

Causal Statistical Decision Theory|What are interventions?

David Johnston

January 26, 2021

Contents

0.1	Theories of causal inference	3
0.1.1	Probability Theory	6
0.1.2	Product Notation	8
0.1.3	String Diagrams	9
0.1.4	Random Variables	15
1	Chapter 3: See-do models	29
1.1	Definition	32
1.1.1	D-causation	35
1.1.2	D-causation vs Limited Unresponsiveness	36
1.1.3	Properties of D-causation	39
1.1.4	Decision sequences and parallel decisions	39
1.1.5	Residual dependence on observations	40
1.1.6	Causal questions and decision functions	41
1.2	Existence of counterfactuals	43
2	Chapter 4: See-do models compared to causal graphical models and potential outcomes	45

0.1 Theories of causal inference

Feedback start here

Beginning in the 1930s, a number of associations between cigarette smoking and lung cancer were established: on a population level, lung cancer rates rose rapidly alongside the prevalence of cigarette smoking. Lung cancer patients were far more likely to have a smoking history than demographically similar individuals without cancer and smokers were around 40 times as likely as demographically similar non-smokers to go on to develop lung cancer. In laboratory experiments, cells which were introduced to tobacco smoke developed *ciliastasis*, and mice exposed to cigarette smoke tars developed tumors(Proctor, 2012). Nevertheless, until the late 1950s, substantial controversy persisted over the question of whether the available data was sufficient to establish that smoking cigarettes *caused* lung cancer. Cigarette manufacturers famously argued against

any possible connection (Oreskes and Conway, 2011) and Roland Fisher in particular argued that the available data was not enough to establish that smoking actually caused lung cancer (Fisher, 1958). Today, it is widely accepted that cigarettes do cause lung cancer, along with other serious conditions such as vascular disease and chronic respiratory disease (World Health Organisation, 2018; Wiblin, 2016).

The question of a causal link between smoking and cancer is a very important one. Individuals who enjoy smoking (or think they might) may wish to avoid smoking if cigarettes pose a severe health risk, so they are interested in knowing whether or not it is so. Potential investors in cigarette manufacturers want to know if the product they are backing is likely to see limited adoption due to health concerns. People holding investments in cigarette manufacturing firms want the world to be such that cigarettes do not pose a substantial health risk, as this increases the value of their investment. Governments and organisations with a responsibility for public health may see themselves as having responsibility to discourage smoking as much as possible if smoking is severely detrimental to health. The costs and benefits of poor decisions about smoking are large: 8 million annual deaths are attributed to cigarette-caused cancer and vascular disease in 2018 (World Health Organisation, 2018) while global cigarette sales were estimated at US\$711 million in 2020, while (Statista, 2020) (a figure which might be substantially larger if cigarettes were not widely believed to be harmful).

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

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

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

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

Theories of causal inference exist to enable formal rather than ad-hoc reasoning about causal questions. Instead of posing informal causal question and

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

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

Defining causal models is difficult. In general, applied theories of causal inference posit that:

1. “ X_i causes X_j ” means that there exist different *ideal interventions* that result in different values of X_i , hold other “causally sufficient” variables constant, do not directly affect X_j but nonetheless entail different values of X_j
2. “ X_i causes X_j ” means that the *counterfactual value* of X_j would be different “if X_i had taken a different value”

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

“Ideal interventions” themselves are difficult to define. The structural model approach of Pearl (2009) defines ideal interventions in terms of “causally sufficient models”. However, most attempts to formalise this definition end up being circular. For example:

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

1

Circularity is a recognised problem with interventional definitions of causation (Woodward, 2016). In Section ??, I further show that assuming ideal interventions always exist leads to counterintuitive conclusions. An alternative approach is to designate certain real-world events – such as flipping coins,

¹Weaker conditions for causal sufficiency are possible, but they are still premised on causal relationships, so circularity stands (Shpitser and Pearl, 2008).

querying random number generators and so forth – as prototypical “ideal interventions”. This approach is rather inflexible, and refuses to offer answers to causal questions that don’t happen to have involve just the right kinds of real world events, typically randomised experiments. However, many causal questions do have apparent answers (Pearl, 2018), and even when gold-standard randomised experimental data is available, it may not permit answers to the original questions of interest (Deaton and Cartwright, 2018; Heckman, 1991).

The difficulty in defining “ideal interventions” is not unprecedented. It is also difficult to provide an account of what it means for data to be “distributed according to probability distribution \mathbb{P} ” (Hájek, 2019), but the usefulness of using probability distributions to model data is widely accepted.

Causal statistical decision theory (CSDT) is a theory of “causal questions” that does not depend on an underlying theory of causation. 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 show that it is also possible to formalise the problem of determining *counterfactual consequences* without appealing to an underlying theory of causation.

A key feature of CSDT is the importance of the *option set*, which is the set of decisions, acts or counterfactual actions under consideration in a given problem. A great deal of work on causal inference defines with the option set implicitly, possibly also relying on default choices such as that of “hard intervention”. We argue that:

- 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 can affect the difficulty of causal questions
- Hard interventions are not a good choice for default option sets

Theorem 0.1.1 (Representation).

Representation theorem: can uniquely define kernel $P^{X|Y}$ with $P^{Z|Y}$ and $P^{X|ZY}$

0.1.1 Probability Theory

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} , $\cup_{i \in \mathbb{N}} B_i \in \mathcal{A}$

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

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 .

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

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

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

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

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

A particular probability measure we will often discuss is the *Dirac measure*. For any $x \in X$, the Dirac measure $\delta_x \in \Delta(\mathcal{X})$ is the probability measure where $\delta_x(A) = 0$ if $x \notin A$ and $\delta_x(A) = 1$ if $x \in A$.

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

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

Extending the subscript notation, for $\mathbb{C} : X \times Y \rightarrow \Delta(\mathcal{Z})$ and $x \in X$ we will write $\mathbb{C}_{x,\cdot}$ for the “curried” map $y \mapsto \mathbb{C}_{x,y}$. If \mathbb{C} is a Markov kernel with respect to $(X \times Y, \mathcal{X} \otimes \mathcal{Y}), (Z, \mathcal{Z})$ then it is straightforward to show that $\mathbb{C}_{x,\cdot}$ is a Markov kernel with respect to $(Y, \mathcal{Y}), (Z, \mathcal{Z})$.

This yields the notational conventions for arbitrary kernel \mathbb{C} :

- \mathbb{C} with no subscripts is a Markov kernel
- $\mathbb{C}_{\cdot,a,b}$ with at least one \cdot subscript is a Markov kernel
- \mathbb{C}_y with no \cdot subscripts is a probability measure

The map $x \mapsto \mathbb{M}_x$ is of type $E \rightarrow \Delta(\mathcal{F})$. We will abuse notation somewhat to write $\mathbb{M} : E \rightarrow \Delta(\mathcal{F})$. In this sense, we view Markov kernels as maps from elements of E to probability measures on \mathcal{F} . This is simply a convention that helps us to think about constructions involving Markov kernels, and it is equally valid to view Markov kernels as maps from elements of \mathcal{F} to measurable functions $E \rightarrow [0, 1]$, a view found in Clerc et al. (2017), or simply in terms of their definition above.

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

Two Markov kernels $\mathbb{M}X \rightarrow \Delta(\mathcal{Y})$ and $\mathbb{N} : X \rightarrow \Delta(\mathcal{Y})$ are equal iff for all $x \in X$, $A \in \mathcal{Y}$

$$\mathbb{M}_x(A) = \mathbb{N}_x(A) \quad (1)$$

We will typically be more concerned with “almost sure” equality than exact equality, which will be defined later.

0.1.2 Product Notation

Probability measures, Markov kernels and measurable functions can be combined to yield new probability measures, Markov kernels or measurable functions. Given $\mu \in \Delta(\mathcal{X})$, $\mathbb{T} : Y \rightarrow T$, $\mathbb{M} : X \rightarrow \Delta(\mathcal{Y})$ and $\mathbb{N} : Y \rightarrow \Delta(\mathcal{Z})$ define:

The **measure-kernel** product $\mu\mathbb{M} : \mathcal{Y} \rightarrow [0, 1]$ where for all $A \in \mathcal{Y}$,

$$\mu\mathbb{M}(A) := \int_X \mathbb{M}_x(A) d\mu(x) \quad (2)$$

The **kernel-function** product $\mathbb{M}\mathbb{T} : X \rightarrow T$ where for all $x \in X$:

$$\mathbb{M}\mathbb{T}(x) := \int_Y T(y) d\mathbb{M}_x(y) \quad (3)$$

The **kernel-kernel** product $\mathbb{M}\mathbb{N} : X \rightarrow \Delta(\mathcal{Z})$ where for all $x \in X$, $A \in \mathcal{Z}$:

$$(\mathbb{M}\mathbb{N})_x(A) := \int_Y \mathbb{N}_y(A) d\mathbb{M}_x(y) \quad (4)$$

All kernel products are associative (Çinlar, 2011). An intuition for this notation can be gained from thinking of probability measures $\mu \in \Delta(\mathcal{X})$ as row vectors, Markov kernels \mathbb{M}, \mathbb{N} as matrices and measurable functions $\mathbb{T} : Y \rightarrow T$ as column vectors and kernel products are vector-matrix and matrix-matrix products. If the X, Y, Z and T are discrete spaces then this analogy is precise.

Finally, the **tensor product** $\mathbb{M} \otimes \mathbb{N} : X \times Y \rightarrow \Delta(\mathcal{Y} \otimes \mathcal{Z})$ yields the kernel that applies \mathbb{M} and \mathbb{N} “in parallel”. For all $x \in X$, $y \in Y$, $G \in \mathcal{Y}$ and $H \in \mathcal{Z}$:

$$(\mathbb{M} \otimes \mathbb{N})_{x,y}(G \times H) := \mathbb{M}_x(G) \mathbb{N}_y(H) \quad (5)$$

0.1.3 String Diagrams

Some constructions are unwieldy in product notation; for example, given $\mu \in \Delta(\mathcal{E})$ and $\mathbb{M} : E \rightarrow (\mathcal{F})$, it is not straightforward to write an expression using kernel products and tensor products that represents the “joint distribution” given by $A \times B \mapsto \int_A \mathbb{M}(x; B) d\mu$.

