ such that $\mathbb{C}_{s,d}^{\mathbf{Y}_1\mathbf{Y}_2|\mathbf{SD}} = \delta_{y_1(s,d)} \otimes \delta_{y_2(s,d)}$.

By the assumption of limited unresponsiveness, for all d, d' such that $y_2(s, d) = y_2(s, d')$, $y_1(s, d) = y_1(s, d')$ also. Define $f : Y_2 \times S \rightarrow Y_1$ by $(s, y_1) \mapsto y_1(s, [y_1(s, \cdot)]^{-1}(y_1(s, d)))$ where $[y_1(s, \cdot)]^{-1}(a)$ is an arbitrary element of $\{d | y_1(s, d) = a\}$. For all s, d , $f(y_1(s, d), s) = y_2(s, d)$. Define $\mathbb{M} : Y_2 \times S \times D \rightarrow \Delta(\mathcal{Y}_1)$ by $(y_2, s, d) \mapsto \delta_{f(y_2, s)}$. \mathbb{M} is a version of $\mathbb{C}^{\mathbf{Y}_1|\mathbf{Y}_2, S, \mathbf{D}}$ because, for all $A \in \mathcal{Y}_2$, $B \in \mathcal{Y}_1$, $s \in S$, $d \in D$:

$$\mathbb{C}_{(s,d)}^{\mathbf{Y}_2|\mathbf{SD}} \mathbb{Y}(\mathbb{M} \otimes \text{Id}) = \int_A \mathbb{M}(y'_2, d, s; B) d\delta_{y_2(s,d)}(y'_2) \quad (6.1)$$

$$= \int_A \delta_{f(y'_2, s)}(B) d\delta_{y_2(s,d)}(y'_2) \quad (6.2)$$

$$= \delta_{f(y_2(s,d), s)}(B) \delta_{y_2(s,d)}(A) \quad (6.3)$$

$$= \delta_{y_1(s,d)}(B) \delta_{y_2(s,d)}(A) \quad (6.4)$$

$$= \delta_{y_2(s,d)} \otimes \delta_{y_1(s,d)}(A \times B) \quad (6.5)$$

\mathbb{M} is clearly constant in D . Therefore $\mathbf{Y}_1 \perp\!\!\!\perp_{\mathbb{C}} D | \mathbf{Y}_2 S$. \square

However, despite limited unresponsiveness implying D -causation, it does not imply D -causation in mixtures of states. Suppose $D = \{0, 1\}$ where 1 stands for “toggle light switch” and 0 stands for “do nothing”. Suppose $S = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ where $[0, 0]$ represents “switch initially off, mains off” the other states generalise this in the obvious way. Finally, $F \in \{0, 1\}$ is the final position of the switch and $L \in \{0, 1\}$ is the final state of the light. We have

define this

$$\mathbb{C}_{d, [i, m]}^{\mathbf{LF}|\mathbf{DS}} = \delta_{(d \text{ XOR } i) \text{ AND } m} \otimes \delta_{(d \text{ XOR } i) \text{ AND } m} \quad (6.6)$$

Within states $[0, 0]$ and $[1, 0]$, the light is always off, so $F = a \implies L = 0$ for any a . In states $[0, 1]$ and $[1, 1]$, $F = 1 \implies L = 1$ and $F = 0 \implies L = 0$. Thus $L \not\prec_F D$. However, suppose we take a mixture of consequence models:

$$\mathbb{C}_\gamma = \frac{1}{4}\mathbb{C}_{\cdot, [0, 0]} + \frac{1}{4}\mathbb{C}_{\cdot, [0, 1]} + \frac{1}{2}\mathbb{C}_{\cdot, [1, 1]} \quad (6.7)$$

$$\mathbb{C}_\gamma^{\mathbf{FL}|\mathbf{D}} = \frac{1}{4} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \otimes \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (6.8)$$

Then

$$[1, 0] \mathbb{C}_\gamma^{\mathbf{FL}|\mathbf{D}} = \frac{1}{4}[0, 1] \otimes [1, 0] + \frac{1}{4}[0, 1] \otimes [0, 1] + \frac{1}{2}[1, 0] \otimes [1, 0] \quad (6.9)$$

$$[1, 0] \mathbb{Y}(\mathbb{C}_\gamma^{\mathbf{F}|\mathbf{D}} \otimes \mathbb{C}_\gamma^{\mathbf{L}|\mathbf{D}}) = (\frac{1}{2}[0, 1] + \frac{1}{2}[1, 0]) \otimes (\frac{1}{4}[0, 1] + \frac{3}{4}[1, 0]) \quad (6.10)$$

$$\implies [1, 0] \mathbb{C}_\gamma^{\mathbf{FL}|\mathbf{D}} \neq [1, 0] \mathbb{Y}(\mathbb{C}_\gamma^{\mathbf{F}|\mathbf{D}} \otimes \mathbb{C}_\gamma^{\mathbf{L}|\mathbf{D}}) \quad (6.11)$$

Thus under the prior γ , F does not D -cause L even though F D -causes L in all states S . The definition of D -causation was motivated by the idea that we could reduce a difficult decision problem with a large set D to a simpler problem with a smaller “effective” set of decisions by exploiting conditional independence. Even if X D -causes Y in every $H \in S$, X does not necessarily D -cause Y in mixtures of states in S . For this reason, we do not say that X D -causes Y in S if X D -causes Y in every $H \in S$, and in this way we differ substantially from Heckerman and Shachter (1995).

define this

Instead, we simply extend the definition of D -causation to mixtures of hypotheses: if $\gamma \in \Delta(\mathbf{H})$ is a mixture of hypotheses, define $\mathbb{C}_\gamma := (\gamma \otimes \mathbf{Id})\mathbb{C}$. Then \mathbf{X} D -causes \mathbf{Y} relative to γ iff $\mathbf{Y} \perp\!\!\!\perp_{\mathbb{C}_\gamma} \mathbf{D} | \mathbf{X}$.

Theorem 6.4.4 shows that under some conditions, D -causation can hold for arbitrary mixtures over subsets of the hypothesis class \mathbf{H} .

Theorem 6.4.4 (Universal D -causation). *If $\mathbf{X} \perp\!\!\!\perp \mathbf{H} | \mathbf{D}$ for all $\mathbf{H}, \mathbf{H}' \in S \subset \mathbf{H}$ and \mathbf{X} D -causes \mathbf{Y} in all $\mathbf{H} \in S$, then \mathbf{X} D -causes \mathbf{Y} with respect to all mixed consequence models \mathbb{C}_γ for all $\gamma \in \Delta(\mathbf{H})$ with $\gamma(S) = 1$.*

Proof. For $\gamma \in \Delta(\mathbf{H})$, define the mixture

$$\mathbb{C}_\gamma := \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}} \text{---} \mathbf{F} \end{array} \quad (6.12)$$

Because $\mathbb{C}_\mathbf{H}^{\mathbf{X}|\mathbf{D}} = \mathbb{C}_{\mathbf{H}'}^{\mathbf{X}|\mathbf{D}}$ for all $\mathbf{H}, \mathbf{H}' \in \mathbf{H}$, we have

$$\begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \mathbf{X} \end{array} \quad \mathbf{H} = \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \mathbf{X} \end{array} \quad \mathbf{H} \quad (6.13)$$

Also

$$\mathbb{C}_\gamma^{\mathbf{X}\mathbf{Y}|\mathbf{D}} = \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}} \text{---} \boxed{\mathbb{F}^{\mathbf{X} \otimes \mathbf{Y}}} \text{---} \begin{array}{l} \mathbf{X} \\ \mathbf{Y} \end{array} \end{array} \quad (6.14)$$

$$= \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}\mathbf{Y}|\mathbf{D}\mathbf{H}}} \text{---} \begin{array}{l} \mathbf{X} \\ \mathbf{Y} \end{array} \end{array} \quad (6.15)$$

$$= \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \boxed{\mathbb{C}^{\mathbf{Y}|\mathbf{X}\mathbf{D}\mathbf{H}}} \text{---} \begin{array}{l} \mathbf{Y} \\ \mathbf{X} \end{array} \end{array} \quad (6.16)$$