An alternative notation known as *string diagrams* provides greater expressive capability than product notation while being more visually clear than integral notation. Cho and Jacobs (2019) provides an extensive introduction to string diagram notation for probability theory.

Key features of string diagrams include:

- String diagrams as they are used in this work can always be interpreted as a mixture of kernel-kernel products and tensor products of Markov kernels
- String diagrams are the subject of a coherence theorem: two string diagrams that differ only by planar deformation are always equal (Selinger, 2010). This also holds for a number of additional transformations detailed below
 - Informally, diagrams that look like they should be the same are in fact the same

Elements of string diagrams

The basic elements of a string diagram are Markov kernels. Diagrams representing Markov kernels can be assembled into larger diagrams by taking regular products or tensor products.

Indiscrete spaces play a key role in string diagrams. An indiscrete space is any one element measurable space $(\{*\}, \{\emptyset, \{*\}\})$ which admits the unique probability measure $\mu : \{\emptyset, \{*\}\} \rightarrow (0, 1)$ given by $\mu(\emptyset) = 0$, $\mu(\{*\}) = 1$. Any probability measure $\mu \in \Delta(\mathcal{X})$ can be interpreted as a Markov kernel $\mu' : \{*\} \rightarrow \Delta(\mathcal{X})$ where $\mu'_* = \mu$ (note that $*$ is the *only* argument μ' can be given).

A Markov kernel $\mathbb{M} : X \rightarrow \Delta(\mathcal{Y})$ can always be represented as a rectangular box with input and output wires labeled with the relevant spaces:

$$X \text{ --- } \boxed{\mathbb{M}} \text{ --- } Y \quad (6)$$

Note that we will later substitute labelling wires with spaces for labelling them with random variable names.

Probability measures are kernels with an indiscrete domain $\mu \in \Delta(\mathcal{X})$ can be written as triangles:

$$\triangleleft \mu \text{ --- } X \quad (7)$$

Note that Eq 7 technically represents a Markov kernel $\mu' : \{*\} \rightarrow \Delta(\mathcal{X})$, but for our purposes this distinction isn't practically important.

We do *not* define kernel-function products for string diagrams. While kernel-kernel products always yield Markov kernels as a result, and measure-kernel products can be reinterpreted as kernel-kernel products, kernel-function products do not admit such a reinterpretation. Cho and Jacobs (2019) defines the operation of *conditioning* using kernel-function products, but this will take extra work to incorporate into our notation which hasn't yet been done.

Elementary operations Kernel-kernel products have a visually similar representations in string diagram notation to the previously introduced product notation. Given $\mathbb{M} : X \rightarrow \Delta(\mathcal{Y})$ and $\mathbb{N} : Y \rightarrow \Delta(\mathcal{Z})$, we have

$$\mathbb{M}\mathbb{N} := X \text{ --- } \boxed{\mathbb{M}} \text{ --- } \boxed{\mathbb{N}} \text{ --- } Z \quad (8)$$

For $\mu \in \Delta(\mathcal{E})$,

$$\mu\mathbb{M} := \triangleleft \mu \text{ --- } \boxed{\mathbb{M}} \text{ --- } Z \quad (9)$$

Tensor products in string diagram notation are represented by vertical juxtaposition. For $\mathbb{O} : Z \rightarrow \Delta(\mathcal{W})$:

$$\mathbb{M} \otimes \mathbb{O} := \begin{array}{c} X \text{ --- } \boxed{\mathbb{M}} \text{ --- } Y \\ Z \text{ --- } \boxed{\mathbb{O}} \text{ --- } W \end{array} \quad (10)$$

A space X is identified with the identity kernel $\text{Id}^X : X \rightarrow \Delta(\mathcal{X})$, $x \mapsto \delta_x$. A bare wire represents an identity kernel or, equivalently, the space given by its labels:

$$\text{Id}^X := X \text{ ————— } X \quad (11)$$

Product spaces $X \times Y$ are identified with tensor products of identity kernels $X \times Y \cong \mathbb{I}^X \otimes \mathbb{I}^Y$. These can be represented either by two parallel wires or by a single wire equipped with appropriate labels:

$$X \times Y \cong \text{Id}^X \otimes \text{Id}^Y := \begin{array}{c} X \text{ --- } X \\ Y \text{ --- } Y \end{array} \quad (12)$$

$$= X \times Y \text{ ————— } X \times Y \quad (13)$$

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:

$$X \text{ --- } \boxed{\mathbb{L}} \text{ --- } Y \quad Z \quad (14)$$

$$\equiv \quad (15)$$

$$X \text{ --- } \boxed{\mathbb{L}} \text{ --- } Y \times Z \quad (16)$$

Markov kernels with special notation A number of Markov kernels are given special notation distinct from the generic “box” above. This notation facilitates intuitive visual representation.

As has already been noted, the identity kernel $\mathbf{Id} : X \rightarrow \Delta(X)$ maps a point x to the measure δ_x that places all mass on the same point:

$$\mathbf{Id} : x \mapsto \delta_x \equiv X \text{ --- } X \quad (17)$$

The identity kernel is an identity under left and right products:

$$(\mathbb{K}\mathbf{Id})_w(A) = \int_X \mathbf{Id}_x(A) d\mathbb{K}_w(x) \quad (18)$$

$$= \int_X \delta_x(A) d\mathbb{K}_w(x) \quad (19)$$

$$= \int_A d\mathbb{K}_w(x) \quad (20)$$

$$= \mathbb{K}_w(A) \quad (21)$$

$$(\mathbf{Id}\mathbb{K})_w(A) = \int_X \mathbb{K}_x(A) d\mathbf{Id}_w(x) \quad (22)$$

$$= \int_X \mathbb{K}_x(A) d\delta_w(x) \quad (23)$$

$$= \mathbb{K}_w(A) \quad (24)$$

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

$$\Upsilon : x \mapsto \delta_{(x,x)} \equiv X \text{ --- } \begin{matrix} X \\ X \end{matrix} \quad (25)$$

The copy map “copies” its arguments to kernels or under the right product:

$$\int_{(\cdot)} X \times X \mathbb{K}_{x',x''}(A) d\Upsilon_x(x',x'') = \int_{(\cdot)} X \times X \mathbb{K}_{x',x''}(A) d\delta_{(x,x)}(x',x'') \quad (26)$$

$$= \mathbb{K}_{x,x}(A) \quad (27)$$

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

$$\sigma := (x, y) \rightarrow \delta_{(y, x)} \equiv \begin{matrix} Y \\ X \end{matrix} \succ \begin{matrix} X \\ Y \end{matrix} \quad (28)$$

Under products are taken with the swap map, arguments are interchanged. For $\mathbb{K} : X \times Y \rightarrow \Delta(\mathcal{Z})$ and $\mathbb{L} : Z \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$, $A \in \mathcal{X}$, $B \in \mathcal{Y}$:

$$(\sigma \mathbb{K})_{y, x}(A) = \int_{(X \times Y)} \mathbb{K}_{x', y'}(A) d\sigma_{(y, x)}(x', y') = \int_{(X \times Y)} \mathbb{K}_{x', y'}(A) d\delta_{(x, y)}(x', y') \quad (29)$$

$$= \mathbb{K}_{x, y}(A) \quad (30)$$

$$(\mathbb{L} \sigma)_z(B \times A) = \int_{X \times Y} \sigma_{x', y'}(B \times A) d\mathbb{L}_z(x', y') \quad (31)$$

$$= \int_{X \times Y} \delta_{(y', x')}(B \times A) d\mathbb{L}_z(x', y') \quad (32)$$

$$= \mathbb{L}_z(A \times B) \quad (33)$$

The discard map $* : X \rightarrow \Delta(\{*\})$ maps every input to δ_* , the unique probability measure on the indiscrete set $\{\emptyset, \{*\}\}$.

$$* : x \mapsto \delta_* \equiv X \longrightarrow * \quad (34)$$

Any measurable function $g : W \rightarrow X$ has an associated Markov kernel $\mathbb{F}^g : W \rightarrow \Delta(\mathcal{X})$ given by $\mathbb{F}^g : w \mapsto \delta_{g(w)}$. Given a probability measure $\mu \in \Delta(\mathcal{W})$, μg is a measure-function product while $\mu \mathbb{F}^g$ is commonly called the pushforward measure $g_{\#} \mu$. We will generalise this slightly to the notion of *pushforward kernels*.

Definition 0.1.2 (Kernel associated with a function). Given a measurable function $g : W \rightarrow X$, define the function induced kernel $\mathbb{F}^g : W \rightarrow \Delta(\mathcal{X})$ to be the the Markov kernel $w \mapsto \delta_{g(w)}$ for all $w \in W$.

Definition 0.1.3 (Pushforward kernel). Given a kernel $\mathbb{M} : V \rightarrow \Delta(\mathcal{W})$ and a measurable function $g : W \rightarrow X$, the *pushforward kernel* $g_{\#} \mathbb{M} : V \rightarrow \Delta(\mathcal{X})$ is the kernel $g_{\#} \mathbb{M}$ such that $(g_{\#} \mathbb{M})_a(B) = \mathbb{M}_a(g^{-1}(B))$ for all $a \in V$, $B \in \mathcal{X}$.