$$\stackrel{\mathbf{Y} \perp\!\!\!\perp_{\mathbb{C}_\gamma} \mathbf{D} | \mathbf{X}\mathbf{H}}{=} \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \boxed{\mathbb{C}^{\mathbf{Y}|\mathbf{X}\mathbf{H}}} \text{---} \begin{array}{l} \mathbf{Y} \\ \mathbf{X} \end{array} \end{array} \quad (6.17)$$

$$\stackrel{6.13}{=} \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \boxed{\mathbb{C}^{\mathbf{Y}|\mathbf{X}\mathbf{H}}} \text{---} \begin{array}{l} \mathbf{Y} \\ \mathbf{X} \end{array} \end{array} \quad (6.18)$$

$$\stackrel{6.13}{=} \begin{array}{c} \triangle \gamma \\ | \\ \text{D} \text{---} \boxed{\mathbb{C}_\gamma^{\mathbf{X}|\mathbf{D}\mathbf{H}}} \text{---} \boxed{\mathbb{C}^{\mathbf{Y}|\mathbf{X}\mathbf{H}}} \text{---} \begin{array}{l} \mathbf{Y} \\ \mathbf{X} \end{array} \end{array} \quad (6.19)$$

Equation 6.19 establishes that $(\gamma \otimes \mathbf{Id}_X \otimes \mathbf{*}_D) \mathbb{C}^{Y|X^H}$ is a version of $\mathbb{C}_\gamma^{Y|XD}$, and thus $Y \perp\!\!\!\perp_{\mathbb{C}_\gamma} D|X$.

This can also be derived from the semi-graphoid rules:

$$H \perp\!\!\!\perp D \wedge H \perp\!\!\!\perp X|D \implies H \perp\!\!\!\perp XD \quad (6.20)$$

$$\implies H \perp\!\!\!\perp D|X \quad (6.21)$$

$$D \perp\!\!\!\perp H|X \wedge D \perp\!\!\!\perp Y|XH \implies D \perp\!\!\!\perp Y|X \quad (6.22)$$

$$\implies Y \perp\!\!\!\perp D|X \quad (6.23)$$

□

6.4.3 Properties of D-causation

If X D-causes Y relative to \mathbb{C}_H , then the following holds:

$$\mathbb{C}_H^{X|D} = D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (6.24)$$

This follows from version (2) of Definition ??:

$$\mathbb{C}_H^{X|D} = D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|XD}} - Y \quad (6.25)$$

$$= D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (6.26)$$

$$= D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (6.27)$$

D-causation is not transitive: if X D-causes Y and Y D-causes Z then X doesn't necessarily D-cause Z .

Pearl's “front door adjustment” and general identification results make use of composing “sub-consequence-kernels” like this. Show, if possible, that Pearl's “sub-consequence-kernels” obey D -causation like relations

Does this “weak D-causation” respect mixing under the same conditions as regular D-causation?

6.4.4 Decision sequences and parallel decisions

Just as observations X can be a sequence of random variables X_1, X_2, \dots , D can be a sequence of “sub-choices” D_1, D_2, \dots . Note that by positing such a sequence there is no requirement that D_1 comes “before” D_2 in any particular sense.

6.5 Existence of counterfactuals

I'm struggling with how to explain this well.

“Counterfactual” or “potential outcomes” models in the causal inference literature are consequence models where choices can be considered in *parallel*.

Before defining parallel choices, we will consider a “counterfactual model” without parallel choices. Consider the following definitions, first from Pearl (2009) pg. 203-204. I have preserved his notation, including not using any special fonts for things called “variables” because this term is used interchangeably with “sets of variables” and using special fonts for variables might give the impression that these should be treated as different things while using special fonts for sets of variables is inconsistent with my usual notation.

The real solution here is that Pearl’s “variable sets” are actually “coupled variables”, see Definition ??, but I’d rather not change his definitions if I can avoid it

put the following inside a quote environment somehow, the regular quote environment fails due to too much markup

““

Definition 7.1.1 (Causal Model) A causal model is a triple $M = \langle U, V, F \rangle$, where:

- (i) U is a set of *background* variables, (also called *exogenous*), that are determined by factors outside the model;
- (ii) V is a set $\{V_1, V_2, \dots, V_n\}$ of variables, called *endogenous*, that are determined by variables in the model – that is, variables in $U \cup V$;
- (iii) F is a set of functions $\{f_1, f_2, \dots, f_n\}$ such that each f_i is a mapping from (the respective domains of) $U_i \cup PA_i$ to V_i , where $U_i \subseteq U$ and $PA_i \subseteq V \setminus V_i$ and the entire set F forms a mapping from U to V . In other words, each f_i in

$$v_i = f_i(pa_i, u_i), \quad i \in 1, \dots, n,$$

assigns a value to V_i that depends on (the values of) a select set of variables in $V \cup U$, and the entire set F has a unique solution $V(u)$.

Definition 7.1.2 (Submodel) Let M be a causal model, X a set of variables in V , and x a particular realization of X . A submodel M_x of M is the causal model

$$M_x = \{U, V, F_x\},$$

where

$$F_x = \{f_i : V_i \notin X\} \cup \{X = x\}.$$

Definition 7.1.3 (Effect of Action) Let M be a causal model, X a set of variables in V , and x a particular realization of X . The effect of action $do(X = x)$ on M is given by the submodel M_x .

Definition 7.1.4 (Potential Response) Let X and Y be two subsets of variables in V . The potential response of Y to action $do(X = x)$, denoted $Y_x(u)$, is the solution for Y of the set of equations F_x , that is, $Y_x(u) = Y_{M_x}(u)$.

Definition 7.1.6 (Probabilistic Causal Model) A probabilistic causal model is a pair $\langle M, P(u) \rangle$, where M is a causal model and $P(u)$ is a probability function defined over the domain of U . ”

Implicitly, Definition 7.1.3 proposes a set of “actions” that have “effects” given by M_x . It’s not entirely clear what this set of actions should be – the definition seems to suggest that there is an action for each “realization” of each variable in V , which would imply that the set of actions corresponds to the range of V . For the following discussion, we will call the set of actions D , whatever it actually contains (we have deliberately chosen to use the same letter as we use to represent choices or actions in see-do models).

Given D , Definition 7.1.3 appears to define a function $h : \mathcal{M} \times D \rightarrow \mathcal{M}$, where \mathcal{M} is the space of causal models with background variables U and endogenous variables V , such that for $M \in \mathcal{M}$, $do(X = x) \in D$, $h(M, do(X = x)) = M_x$.

Definition 7.1.4 then appears to define a function $Y(\cdot) : D \times U \rightarrow Y$ (distinct from Y , which appears to be a function $U \rightarrow \text{something}$) and calls $Y(\cdot)$ the “potential response”. We could always consider the variable $\mathbf{V} := \bigotimes_{i \in [n]} V_i$ and define the “total potential response” $\mathbf{g} := \mathbf{V}(\cdot)$, which captures the potential responses of any subset of variables in V .

From this, we might surmise that in the Pearlean view, it is necessary that a “counterfactual” or “potential response” model has a probability measure P on background variables U , a set of actions D and a *deterministic* potential response function $\mathbf{g} : D \times U \rightarrow V$.

Pearl’s model also features a second deterministic function $\mathbf{f} : U \rightarrow Y$, and G is derived from F via the equation modifications permitted by D . It is straightforward to show that an arbitrary function $\mathbf{f} : U \rightarrow Y$ can be constructed from Pearl’s set of functions f_i , and if D may modify the set F arbitrarily, then it appears that \mathbf{g} can in principle be an arbitrary function $D \times U \rightarrow Y$ (though many possible choices would be quite unusual).

Pearl’s counterfactual model seems to essentially be a deterministic map $\mathbf{g} : D \times U \rightarrow V$ along with a probability measure P on U . Putting these together and marginalising over U (as we might expect we want to do with “background variables”) simply yields a consequence map $D \rightarrow \Delta(\mathcal{V})$, which doesn’t seem to have any special counterfactual properties.