Lemma 0.1.4 (Pushforward kernels are functional kernel products). *Given a kernel $\mathbb{M} : V \rightarrow \Delta(\mathcal{W})$ and a measurable function $g : W \rightarrow X$, $g_{\#} \mathbb{M} = \mathbb{M} \mathbb{F}^g$.*

Proof. for any $a \in V$, $B \in \mathcal{X}$:

$$(\mathbb{M}\mathbb{F}^g)_a(B) = \int_W \delta_{g(y)}(B) d\mathbb{M}_a(y) \quad (35)$$

$$= \int_W \delta_y(g^{-1}(B)) d\mathbb{M}_a(y) \quad (36)$$

$$= \int_{g^{-1}(B)} d\mathbb{M}_a(y) \quad (37)$$

$$= (g_{\#}\mathbb{M})_a(B) \quad (38)$$

□

Working With String Diagrams

todo:

- Infinite copy map
- De Finetti's representation theorem

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

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

$$\begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} := \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} \quad (39)$$

$$\begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array}^* = \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array}^* = \text{---} \quad (40)$$

$$\begin{array}{c} \text{X} \text{---} \text{---} \text{---} \\ \text{X} \text{---} \text{---} \text{---} \end{array} = \begin{array}{c} \text{---} \text{---} \text{---} \\ \text{---} \text{---} \text{---} \end{array} \quad (41)$$

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

$$\text{---} \boxed{\mathbb{A}} \text{---} * = \text{---} * \quad (42)$$

Combining 40 and 42 we can derive the following: integrating $\mathbb{A} : X \rightarrow \Delta(\mathcal{Y})$ with respect to $\mu \in \Delta(\mathcal{X})$ and then discarding the output of \mathbb{A} leaves us with μ :

$$\begin{array}{c} \triangleleft \mu \text{---} \end{array} \text{---} \boxed{\mathbb{A}} \text{---} * = \begin{array}{c} \triangleleft \mu \text{---} \end{array} \text{---} * = \begin{array}{c} \triangleleft \mu \text{---} \end{array} \quad (43)$$

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

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

$$X \times Y \text{---} * = \begin{array}{c} X \text{---} * \\ Y \text{---} * \end{array} \quad (44)$$

$$X \times Y \text{---} \begin{array}{c} X \times Y \\ X \times Y \end{array} = \begin{array}{c} X \\ Y \end{array} \text{---} \begin{array}{c} X \\ Y \end{array} \quad (45)$$

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

$$\text{---} \boxed{\mathbb{A}} \text{---} \begin{array}{c} \text{---} \\ \text{---} \end{array} = \begin{array}{c} \text{---} \boxed{\mathbb{A}} \text{---} \\ \text{---} \boxed{\mathbb{A}} \text{---} \end{array} \quad (46)$$

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

Definition 0.1.5 (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 $x \in E$, $B \in \mathcal{F}$.

Theorem 0.1.6 (Copy map commutes for deterministic kernels (Fong, 2013)). Equation 46 holds iff \mathbb{A} is deterministic.

Examples

Given $\mu \in \Delta(X)$, $\mathbb{K} : X \rightarrow \Delta(Y)$, $A \in \mathcal{X}$ and $B \in \mathcal{Y}$:

$$A \times B \mapsto \int_A \mathbb{K}(x; B) d\mu(x) \quad (47)$$

$$\equiv \quad (48)$$

$$\mu^\vee(\mathbf{Id}_X \otimes \mathbb{K}) \quad (49)$$

$$\equiv \quad (50)$$

$$\begin{array}{c} \text{---} \mu \text{---} \text{---} X \\ \text{---} \text{---} \boxed{\mathbb{K}} \text{---} Y \end{array} \quad (51)$$

Cho and Jacobs (2019) calls this operation “integrating \mathbb{K} with respect to μ ”.

Given $\nu \in \Delta(\mathcal{X} \otimes \mathcal{Y})$, define the marginal $\nu^Y \in \Delta(\mathcal{Y}) : B \mapsto \mu(X \times B)$ for $B \in \mathcal{Y}$. Say that ν^Y is obtained by marginalising over “ X ” (a notion that can be made more precise by assigning names to wires). Then

$$\nu(* \otimes \text{Id}^Y) = \begin{array}{c} \text{---} \nu \text{---} * \\ \text{---} \text{---} Y \end{array} \quad (52)$$

$$\nu(* \otimes \text{Id}^Y)(B) := \nu(* \otimes \text{Id}^Y)(B \times \{*\}) \quad (53)$$

$$= \int_{X \times Y} \text{Id}_y^Y(B) *_{\{*\}} d\nu(x, y) \quad (54)$$

$$= \int_{X \times Y} \delta_y(B) \delta_{\{*\}} d\nu(x, y) \quad (55)$$

$$= \int_{X \times B} d\nu(x, y) \quad (56)$$

$$= \nu(X \times B) \quad (57)$$

$$= \nu^Y(B) \quad (58)$$

Thus the action of the erasing wire “ X ” is equivalent to marginalising over “ X ”.

Consider the result of marginalising 51 over “ X ”:

$$\begin{array}{c} \text{---} \mu \text{---} \text{---} * \\ \text{---} \text{---} \boxed{\mathbb{A}} \text{---} Y \end{array} \quad (59)$$

$$= \begin{array}{c} \text{---} \mu \text{---} \boxed{\mathbb{A}} \text{---} Y \end{array} \quad (60)$$

0.1.4 Random Variables

The summary of this section is:

- Random variables are usually defined as measurable functions on a *probability space*
- It's possible to define them as measurable functions on a *Markov kernel space* instead
- It is useful to label wires with random variable names instead of names of spaces

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

We use the same approach here, with a couple of additions. We are interested in variables whose outcomes depend both on random processes and decisions. Suppose that given a particular distribution over decision variables, a probability distribution over the decision variables and random variables is obtained. Such a model is described by a Markov kernel rather than a probability distribution. We therefore investigate *Markov kernel spaces*.

In the graphical notation that we are using, random variables can be thought of as a means of assigning unambiguous names to each wire in a set of diagrams. In order to do this, it is necessary to suppose that all diagrams in the set describe properties of an *ambient Markov kernel* or *ambient probability measure*. Consider the following example with the ambient probability measure $\mu \in \Delta(\mathcal{X} \otimes \mathcal{X})$. Suppose we have a Markov kernel $\mathbb{K} : X \rightarrow \Delta(\mathcal{X})$ such that the following holds:

$$\triangleleft_{\mu} \frac{X}{X} = \triangleleft_{\mu} * \boxed{\mathbb{K}} \frac{X}{X} \quad (61)$$

Suppose that we also assign the names X_1 to the upper output wire and X_2 to the lower output wire in the diagram of μ :

$$\triangleleft_{\mu} \frac{X_1}{X_2} \quad (62)$$

Then it seems sensible to call \mathbb{K} “the probability of X_2 given X_1 ”. We will make this precise, and it will match the usual notion of the probability of one variable given another (see Çinlar (2011) for a definition of this usual notion).

Definition 0.1.7 (Probability space, Markov kernel space). A *Markov kernel space* $(\mathbb{K}, (D, \mathcal{D}), (\Omega, \mathcal{F}))$ is a Markov kernel $\mathbb{K} : D \rightarrow \Delta(\mathcal{D} \otimes \mathcal{F})$, called the *ambient kernel*, along with the sample space (Ω, \mathcal{F}) and the domain (D, \mathcal{D}) . We suppose that \mathbb{K} is such that there exists a *fundamental kernel* \mathbb{K}_0 satisfying

$$\mathbb{K} := \text{---} \boxed{\mathbb{K}_0} \text{---} \quad (63)$$

For brevity, we will omit the σ -algebras in further definitions of Markov kernel spaces: (\mathbb{K}, D, Ω) .

A *probability space* $(\mathbb{P}, \Omega, \mathcal{F})$ is a probability measure $\mathbb{P} : \Delta(\Omega)$, which we call the *ambient measure*, along with the *sample space* Ω and the *events* \mathcal{F} . A probability space is equivalent to a Markov kernel space with domain $D = \{*\}$ - note that $\Omega \times \{*\} \cong \Omega$.

Definition 0.1.8 (Random variable). Given a Markov kernel space (\mathbb{K}, D, Ω) , a random variable \mathbf{X} is a measurable function $\Omega \times D \rightarrow E$ for arbitrary measurable E .

Definition 0.1.9 (Domain variable). Given a Markov kernel space (\mathbb{K}, D, Ω) , the *domain variable* $\mathbf{D} : \Omega \times D \rightarrow D$ is the distinguished random variable $\mathbf{D} : (x, d) \mapsto d$.

Unlike random variables on probability spaces, random variables on Markov kernel spaces do not generally have unique marginal distributions. An analogous operation of *marginalisation* can be defined, but the result is generally a Markov kernel. We will define marginalisation via coupled tensor products.

Definition 0.1.10 (Coupled tensor product \otimes). Given two Markov kernels \mathbb{M} and \mathbb{N} or functions f and g with shared domain E , let $\mathbb{M} \otimes \mathbb{N} := \vee(\mathbb{M} \otimes \mathbb{N})$ and $f \otimes g := \vee(f \otimes g)$ where these expressions are interpreted using standard product notation. Graphically:

$$\mathbb{M} \otimes \mathbb{N} := \begin{array}{c} E \text{---} \begin{array}{c} \boxed{\mathbb{M}} \text{---} \mathbf{X} \\ \boxed{\mathbb{N}} \text{---} \mathbf{Y} \end{array} \end{array} \quad (64)$$

$$f \otimes g := \begin{array}{c} E \text{---} \begin{array}{c} \triangle f \\ \triangle g \end{array} \end{array} \quad (65)$$

The operation denoted by \otimes is associative (Lemma 0.1.11), so we can without ambiguity write $f \otimes g \otimes h = (f \otimes g) \otimes h = f \otimes (g \otimes h)$ for finite groups of functions or Markov kernels sharing a domain.

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

Lemma 0.1.11 (\otimes is associative). For Markov kernels $\mathbb{L} : E \rightarrow \delta(\mathcal{F})$, $\mathbb{M} : E \rightarrow \delta(\mathcal{G})$ and $\mathbb{N} : E \rightarrow \delta(\mathcal{H})$, $(\mathbb{L} \otimes \mathbb{M}) \otimes \mathbb{N} = \mathbb{L} \otimes (\mathbb{M} \otimes \mathbb{N})$.