In order to pose counterfactual questions, Pearl introduces the idea of holding U fixed:

““

Definition 7.1.5 (Counterfactual) Let X and Y be two subsets of variables in V . The counterfactual sentence “ Y would be y (in situation u), had X been x ” is interpreted as the equality $Y_x(u) = y$, with $Y_x(u)$ being the potential response of Y to $X = x$.’ ”

Holding U fixed allows SCM counterfactual models to answer questions about what would have happened if we had taken different actions given the same background context. For example, we can compare $Y_x(u)$ with $Y_{x'}(u)$ and interpret the comparison as telling us what would have happened in the same situation u if we did x and, at the same time, what would happen if we did x' . It is the ability to consider different actions “in exactly the same situation” that makes these models “counterfactual”.

One obvious question is: does \mathbf{g} have to be deterministic? While SCMs are defined in terms of deterministic functions with noise arguments, it’s not clear that this is a necessary feature of counterfactual models. If \mathbf{g} were properly stochastic, what is the problem with considering $\mathbf{g}(x, u)$ and $\mathbf{g}(x', u)$ to represent what would happen in a fixed situation u if I did x and if I did x' respectively? In fact, a nondeterministic \mathbf{g} arguably fails to capture a key intuition of taking actions “in exactly the same situation”. If I want to know the result of doing action x and, in exactly the same situation, the result of doing action x , then one might intuitively think that the result should always be *deterministically the same*. This property, which we call *deterministic reproducibility*, does not hold if we consider a nondeterministic potential response map \mathbf{g} .

This idea of doing x and, in the same situation, doing x doesn’t render very well in English. Furthermore, even though deterministic reproducibility seems to be an important property of counterfactual SCMs, they don’t help very much to elucidate the idea. “If I take action x in situation U I get $V_x(u)$ and if I take action x in situation U I get $V_x(u)$ ” is just a redundant repetition. It seems that we want some way to express the idea of having two copies of $V_x(u)$ or, more generally, having multiple copies of a potential response function in such a way that we can make comparisons between their results.

The idea that we need *can* be clearly expressed with a see-do model.

References

- A. V. Banerjee, S. Chassang, and E. Snowberg. Chapter 4 - Decision Theoretic Approaches to Experiment Design and External Validity. We thank Esther Duflo for her leadership on the handbook and for extensive comments on earlier drafts. Chassang and Snowberg gratefully acknowledge the support of NSF grant SES-1156154. In Abhijit Vinayak Banerjee and Esther Duflo, editors, *Handbook of Economic Field Experiments*, volume 1 of *Handbook of Field Experiments*, pages 141–174. North-Holland, January 2017. doi: 10.1016/bs.hefe.2016.08.005. URL <https://www.sciencedirect.com/science/article/pii/S2214658X16300071>.