Proof.

$$\mathbb{L} \otimes (\mathbb{M} \otimes \mathbb{N}) = \begin{array}{c} E \text{ --- } \begin{array}{l} \boxed{\mathbb{L}} \text{ --- } F \\ \boxed{\mathbb{M}} \text{ --- } G \\ \boxed{\mathbb{N}} \text{ --- } H \end{array} \end{array} \quad (66)$$

$$= \begin{array}{c} E \text{ --- } \begin{array}{l} \boxed{\mathbb{L}} \text{ --- } F \\ \boxed{\mathbb{M}} \text{ --- } G \\ \boxed{\mathbb{N}} \text{ --- } H \end{array} \end{array} \quad (67)$$

$$= (\mathbb{L} \otimes \mathbb{M}) \otimes \mathbb{N} \quad (68)$$

This follows directly from Equation 39. \square

Definition 0.1.12 (Marginal distribution, marginal kernel). Given a probability space $(\mathbb{P}, \Omega, \mathcal{F})$ and the random variable $\mathbf{X} : \Omega \rightarrow G$ the *marginal distribution* of \mathbf{X} is the probability measure $\mathbb{P}^{\mathbf{X}} := \mathbb{P}\mathbb{F}^{\mathbf{X}}$.

See Lemma 0.1.4 for the proof that this matches the usual definition of marginal distribution.

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

Definition 0.1.13 (Joint distribution, joint kernel). Given a probability space $(\mathbb{P}, \Omega, \mathcal{F})$ and the random variables $\mathbf{X} : \Omega \rightarrow G$ and $\mathbf{Y} : \Omega \rightarrow H$, the *joint distribution* of \mathbf{X} and \mathbf{Y} , $\mathbb{P}^{\mathbf{X}\mathbf{Y}} \in \Delta(\mathcal{G} \otimes \mathcal{H})$, is the marginal distribution of $\mathbf{X} \otimes \mathbf{Y}$. That is, $\mathbb{P}^{\mathbf{X}\mathbf{Y}} := \mathbb{P}\mathbb{F}^{\mathbf{X} \otimes \mathbf{Y}}$.

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

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

Joint distributions and kernels have a nice visual representation, as a result of Lemma 0.1.14 which follows.

Lemma 0.1.14 (Product marginalisation interchange). *Given two functions, the kernel associated with their coupled product is equal to the coupled product of the kernels associated with each function.*

Given $\mathbf{X} : \Omega \rightarrow G$ and $\mathbf{Y} : \Omega \rightarrow H$, $\mathbb{F}^{\mathbf{X} \otimes \mathbf{Y}} = \mathbb{F}^{\mathbf{X}} \otimes \mathbb{F}^{\mathbf{Y}}$

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

$$\mathbb{F}^{\mathbf{X} \otimes \mathbf{Y}}(a; B \times C) = \delta_{\mathbf{X}(a), \mathbf{Y}(a)}(B \times C) \quad (69)$$

$$= \delta_{\mathbf{X}(a)}(B) \delta_{\mathbf{Y}(a)}(C) \quad (70)$$

$$= (\delta_{\mathbf{X}(a)} \otimes \delta_{\mathbf{Y}(a)})(B \times C) \quad (71)$$

$$= \mathbb{F}^{\mathbf{X}} \otimes \mathbb{F}^{\mathbf{Y}} \quad (72)$$

Equality follows from the monotone class theorem. \square

Corollary 0.1.15. *Given a Markov kernel space (\mathbb{K}, Ω, D) and random variables $\mathsf{X} : \Omega \times D \rightarrow X$, $\mathsf{Y} : \Omega \times D \rightarrow Y$, the following holds:*

$$D \text{ --- } \boxed{\mathbb{K}^{\mathsf{XY}|D}} \text{ --- } \begin{matrix} X \\ Y \end{matrix} = D \text{ --- } \boxed{\mathbb{K}} \text{ --- } \begin{matrix} \boxed{\mathbb{F}^{\mathsf{X}}} \text{ --- } X \\ \boxed{\mathbb{F}^{\mathsf{Y}}} \text{ --- } Y \end{matrix} \quad (73)$$

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

Definition 0.1.16 (Wire labels - joint kernels). Suppose we have a Markov kernel space (\mathbb{K}, D, Ω) , random variables $\mathsf{X} : \Omega \times D \rightarrow X$, $\mathsf{Y} : \Omega \times D \rightarrow Y$ and a Markov kernel $\mathbb{L} : D \rightarrow \Delta(\mathcal{X} \times \mathcal{Y})$. The following *output labelling* of \mathbf{L} :

$$D \text{ --- } \boxed{\mathbb{L}} \text{ --- } \begin{matrix} \mathsf{X} \\ \mathsf{Y} \end{matrix} \quad (74)$$

is *valid* iff

$$\mathbb{L} = \mathbb{K}_{\mathsf{XY}|D} \quad (75)$$

and

$$D \text{ --- } \boxed{\mathbb{L}} \text{ --- } \begin{matrix} \mathsf{X} \\ * \end{matrix} = \mathbb{K}^{\mathsf{X}|D} \quad (76)$$

and

$$D \text{ --- } \boxed{\mathbb{L}} \text{ --- } \begin{matrix} * \\ \mathsf{Y} \end{matrix} = \mathbb{K}^{\mathsf{Y}|D} \quad (77)$$

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

$$D \text{ --- } \boxed{\mathbb{L}} \text{ --- } \begin{matrix} W \\ Z \end{matrix} \times Y \quad (78)$$

For *this* diagram, properties 76 and 77 do not hold, even though 75 does.

Lemma 0.1.17 (Output label assignments exist). *Given Markov kernel space (\mathbb{K}, D, Ω) , random variables $\mathsf{X} : \Omega \times D \rightarrow X$ and $\mathsf{Y} : \Omega \times D \rightarrow Y$ then there exists a diagram of $\mathbb{L} := \mathbb{K}^{\mathsf{XY}|D}$ with a valid output labelling assigning X and Y to the output wires.*

Proof. By definition, \mathbb{L} has signature $D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$. Thus, by the rule that tensor product spaces can be represented by parallel wires, we can draw

$$D - \boxed{\mathbb{L}} - \begin{array}{c} X \\ Y \end{array} \quad (79)$$

By Corollary 0.1.15, we have

$$D - \boxed{\mathbb{L}} - \begin{array}{c} X \\ Y \end{array} = D - \boxed{\mathbb{K}} - \left(\begin{array}{c} \boxed{\mathbb{F}^X} - X \\ \boxed{\mathbb{F}^Y} - Y \end{array} \right) \quad (80)$$

Therefore

$$D - \boxed{\mathbb{K}} - \left(\begin{array}{c} \boxed{\mathbb{F}^X} - X \\ \boxed{\mathbb{F}^Y} - * \end{array} \right) = \mathbb{K}\mathbb{F}^X \quad (81)$$

$$= \mathbb{K}^{X|D} \quad (82)$$

$$D - \boxed{\mathbb{K}} - \left(\begin{array}{c} \boxed{\mathbb{F}^X} - * \\ \boxed{\mathbb{F}^Y} - Y \end{array} \right) = \mathbb{K}\mathbb{F}^Y \quad (83)$$

$$= \mathbb{K}^{Y|D} \quad (84)$$

□

In all further work, wire labels will be used without special colouring.

Definition 0.1.18 (Disintegration). Given a probability space $(\mathbb{P}, \Omega, \mathcal{F})$, and random variables X and Y , we say that $\mathbb{M} : E \rightarrow \Delta(\mathcal{F})$ is a *Y on X disintegration* of \mathbb{P} iff

$$\triangleleft \mathbb{P}^{XY} = \triangleleft \mathbb{P}^X \begin{array}{c} * \\ \boxed{\mathbb{M}} - \begin{array}{c} X \\ Y \end{array} \end{array} \quad (85)$$

\mathbb{M} is a version of $\mathbb{P}^{Y|X}$, “the probability of Y given X ”. Let $\mathbb{P}^{\{Y|X\}}$ be the set of all kernels that satisfy 85 and $\mathbb{P}^{Y|X}$ an arbitrary member of $\mathbb{P}^{Y|X}$.

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

$$- \boxed{\mathbb{K}^{YX|D}} - \begin{array}{c} X \\ Y \end{array} = \begin{array}{c} \boxed{\mathbb{K}^{YX|D}} - * \\ \boxed{\mathbb{M}} - \begin{array}{c} X \\ Y \end{array} \end{array} \quad (86)$$

Write $\mathbb{K}^{\{Y|XD\}}$ for the set of kernels satisfying 86 and $\mathbb{K}^{Y|XD}$ for an arbitrary member of $\mathbb{K}^{\{Y|XD\}}$.

Definition 0.1.19 (Wire labels – input). An input wire is *connected* to an output wire if it is possible to trace a path from the start of the input wire to the end of the output wire without passing through any boxes, erase maps or right facing triangles.

If an input wire is connected to an output wire and that output wire has a valid label X , then it is valid to label the input wire with X .

For example, if the following are valid output labels with respect to (\mathbb{P}, Ω) :

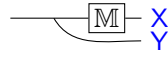

(87)

i.e. if $\mathbb{L} \in \mathbb{P}^{XY|Y}$, then the following is a valid input label:


(88)

An input wire in a diagram for \mathbb{M} may be labeled X *if and only if* copy and identity maps can be inserted to yield a diagram in which the input wire labeled X is connected to an output wire with valid label X .

So, if $\mathbb{M} \in \mathbb{P}^{X|Y}$, then it is straightforward to show that


 $\in \mathbb{P}^{XY|Y}$
(89)

and hence the output labels are valid. Diagram 89 is constructed by taking the product of the copy map with $\mathbb{M} \otimes \text{Id}$. Thus it is valid to label \mathbb{M} with


(90)

Lemma 0.1.20 (Labeling of disintegrations). *Given a kernel space (\mathbb{K}, D, Ω) , random variables X and Y , domain variable D and disintegration $\mathbb{L} \in \mathbb{K}^{Y|XD}$, there is a diagram of \mathbb{L} with valid input labels X and D and valid output label Y .*

Proof. Note that for any variable $W : \Omega \times D \rightarrow W$ and the domain variable

$D : \Omega \times D \rightarrow D$ we have by definition of \mathbb{K} :

$$\text{---} \boxed{\mathbb{K}^{WD|D}} \text{---} \begin{matrix} W \\ D \end{matrix} = \begin{matrix} & & \boxed{\mathbb{F}^W} & W \\ & \boxed{\mathbb{K}_0} & \swarrow & \\ & & \boxed{\mathbb{F}^D} & D \end{matrix} \quad (91)$$

$$= \begin{matrix} & & \boxed{\mathbb{F}^W} & W \\ & \boxed{\mathbb{K}_0} & \swarrow & \\ & & & D \end{matrix} \quad (92)$$

$$= \begin{matrix} & & \boxed{\mathbb{F}^W} & W \\ & \boxed{\mathbb{K}_0} & \swarrow & \\ & & & D \end{matrix} \quad (93)$$

$$= \begin{matrix} & & \boxed{\mathbb{F}^W} & W \\ & \boxed{\mathbb{K}} & \swarrow & \\ & & & D \end{matrix} \quad (94)$$

$$= \begin{matrix} & & \boxed{\mathbb{K}^{W|D}} & W \\ & & \swarrow & \\ & & & D \end{matrix} \quad (95)$$

□

We use the informal convention of labelling wires in quote marks “X” if that wire is “supposed to” carry the label X but the label may not be valid.

Theorem 0.1.21 (Iterated disintegration). *Given a kernel space (\mathbb{K}, D, Ω) , random variables X, Y and Z and domain variable D ,*

$$\begin{matrix} \text{“D”} \\ \text{“X”} \end{matrix} \text{---} \boxed{\mathbb{K}^{Y|XD}} \text{---} \boxed{\mathbb{K}^{Z|XYD}} \text{---} \begin{matrix} \text{“Z”} \\ \text{“Y”} \end{matrix} \in \mathbb{K}^{\{ZY|XD\}} \quad (96)$$

Equivalently, for $d \in D$ and $x \in X$, $A \in \mathcal{Y}$, $B \in \mathcal{Z}$,

$$(d, x; A, B) \mapsto \int_A \mathbb{K}_{(x,y,d)}^{Z|XYD}(B) d\mathbb{K}_{(x,d)}^{Y|XD}(y) \in \mathbb{K}^{\{ZY|XD\}} \quad (97)$$

Proof.

write this up

□

The existence of disintegrations of standard measurable probability spaces is well known.

Theorem 0.1.22 (Disintegration existence - probability space). *Given a probability measure $\mu \in \Delta(\mathcal{X} \otimes \mathcal{Y})$, if (F, \mathcal{F}) is standard then a disintegration $\mathbb{K} : X \rightarrow \Delta(\mathcal{Y})$ exists (Çinlar, 2011).*

In particular, if for all $x \in X$, $\mathbb{P}^X(X \in \{x\}) > 0$, then $\mathbb{P}_x^{Y|X}(Y \in A) = \frac{\mathbb{P}^{XY}(Y \in A \ \& \ X \in \{x\})}{\mathbb{P}^X(X \in \{x\})}$.

For Markov kernel spaces, we make the simplifying assumption that the domain space D is a discrete space. Given this assumption, there exists a positive definite probability $\mu \in \Delta(\mathcal{D})$. That is, for every $d \in D$, $\mu(\{d\}) > 0$. Given this assumption, for every Markov kernel space (\mathbb{K}, D, Ω) there is a probability space $(\mathbb{P}, \Omega \times D)$ such that \mathbb{K} can be uniquely defined as a disintegration of \mathbb{P} . For uncountable D , even if it is standard measurable, this is not possible (Hájek, 2003).

Definition 0.1.23 (Relative probability space).

better name

Given a Markov kernel space (\mathbb{K}, D, Ω) and a positive definite measure $\mu \in \Delta(\mathcal{D})$, $(\mu\mathbb{K}, \Omega \times D)$ is a *relative probability space*.

For any random variable $X : \Omega \times D \rightarrow X$ on (\mathbb{K}, D, Ω) , its relative on $(\mu\mathbb{K}, \Omega \times D)$ is given by the same measurable function, and we give it the same name X .

Lemma 0.1.24 (Agreement of disintegrations). *Given a Markov kernel space (\mathbb{K}, D, Ω) , any relative probability space $(\mu\mathbb{K}, \Omega \times D)$ and any random variables $X : \Omega \times D \rightarrow X$, $Y : \Omega \times D \rightarrow Y$, $\mathbb{K}^{\{Y|XD\}} = (\mu\mathbb{K})^{\{Y|XD\}}$ (note that this set equality).*

Proof. Define $\mathbb{P} := \mu\mathbb{K}$ and let \mathbb{M} be an arbitrary version of $\mathbb{K}^{\{Y|XD\}}$. Then

$$\begin{array}{c}
 \begin{array}{c} \triangleleft \\ \text{P}^{XYD} \end{array} \begin{array}{l} X \\ Y \\ D \end{array} = \begin{array}{c} \triangleleft \\ \mu \end{array} \begin{array}{c} \boxed{\mathbb{K}^{XY|D}} \\ X \\ Y \\ D \end{array} \quad (98)
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{c} \triangleleft \\ \mu \end{array} \begin{array}{c} \boxed{\mathbb{K}^{X|D}} \\ X \\ D \end{array} \begin{array}{c} \boxed{\mathbb{M}} \\ Y \\ D \end{array} = \quad (99)
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{c} \triangleleft \\ \text{P}^{XD} \end{array} \begin{array}{c} \boxed{\mathbb{M}} \\ X \\ Y \\ D \end{array} = \quad (100)
 \end{array}$$

Thus $\mathbb{M} \in \mathbb{P}^{\{Y|XD\}}$.

Let \mathbb{N} be an arbitrary version of $\mathbb{P}^{\{Y|XD\}}$. To show that $\mathbb{N} \in \mathbb{K}^{\{Y|XD\}}$, we will show for all $d \in D$

$$\mathbb{Q} := \quad (101)$$

$$= \mathbb{K}_d^{\text{XYD}|D} \quad (102)$$

For $A \in \mathcal{X}, B \in \mathcal{Y}, d \in D$, we have $\mathbb{Q}(A \times B \times \emptyset) = 0 = \mathbb{K}_d^{\text{XYD}|D}(A \times B \times \emptyset)$, and for $\{d\} \in \mathcal{D}$ we have $\mu(\{d\}) > 0$ so:

$$\mathbb{Q}(A \times B \times \{d\}) = \int_{X^2} \int_X \int_{D^3} \mathbb{N}_{d'',x'}(A) \mathbf{Id}_{x''}(B) \mathbf{Id}_{d'''}(\{d\}) d\gamma_d(d', d'', d''') d\mathbb{K}_d^{\text{X}|D}(x) d\gamma_x(x', x'') \quad (103)$$

$$= \delta_d(\{d\}) \int_X \mathbb{N}_{d,x}(A) \delta_x(B) d\mathbb{K}_d^{\text{X}|D}(x) \quad (104)$$

$$= \frac{1}{\mu(\{d\})} \int_{\{d\}} d\mu(d') \int_X \mathbb{N}_{d,x}(A) \delta_x(B) d\mathbb{K}_d^{\text{X}|D}(x) \quad (105)$$

$$= \frac{1}{\mu(\{d\})} \int_D \int_X \mathbb{N}_{d,x}(A) \delta_{d'}(\{d\}) \delta_x(B) d\mathbb{K}_d^{\text{X}|D}(a) d\mu(d') \quad (106)$$

$$= \frac{1}{\mu(\{d\})} \int_D \int_X \mathbb{N}_{d,x}(A) \delta_{d'}(\{d\}) \delta_x(B) d\mathbb{K}_{d'}^{\text{X}|D}(a) d\mu(d') \quad (107)$$

$$= \frac{1}{\mu(\{d\})} \mathbb{P}^{\text{XYD}}(A \times B \times \{d\}) \quad (108)$$

$$= \frac{1}{\mu(\{d\})} \int_D \mathbb{K}_{d'}^{\text{XYD}|D}(A \times B \times \{d\}) d\mu(d') \quad (109)$$

$$= \frac{1}{\mu(\{d\})} \int_D \mathbb{K}_{d'}^{\text{XY}|D}(A \times B) \delta_{d'}(\{d\}) d\mu(d') \quad (110)$$

$$= \mathbb{K}_d^{\text{XY}|D}(A \times B) \quad (111)$$

$$= \mathbb{K}_d^{\text{XY}|D}(A \times B) \delta_d(\{d\}) \quad (112)$$

$$= \int_D \mathbb{K}_{d'}^{\text{XY}}(A \times B) \delta_{d''}(\{d\}) d\gamma_d(d', d'') \quad (113)$$

$$= \mathbb{K}_d^{\text{XYD}|D}(A \times B \times \{d\}) \quad (114)$$

Equality follows from the monotone class theorem. Thus $\mathbb{N} \in \mathbb{K}^{\{\text{Y}|XD\}}$. \square

Thus any kernel conditional probability $\mathbb{K}^{\text{Y}|XD}$ can equally well be considered a regular conditional probability $\mathbb{P}^{\text{Y}|XD}$ for a related probability space $(\mathbb{P}, \Omega \times D)$ under the obvious identification of random variables, provided D is countable. Note that any conditional probability $\mathbb{P}^{\text{Y}|X}$ that is *not* conditioned on D is undefined in the kernel space (\mathbb{K}, D, Ω) .

$Y \rightarrow \Delta(\mathcal{X})$ be the map $x \mapsto \mathbb{M}_{(x, d_0)}$. By constancy in D , $\mathbb{M} = * \otimes N$. We wish to show $\mathbb{P}^{X|Y} \underline{\otimes} \mathbb{P}^{D|Y} \in \mathbb{P}^{X \cup D|Y}$. By Theorem 0.1.21, we have

(119)

Definition 0.1.28 (Conditional probability existence). Given a kernel space (\mathbb{K}, D, Ω) and random variables X, Y , we say $\mathbb{K}^{Y|X}$ exists if $Y \perp\!\!\!\perp_{\mathbb{K}} D|X$. If $\mathbb{K}^{Y|X}$ exists then it is by definition equal to $\mathbb{P}^{Y|X}$ for any related probability space $(\mathbb{P}, \Omega \times D)$.

Note that $\mathbb{K}^{Y|XD}$ always exists.

Definition 0.1.29 (Conditional Independence). Given a kernel space (\mathbb{K}, D, Ω) , relative probability space $(\mathbb{P}, \Omega \times D)$, variables X, Y and Z , X is *conditionally independent* of Z given Y , written $X \perp\!\!\!\perp_{\mathbb{K}} Z|Y$ if $\mathbb{K}^{X|YZ}$ exists and any of the following equivalent conditions hold:

Almost sure equality

- $\mathbb{P}^{XZ|Y} \sim \mathbb{P}^{X|Y} \underline{\otimes} \mathbb{P}^{Z|Y}$
- For any version of $\mathbb{P}^{X|Y}$, $\mathbb{P}^{X|Y} \otimes *_Z$ is a version of $\mathbb{K}^{X|YZ}$
- There exists a version of $\mathbb{K}^{X|YZ}$ constant in Z

Lemma 0.1.30 (Diagrammatic consequences of labels). *In general, diagram labels are “well behaved” with regard to the application of any of the special Markov kernels: identities 17, swaps 28, discards 34 and copies 25 as well as with respect to the coherence theorem of the CD category. They are not “well behaved” with respect to composition.*

Fix some Markov kernel space (\mathbb{K}, D, Ω) and random variables X, Y, Z taking values in X, Y, Z respectively. $\text{Sat} :$ indicates that a labeled diagram satisfies definitions 0.1.16 and 0.1.19 with respect to (\mathbb{K}, D, Ω) and X, Y, Z . The following always holds:

$$\text{Sat} : X \text{ --- } X \tag{120}$$

and the following implications hold:

$$\text{Sat} : Z - \boxed{\mathbb{K}} - \begin{array}{c} X \\ \diagdown \\ Y \end{array} \implies \text{Sat} : Z - \boxed{\mathbb{K}} - \begin{array}{c} X \\ * \end{array} \quad (121)$$

$$\text{Sat} : Z - \boxed{\mathbb{K}} - \begin{array}{c} X \\ \diagdown \\ Y \end{array} \implies \text{Sat} : Z - \boxed{\mathbb{K}} - \begin{array}{c} Y \\ \diagup \\ X \end{array} \quad (122)$$

$$\text{Sat} : Z - \boxed{\mathbb{L}} - X \implies \text{Sat} : Z - \boxed{\mathbb{L}} - \begin{array}{c} X \\ \diagup \\ X \end{array} \quad (123)$$

$$\text{Sat} : Z - \boxed{\mathbb{K}} - Y \implies \text{Sat} : \begin{array}{c} Z \\ \diagdown \\ \boxed{\mathbb{K}} - Y \end{array} \quad (124)$$

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

- $\mathbb{K} \text{Id} \otimes * (w; A) = \int_{X \times Y} \delta_x(A) \mathbb{1}_Y(y) d\mathbb{K}_w(x, y) = \mathbb{K}_w(A \times Y) = \mathbb{P}_{X|Z}(w; A)$
- $\int_{X \times Y} \delta_{\text{swap}(x, y)}(A \times B) d\mathbb{K}_w(x, y) = \mathbb{P}_{YX|Z}(w; A \times B)$
- $\mathbb{K} \Upsilon(w; A \times B) = \int_X \delta_{x, x}(A \times B) d\mathbb{K}_w(x) = \mathbb{P}_{XX|Z}(w; A \times B)$

124: Suppose \mathbb{K} is a version of $\mathbb{P}_{Y|Z}$. Then

$$\mathbb{P}_{ZY} = \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \diagdown \\ \boxed{\mathbb{K}} - \begin{array}{c} Z \\ \diagdown \\ Y \end{array} \end{array} \quad (125)$$

$$\mathbb{P}_{ZZY} = \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \diagdown \\ \boxed{\mathbb{K}} - \begin{array}{c} Z \\ \diagup \\ Z \\ \diagdown \\ Y \end{array} \end{array} \quad (126)$$

$$= \begin{array}{c} \triangleleft \mathbb{P}_Z \\ \diagdown \\ \boxed{\mathbb{K}} - \begin{array}{c} Z \\ \diagup \\ Z \\ \diagdown \\ Y \end{array} \end{array} \quad (127)$$

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

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

$$\text{Sat} : Z - \boxed{\mathbb{K}} - Y, Y - \boxed{\mathbb{L}} - X \implies \text{Sat} : Z - \boxed{\mathbb{K}} - \boxed{\mathbb{L}} - X \quad (128)$$

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

Chapter 1

Chapter 3: See-do models

These are “todo” notes. All such notes that involve theoretical development are also collected in an unordered list of outstanding theoretical questions

The basic claim of this chapter is that see-do models are the basic type of thing that everyone who is studying “causal inference” is working with, even if they don’t know it themselves

Consider the following problem: you are presented with a collection \mathcal{H} of hypotheses about how the world might function and a vector \mathbf{x} of observational data which you know could have taken values in some space X . You want to determine which hypothesis $H \in \mathcal{H}$ best describes the world. However you ultimately solve the problem, the next step you take will probably be to determine for each $H \in \mathcal{H}$ a probability distribution $\mathbb{P}_H \in \Delta(\mathcal{X})$ that indicates how likely you would be to observe the various elements of X were H in fact the case. This is a *statistical model* – an indexed set of probability distributions $\{\mathbb{P}_H | H \in \mathcal{H}\}$. Statistical models are ubiquitous in the field of statistics – they are found in statistical decision theory where the elements of \mathcal{H} are typically called “states” (Wald, 1950), in Bayesian inference where the elements of \mathcal{H} may be called “parameters” (Freedman, 1963) and in frequentist inference where elements of \mathcal{H} they may be called “hypotheses” (Fisher, 1992).

These different approaches to statistics may have different notions of what the “best hypothesis” H is, may employ different estimation methods and may not even agree about what “distributed according to \mathbb{P}_H ” means. Nonetheless, the interpretation of the statistical model in each case is roughly the same: supposing $H \in \mathcal{H}$ is true, the data will be distributed according to \mathbb{P}_H . A statistical model takes a hypothesis and tells you what you are likely to *see*.

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; Savage, 1972) 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).

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 (1972) 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 con-

sequence model and their coupling, making them “two picture” causal models. Causal Bayesian Networks and Single World Interention 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).

1.1 Definition

Terminology question: The variables H and D aren’t necessarily *random* in the commonsense understanding of the word. They’re also defined on a *kernel space* rather than a *probability space*. They’re currently called random variables simply by virtue of being measurable functions on the outcome space. I’m not a huge fan of “quasirandom variables”, but it does capture the idea that these things are very similar to random variables but not exactly the same.

Definition 1.1.1 (See-Do model). A *see-do model* $\langle \mathbb{T}, H, D, X, Y \rangle$ is a kernel space (Definition 0.1.7) $(\mathbb{T}, H \times D, X \times Y)$ along with four random variables: the *hypothesis* $H : H \times D \times X \times Y \rightarrow H$, the *choice* $D : H \times D \times X \times Y \rightarrow D$, the *observations* $X : H \times D \times X \times Y \rightarrow X$ and the *consequences* $Y : H \times D \times X \times Y \rightarrow Y$, all given by the obvious projection maps.

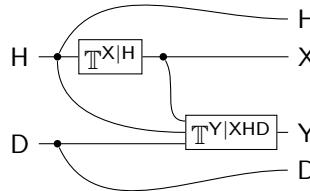
The spaces H , D , X and Y are the hypothesis, choice, observation and consequence spaces respectively.

A see-do model has the additional property that, holding the hypothesis fixed, the observations are independent of the choices - i.e. $X \perp\!\!\!\perp_{\mathbb{T}} D|H$. We require that $H \times D$ is countable.

Theorem 1.1.2 (Observation and Consequence models). *Any see-do model (\mathbb{T}, H, D, X, Y) can be uniquely represented by the following pair of Markov kernels:*

- The observation map $\mathbb{T}^{X|H}$
- The consequence model $\mathbb{T}^{Y|XHD}$

Furthermore



$$\mathbb{T} = \tag{1.1}$$

Maybe moves proofs out of main text

Proof. By 0.1.1,

$$\mathbb{T} = \quad (1.2)$$

By the assumption $X \perp\!\!\!\perp_{\mathbb{T}} D|H$ and version 2 of conditional independence from Theorem 0.1.27,

$$\mathbb{T} = \quad (1.3)$$

$$= \quad (1.4)$$

□

Definition 1.1.3 (Consequence map). Given a see-do model (\mathbb{T}, H, D, X, Y) , a *consequence map* is a map $\mathbb{C} : D \rightarrow \Delta(\mathcal{Y})$ where D is a choice set and Y is a consequence set.

The consequence model evaluated at any particular hypothesis $h \in H$, $\mathbb{T}_{\cdot, h, \cdot}^{Y|XHD}$ is a consequence map.

Not quite sure if this is the right place for the following definition

Definition 1.1.4 (Bayesian See-Do Model). A Bayesian See-Do Model $\langle \mathbb{U}, D, X, Y \rangle$ is a Markov kernel space $(\mathbb{U}, D, X \times Y)$ with the property $X \perp\!\!\!\perp_{\mathbb{U}} D$, along with choices D , observations X and consequences Y , defined as before.

$$\mathbb{U} := (\gamma \otimes \text{Id}^D) \mathbb{T} \quad (1.5)$$

Maybe moves proofs out of main text

By definition

$$= (\gamma \otimes \text{Id}^D) \mathbb{T}\mathbb{F}^{\mathbf{x}} \quad (1.7)$$

Which implies $X \perp\!\!\!\perp_{\mathbb{U}} D$ by version (2) of conditional independence (Theorem 0.1.27). \square

Suppose we are betting on the outcome of the flip of a possibly biased coin with payout 1 for a correct guess and 0 for an incorrect guess, and we are given N previous flips of the coin to inspect. This situation can be modeled by a hypothesis sufficient see-do model. Define $\mathbb{B} : (0, 1) \rightarrow \Delta(\{0, 1\})$ by $\mathbb{B} : \mathbf{H} \mapsto \text{Bernoulli}(\mathbf{H})$. Then define \mathbb{T} by:

- Choice set: $D = \{0, 1\}$
- Observation set: $X = \{0, 1\}^N$
- Consequence set: $Y = \{0, 1\}$

- Hypothesis set: $H = (0, 1)$
- Observation map: $\mathbb{T}^{\mathbf{X}|\mathbf{H}} : \mathcal{Y}^N \rightarrow \mathbb{B}$
- Consequence model: $\mathbb{T}^{\mathbf{Y}|\mathbf{D}\mathbf{H}} : (h, d) \mapsto \text{Bernoulli}(1 - |d - h|)$

In this model, the chance \mathbf{H} of the coin landing on heads is as much as we can hope to know about the success of our bet. \mathbf{H} may be inferred from observation by some standard method, and

1.1.1 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 0.1.4 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 $\mathbb{C} := \mathbb{T}^{\mathbf{Y}|\mathbf{H}\mathbf{D}}$.

Definition 1.1.6 (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 $\mathbf{Y}_2 \perp\!\!\!\perp_{\mathbb{C}} \mathbf{D} | \mathbf{Y}_1 \mathbf{H}$.

1.1.2 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 1.1.7 (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 1.1.8 (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 $(s, y_1) \mapsto y(s, [y_1(s, \cdot)]^{-1}(y_1(s, d)))$ 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 $\mathbb{C}^{\mathbf{Y}_1|\mathbf{Y}_2, \mathbf{S}, \mathbf{D}}$ because, for all $A \in \mathcal{Y}_2$, $B \in \mathcal{Y}_1$, $s \in S$, $d \in D$:

$$\mathbb{C}_{(s,d)}^{\mathbf{Y}_2|\mathbf{SD}} \mathbf{Y}(\mathbb{M} \otimes \text{Id}) = \int_A \mathbb{M}(y'_2, d, s; B) d\delta_{y_2(s,d)}(y'_2) \quad (1.10)$$

$$= \int_A \delta_{f(y'_2, s)}(B) d\delta_{y_2(s,d)}(y'_2) \quad (1.11)$$

$$= \delta_{f(y_2(s,d), s)}(B) \delta_{y_2(s,d)}(A) \quad (1.12)$$

$$= \delta_{y_1(s,d)}(B) \delta_{y_2(s,d)}(A) \quad (1.13)$$

$$= \delta_{y_2(s,d)} \otimes \delta_{y_1(s,d)}(A \times B) \quad (1.14)$$

\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

$$\mathbb{C}_{d, [i, m]}^{\text{LF}|\mathbf{DS}} = \delta_{(d \text{ XOR } i) \text{ AND } m} \otimes \delta_{(d \text{ XOR } i) \text{ AND } m} \quad (1.15)$$

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:

$$\mathbb{C}_\gamma = \frac{1}{4}\mathbb{C}_{\cdot, [0, 0]} + \frac{1}{4}\mathbb{C}_{\cdot, [0, 1]} + \frac{1}{2}\mathbb{C}_{\cdot, [1, 1]} \quad (1.16)$$

$$\mathbb{C}_\gamma^{\text{FL}|D} = \frac{1}{4} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \otimes \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (1.17)$$

Then

$$[1, 0]\mathbb{C}_\gamma^{\text{FL}|D} = \frac{1}{4}[0, 1] \otimes [1, 0] + \frac{1}{4}[0, 1] \otimes [0, 1] + \frac{1}{2}[1, 0] \otimes [1, 0] \quad (1.18)$$

$$[1, 0]\mathbf{Y}(\mathbb{C}_\gamma^{\text{F}|D} \otimes \mathbb{C}_\gamma^{\text{L}|D}) = (\frac{1}{2}[0, 1] + \frac{1}{2}[1, 0]) \otimes (\frac{1}{4}[0, 1] + \frac{3}{4}[1, 0]) \quad (1.19)$$

$$\implies [1, 0]\mathbb{C}_\gamma^{\text{FL}|D} \neq [1, 0]\mathbf{Y}(\mathbb{C}_\gamma^{\text{F}|D} \otimes \mathbb{C}_\gamma^{\text{L}|D}) \quad (1.20)$$

Thus under hypothesis mixture 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 X D -causes Y relative to γ iff $Y \perp\!\!\!\perp_{\mathbb{C}_\gamma} D|X$.

Theorem 1.1.9 shows that under some conditions, D -causation can hold for arbitrary mixtures over subsets of the hypothesis class \mathbf{H} .

Theorem 1.1.9 (Universal D -causation). *If $X \perp\!\!\!\perp H|D$ for all $H, H' \in S \subset \mathbf{H}$ and X D -causes Y in all $H \in S$, then X D -causes 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

$$\mathbb{C}_\gamma := \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}} \text{---} F \end{array} \quad (1.21)$$

Because $\mathbb{C}_H^{X|D} = \mathbb{C}_{H'}^{X|D}$ for all $H, H' \in \mathbf{H}$, we have

$$\begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{X|DH}} \text{---} X \end{array} \quad H = \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{X|DH}} \text{---} X \end{array} \quad H \quad (1.22)$$

Also

$$\mathbb{C}_\gamma^{XY|D} = \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}} \text{---} \boxed{\mathbb{F}^{X \otimes Y}} \text{---} \begin{array}{l} X \\ Y \end{array} \end{array} \quad (1.23)$$

$$= \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{XY|DH}} \text{---} \begin{array}{l} X \\ Y \end{array} \end{array} \quad (1.24)$$

$$= \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{X|DH}} \text{---} \boxed{\mathbb{C}^{Y|XDH}} \text{---} \begin{array}{l} Y \\ X \end{array} \end{array} \quad (1.25)$$

$$\begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{X|DH}} \text{---} \boxed{\mathbb{C}^{Y|XH}} \text{---} \begin{array}{l} Y \\ X \end{array} \end{array} \quad Y \perp\!\!\!\perp_{\mathbb{C}_\gamma} D|XH \quad (1.26)$$

$$\stackrel{1.22}{=} \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}^{X|DH}} \text{---} \boxed{\mathbb{C}^{Y|XH}} \text{---} \begin{array}{l} Y \\ X \end{array} \end{array} \quad (1.27)$$

$$\stackrel{1.22}{=} \begin{array}{c} \triangle \gamma \\ | \\ D \text{---} \boxed{\mathbb{C}_\gamma^{X|DH}} \text{---} \boxed{\mathbb{C}^{Y|XH}} \text{---} \begin{array}{l} Y \\ X \end{array} \end{array} \quad (1.28)$$

Equation 1.28 establishes that $(\gamma \otimes \mathbf{Id}_X \otimes \dagger_D) \mathbb{C}^{Y|XH}$ 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:

$$H \perp\!\!\!\perp D \wedge H \perp\!\!\!\perp X|D \implies H \perp\!\!\!\perp XD \quad (1.29)$$

$$\implies H \perp\!\!\!\perp D|X \quad (1.30)$$

$$D \perp\!\!\!\perp H|X \wedge D \perp\!\!\!\perp Y|XH \implies D \perp\!\!\!\perp Y|X \quad (1.31)$$

$$\implies Y \perp\!\!\!\perp D|X \quad (1.32)$$

□

1.1.3 Properties of D-causation

If X D-causes Y relative to \mathbb{C}_H , then the following holds:

$$\mathbb{C}_H^{X|D} = D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (1.33)$$

This follows from version (2) of Definition 0.1.29:

$$\mathbb{C}_H^{X|D} = D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|XD}} - Y \quad (1.34)$$

$$= D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (1.35)$$

$$= D - \boxed{\mathbb{C}^{X|D}} - \boxed{\mathbb{C}^{Y|X}} - Y \quad (1.36)$$

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?

1.1.4 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 not requirement that D_1 comes “before” D_2 or even that they have any “before” and “after” relations at all.

Define parallel decisions, show that they induce potential outcomes

1.1.5 Residual dependence on observations

Definition 1.1.10 (Hypothesis sufficiency). The hypothesis H is *sufficient* for a see-do model if the consequence model has no dependence on observations X conditional on H . That is, $Y \perp\!\!\!\perp_{\mathbb{T}} X | DH$.

A hypothesis sufficient see-do model can be specified with:

- Hypothesis space H , choices D , observations X and consequences Y
- Observation map $\mathbb{T}^{X|H}$
- Reduced consequence model $\mathbb{T}^{Y|HD}$

Given observations X , assumed to be an IID sequence X_1, X_2, \dots conditional on H , a common “causal inference problem” is to estimate the “true” distribution of observations $\mathbb{T}_{h^*}^{X_i|H}$ and from this to estimate the consequence model $\mathbb{T}_{h^*}^{Y|HD}$, if this is possible. This problem only makes sense if hypothesis sufficiency is assumed – once h^* is given, the consequence model of interest has no further dependence on X . We show that all decision problems can be modeled by a hypothesis sufficient see-do model.

Examples of hypothesis sufficient and insufficient see-do models

Recall the previous example: suppose we are betting on the outcome of the flip of a possibly biased coin with payout 1 for a correct guess and 0 for an incorrect guess, and we are given N previous flips of the coin to inspect. This situation can be modeled by a hypothesis sufficient see-do model. Define $\mathbb{B} : (0, 1) \rightarrow \Delta(\{0, 1\})$ by $\mathbb{B} : H \mapsto \text{Bernoulli}(H)$. Then define ${}^1\mathbb{T}$ by:

- $D = \{0, 1\}$
- $X = \{0, 1\}^N$
- $Y = \{0, 1\}$
- $H = (0, 1)$
- ${}^1\mathbb{T}^{X|H} : \varphi^N \mathbb{B}$
- ${}^1\mathbb{T}^{Y|DH} : (h, d) \mapsto \text{Bernoulli}(1 - |d - h|)$

In this model, the chance H of the coin landing on heads is as much as we can hope to know about how our bet will work out.

Suppose instead that in addition to the N prior flips, we manage to look at the outcome of the flip on which we will bet. In this case, the situation can be modeled by the following hypothesis insufficient see-do model ${}^2\mathbb{T}$:

- $D = \{0, 1\}$
- $X = \{0, 1\}^{N+1}$

- $Y = \{0, 1\}$
- $H = (0, 1)$
- ${}^2\mathbb{T}^{X|H} : \varphi^{N+1}\mathbb{B}$
- ${}^2\mathbb{T}^{Y|XHD} : (h, \mathbf{x}, d) \mapsto \delta_{1-|d-x_{N+1}|}$

In this case, even if we are told the value of H , we still benefit from using the observed data when making our decision.

It appears that it might be possible to model the second situation with a hypothesis sufficient model by including the result of the $N + 1$ th flip in the hypothesis. Define the new hypothesis space $H' = (0, 1) \times \{0, 1\}$ and define ${}^3\mathbb{T}$ by:

- $D = \{0, 1\}$
- $X = \{0, 1\}^{N+1}$
- $Y = \{0, 1\}$
- $H' = (0, 1) \times \{0, 1\}$
- ${}^3\mathbb{T}^{X|H'} : (\varphi^N\mathbb{B} \otimes \delta_{x_{N+1}})$
- ${}^3\mathbb{T}^{Y|H'D} : (h, x_{N+1}, d) \mapsto \delta_{1-|d-x_{N+1}|}$

However, X_{N+1} is related to the previous flips $X_{<N}$ and ${}^3\mathbb{T}$ ignores this fact. In particular, given any $H' = (h, _)$, X_{N+1} as well as X_i , $i \leq N$ should all distributed according to $\text{Bernoulli}(h)$. Thus ${}^2\mathbb{T}$ is preferable to ${}^3\mathbb{T}$ because it represents more of the knowledge we have about the problem.

If a see-do model is employed in a *decision problem* – defined in the next section – there is an alternative way to avoid hypothesis insufficiency that does not require throwing out some of the model structure.

The importance of this is that counterfactual questions are usually *not* decision problems and so they do not have the possibility of avoiding insufficiency available; also *in practice* counterfactual problems are usually hypothesis insufficient while decision problems are usually not.

1.1.6 Causal questions and decision functions

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

1. Association: How are W and Z related? How would observing W change my beliefs about Z ?
2. Intervention: What would happen if I do ... ? How can I make ... happen?
3. Counterfactual: What if I had done ... instead of what I actually did?

Causal decision problems are, roughly speaking, “interventional” problems. In English, a causal decision problem roughly asks

Given that I have data X and I know which values of Y I would like to see and some knowledge about how the world works, which of my available choices D should I select?

This type of question presupposes somewhat more than Pearl’s prototypical interventional questions. First, it supposes that we have *preferences* over the values that Y might take, which we need not have to answer the question “What would happen if I do ...?”. Secondly, and crucially to our theory, causal decision problem suppose that we are given data and a set of choices.

We will return to the question of preferences. For now, we will focus on the idea that a causal decision problem is about selecting a choice given data. That is, however the selection is made, the answer to a causal decision problem is always a *decision function* $\mathbb{D} : X \rightarrow \Delta(\mathcal{D})$.

Avoiding insufficiency with decision functions

Show that a decision problem with a hypothesis insufficient model induces an equivalent decision problem with a hypothesis sufficient model with an expanded set of choices, subject to some conditions.

Decision rules

See-do models encode the relationship between observed data and consequences of decisions. In order to actually make decisions, we also require preferences over consequences. We suppose that a *utility function* is given, and evaluate the desirability of consequences using *expected utility*. A see-do model along with a utility allows us to evaluate the desirability of *decisions rules* according to each hypothesis.

Definition 1.1.11 (Utility function). Given a See-Do Model $\mathbb{T} : \mathcal{H} \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$, a *utility function* u is a measurable function $Y \rightarrow \mathbb{R}$.

Definition 1.1.12 (Expected utility). Given a utility function $u : Y \rightarrow \mathbb{R}$ and probability measures $\mu, \nu \in \Delta(\mathcal{Y})$, the *expected utility* of μ is $\mathbb{E}_\mu[u]$.

μ is *preferred* to ν if $\mathbb{E}_\mu[u] \geq \mathbb{E}_\nu[u]$, and *strictly preferred* if $\mathbb{E}_\mu[u] > \mathbb{E}_\nu[u]$.

Definition 1.1.13 (Decision rule). Given a see-to map $\mathbb{T} : \mathcal{H} \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$, a *decision rule* is a Markov kernel $X \rightarrow \Delta(\mathcal{D})$. A *deterministic decision rule* is a decision rule that is deterministic.

Define deterministic Markov kernels

Expected utility together with a decision rule gives rise to the definition of *risk*, which connects CSDT to classical statistical decision theory (SDT). For historical reasons, risks are minimised while utilities are maximised.

Definition 1.1.14 (Risk). Given a see-to map $\mathbb{T} : \mathbf{H} \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$, a utility $u : Y \rightarrow \mathbb{R}$ and the set of decision rules \mathcal{U} , the *risk* is a function $l : \mathbf{H} \times \mathcal{U} \rightarrow \mathbb{R}$ given by

$$R(\mathbf{H}, \mathbb{U}) := - \int_X \mathbb{U}_x \mathbb{T}_{\cdot, x, \mathbf{H}}^{Y|DXH} u d\mathbb{T}_{\mathbf{H}}^{X|H}(x) \quad (1.37)$$

for $\mathbf{H} \in \mathbf{H}$, $\mathbb{U} \in \mathcal{U}$. Here $\mathbb{U}_x \mathbb{T}_{\cdot, x, \mathbf{H}}^{Y|DXH} u$ is the product of the measure \mathbb{U}_x , the kernel $\mathbb{T}_{\cdot, x, \mathbf{H}}^{Y|DXH} : D \rightarrow \Delta(\mathcal{Y})$ and the function u .

The loss induces a partial order on decision rules. If for all \mathbf{H} , $l(\mathbf{H}, \mathbb{U}) \leq l(\mathbf{H}, \mathbb{U}')$ then \mathbb{U} is at least as good as \mathbb{U}' . If, furthermore, there is some \mathbf{H}_0 such that $l(\mathbf{H}_0, \mathbb{U}) < l(\mathbf{H}_0, \mathbb{U}')$ then \mathbb{U} is preferred to \mathbb{U}' .

Definition 1.1.15 (Induced statistical decision problem). A see-do model $\mathbb{T} : \mathbf{H} \times D \rightarrow \Delta(\mathcal{X} \otimes \mathcal{Y})$ along with a utility u induces the *statistical decision problem* $(\mathbf{H}, \mathcal{U}, R)$ with states \mathbf{H} , decisions \mathcal{U} and risks R .

Statistical decision problems usually define the risk via the loss, but it is only possible to define a loss with a hypothesis sufficient model. We don't actually need a loss, though: the complete class theorem still holds via the induced risk and Bayes risk

1.2 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 *parallelized*.

Before defining parallel choices, we will consider a “counterfactual model” without parallel choices. Consider the following definitions, first from Pearl (2009) pg. 203-204. Note that 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 0.1.10, but I'd rather not change his definitions if I can avoid it

put the following inside a quote environment somehow

““”

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, $Y_x(u) = Y_{M_x}(u)$.
””

Implicitly, Definition 7.1.3 proposes a set of “actions” here 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 , but this would usually rule out parallel actions which, as we will show, are the core of counterfactual models. 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).

Definition 7.1.4 then appears to define a function $D \rightarrow Y$

is just a Markov kernel $\mathbb{C} : D \times Y^{|D|} \rightarrow \Delta(Y)$ (where the choice set D and the consequence set Y is defined by some underlying see-do model $(\mathbb{T}, \mathbb{H}, \mathbb{D}, \mathbb{X}, \mathbb{Y})$).

Chapter 2

Chapter 4: See-do models compared to causal graphical models and potential outcomes

References

- Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- Kenta Cho and Bart Jacobs. Disintegration and Bayesian inversion via string diagrams. *Mathematical Structures in Computer Science*, 29(7):938–971, August 2019. ISSN 0960-1295, 1469-8072. doi: 10.1017/S0960129518000488.
- Florence Clerc, Fredrik Dahlqvist, Vincent Danos, and Ilias Garnier. Pointless learning. *20th International Conference on Foundations of Software Science and Computation Structures (FoSSaCS 2017)*, March 2017. doi: 10.1007/978-3-662-54458-7_21. URL [https://www.research.ed.ac.uk/portal/en/publications/pointless-learning\(694fb610-69c5-469c-9793-825df4f8ddec\).html](https://www.research.ed.ac.uk/portal/en/publications/pointless-learning(694fb610-69c5-469c-9793-825df4f8ddec).html).
- A. P. Dawid. Influence Diagrams for Causal Modelling and Inference. *International Statistical Review*, 70(2):161–189, 2002. ISSN 1751-5823. doi: 10.1111/j.1751-5823.2002.tb00354.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1751-5823.2002.tb00354.x>.
- A. Philip Dawid. Decision-theoretic foundations for statistical causality. *arXiv:2004.12493 [math, stat]*, April 2020. URL <http://arxiv.org/abs/2004.12493>. arXiv: 2004.12493.
- Angus Deaton and Nancy Cartwright. Understanding and misunderstanding randomized controlled trials. *Social Science & Medicine*, 210:2–21, August

2018. ISSN 0277-9536. doi: 10.1016/j.socscimed.2017.12.005. URL <http://www.sciencedirect.com/science/article/pii/S0277953617307359>.
- 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, 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.
- D. Heckerman and R. Shachter. Decision-Theoretic Foundations for Causal Reasoning. *Journal of Artificial Intelligence Research*, 3:405–430, December 1995. ISSN 1076-9757. doi: 10.1613/jair.202. URL <https://www.jair.org/index.php/jair/article/view/10151>.
- James J. Heckman. Randomization and Social Policy Evaluation. SSRN Scholarly Paper ID 995151, Social Science Research Network, Rochester, NY, July 1991. URL <https://papers.ssrn.com/abstract=995151>.
- Alan Hájek. What Conditional Probability Could Not Be. *Synthese*, 137(3): 273–323, December 2003. ISSN 1573-0964. doi: 10.1023/B:SYNT.0000004904.91112.16. URL <https://doi.org/10.1023/B:SYNT.0000004904.91112.16>.