- Elias Bareinboim, Juan D. Correa, Duligur Ibeling, and Thomas Icard. On Pearl's hierarchy and the foundations of causal inference. Technical report, 2020. URL <https://causalai.net/r60.pdf>.
- Christopher Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, New York, 2006. ISBN 978-0-387-31073-2. URL <https://www.springer.com/gp/book/9780387310732>.
- Vladimir Bogachev and Ilya Malofeev. Kantorovich problems and conditional measures depending on a parameter. *Journal of Mathematical Analysis and Applications*, 486:123883, June 2020. doi: 10.1016/j.jmaa.2020.123883.
- Ethan D. Bolker. Functions Resembling Quotients of Measures. *Transactions of the American Mathematical Society*, 124(2):292–312, 1966. ISSN 0002-9947. doi: 10.2307/1994401. URL <https://www.jstor.org/stable/1994401>. Publisher: American Mathematical Society.
- Ethan D. Bolker. A Simultaneous Axiomatization of Utility and Subjective Probability. *Philosophy of Science*, 34(4):333–340, 1967. ISSN 0031-8248. URL <https://www.jstor.org/stable/186122>. Publisher: [The University of Chicago Press, Philosophy of Science Association].
- 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.
- Panayiota Constantinou and A. Philip Dawid. Extended conditional independence and applications in causal inference. *The Annals of Statistics*, 45(6): 2618–2653, 2017. ISSN 0090-5364. URL <http://www.jstor.org/stable/26362953>. Publisher: Institute of Mathematical Statistics.
- A. P. Dawid. Causal Inference without Counterfactuals. *Journal of the American Statistical Association*, 95(450):407–424, June 2000. ISSN 0162-1459. doi: 10.1080/01621459.2000.10474210. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.2000.10474210>. Publisher: Taylor & Francis _eprint: <https://www.tandfonline.com/doi/pdf/10.1080/01621459.2000.10474210>.
- A. 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.
- Bruno de Finetti. Foresight: Its Logical Laws, Its Subjective Sources. In Samuel Kotz and Norman L. Johnson, editors, *Breakthroughs in Statistics: Foundations and Basic Theory*, Springer Series in Statistics, pages 134–174. Springer, New York, NY, [1937] 1992. ISBN 978-1-4612-0919-5. doi: 10.1007/978-1-4612-0919-5_10. URL https://doi.org/10.1007/978-1-4612-0919-5_10.
- R.P. Feynman. *The Feynman lectures on physics*. Le cours de physique de Feynman. Interditions, France, 1979. ISBN 978-2-7296-0030-3. INIS Reference Number: 13648743.
- R. A. Fisher. Statistical Methods for Research Workers. In Samuel Kotz and Norman L. Johnson, editors, *Breakthroughs in Statistics: Methodology and Distribution*, Springer Series in Statistics, pages 66–70. Springer, New York, NY, [1925] 1992. ISBN 978-1-4612-4380-9. doi: 10.1007/978-1-4612-4380-9_6. URL https://doi.org/10.1007/978-1-4612-4380-9_6.
- 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.
- David A. Freedman. On the Asymptotic Behavior of Bayes’ Estimates in the Discrete Case. *Annals of Mathematical Statistics*, 34(4):1386–1403, December 1963. ISSN 0003-4851, 2168-8990. doi: 10.1214/aoms/1177703871. URL <https://projecteuclid.org/euclid.aoms/1177703871>. Publisher: Institute of Mathematical Statistics.
- Tobias Fritz. A synthetic approach to Markov kernels, conditional independence and theorems on sufficient statistics. *Advances in Mathematics*, 370:107239, August 2020. ISSN 0001-8708. doi: 10.1016/j.aim.2020.107239. URL <https://www.sciencedirect.com/science/article/pii/S0001870820302656>.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Sander Greenland and James M Robins. Identifiability, Exchangeability, and Epidemiological Confounding. *International Journal of Epidemiology*, 15(3): 413–419, September 1986. ISSN 0300-5771. doi: 10.1093/ije/15.3.413. URL <https://doi.org/10.1093/ije/15.3.413>.

- J. Y. Halpern. A Counter Example to Theorems of Cox and Fine. *Journal of Artificial Intelligence Research*, 10:67–85, February 1999. ISSN 1076-9757. doi: 10.1613/jair.536. URL <https://www.jair.org/index.php/jair/article/view/10223>.
- 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>.
- Miguel A. Hernán and James M Robins. *Causal Inference: What If*. Chapman & Hall/CRC, 2020. URL <https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>.
- Eric Horvitz, David Heckerman, and Curtis Langlotz. A Framework for Comparing Alternative Formalisms for Plausible Reasoning. January 1986. URL <https://openreview.net/forum?id=rJNeX0gdbr>.
- Guido W. Imbens and Donald B. Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press, Cambridge, 2015. ISBN 978-0-521-88588-1. doi: 10.1017/CBO9781139025751. URL <https://www.cambridge.org/core/books/causal-inference-for-statistics-social-and-biomedical-sciences/71126BE90C58F1A431FE9B2DD07938AB>.
- Bart Jacobs, Aleks Kissinger, and Fabio Zanasi. Causal Inference by String Diagram Surgery. In Mikoaj Bojaczek and Alex Simpson, editors, *Foundations of Software Science and Computation Structures*, Lecture Notes in Computer Science, pages 313–329. Springer International Publishing, 2019. ISBN 978-3-030-17127-8.
- Richard C. Jeffrey. *The Logic of Decision*. University of Chicago Press, July 1965. ISBN 978-0-226-39582-1.
- Olav Kallenberg. The Basic Symmetries. In *Probabilistic Symmetries and Invariance Principles*, Probability and Its Applications, pages 24–68. Springer, New York, NY, 2005. ISBN 978-0-387-28861-1. doi: 10.1007/0-387-28861-9_2. URL https://doi.org/10.1007/0-387-28861-9_2.
- 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.

L. Le Cam. Comparison of Experiments - A Short Review.pdf. *IMS Lecture Notes - Monograph Series*, 30, 1996.

David Lewis. Causal decision theory. *Australasian Journal of Philosophy*, 59(1): 5–30, March 1981. ISSN 0004-8402. doi: 10.1080/00048408112340011. URL <https://doi.org/10.1080/00048408112340011>.

Karl Menger. Random Variables from the Point of View of a General Theory of Variables. In Karl Menger, Bert Schweizer, Abe Sklar, Karl Sigmund, Peter Gruber, Edmund Hlawka, Ludwig Reich, and Leopold Schmetterer, editors, *Selecta Mathematica: Volume 2*, pages 367–381. Springer, Vienna, 2003. ISBN 978-3-7091-6045-9. doi: 10.1007/978-3-7091-6045-9_31. URL https://doi.org/10.1007/978-3-7091-6045-9_31.

Dennis Nilsson and Steffen L. Lauritzen. Evaluating Influence Diagrams using LIMIDs. *arXiv:1301.3881 [cs]*, January 2013. URL <http://arxiv.org/abs/1301.3881>. arXiv: 1301.3881.

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 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.

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.

Thomas S. Richardson, Robin J. Evans, James M. Robins, and Ilya Shpitser. Nested Markov Properties for Acyclic Directed Mixed Graphs. *arXiv:1701.06686 [stat]*, January 2017. URL <http://arxiv.org/abs/1701.06686>. arXiv: 1701.06686.

- 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>.
- L. J. Savage. The theory of statistical decision. *Journal of the American Statistical Association*, 46:55–67, 1951. ISSN 1537-274X(Electronic),0162-1459(Print). doi: 10.2307/2280094.
- Leonard J. Savage. *The Foundations of Statistics*. Wiley Publications in Statistics, 1954.
- P. Selinger. A Survey of Graphical Languages for Monoidal Categories. In Bob Coecke, editor, *New Structures for Physics*, Lecture Notes in Physics, pages 289–355. Springer, Berlin, Heidelberg, 2011. ISBN 978-3-642-12821-9. doi: 10.1007/978-3-642-12821-9_4. URL https://doi.org/10.1007/978-3-642-12821-9_4.
- Eyal Shahar. The association of body mass index with health outcomes: causal, inconsistent, or confounded? *American Journal of Epidemiology*, 170(8):957–958, October 2009. ISSN 1476-6256. doi: 10.1093/aje/kwp292.
- 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>.
- Brian Skyrms. Causal Decision Theory. *The Journal of Philosophy*, 79(11):695–711, November 1982. doi: 10.2307/2026547. URL https://www.pdcnet.org/pdc/bvdb.nsf/purchase?openform&fp=jphil&id=jphil_1982_0079_0011_0695_0711.
- Peter Spirtes, Clark N Glymour, Richard Scheines, David Heckerman, Christopher Meek, Gregory Cooper, and Thomas Richardson. *Causation, prediction, and search*. MIT press, 2000.
- Statista. Cigarettes - worldwide | Statista Market Forecast, 2020. URL <https://www.statista.com/outlook/50010000/100/cigarettes/worldwide>.
- Katie Steele and H. Orri Stefánsson. Decision Theory. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2020 edition, 2020. URL <https://plato.stanford.edu/archives/win2020/entries/decision-theory/>.
- Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer Science & Business Media, June 2013. ISBN 978-1-4757-3264-1. Google-Books-ID: EqACAAAQBAJ.

- J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Theory of games and economic behavior. Princeton University Press, Princeton, NJ, US, 1944.
- N. N. Vorobev. Consistent Families of Measures and Their Extensions. *Theory of Probability & Its Applications*, 7(2), 1962. doi: 10.1137/1107014. URL http://www.mathnet.ru/php/archive.phtml?wshow=paper&jrnid=tvtp&paperid=4710&option_lang=eng.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England, 1950.
- Peter Walley. *Statistical Reasoning with Imprecise Probabilities*. Chapman & Hall, 1991.
- 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>.

Appendix: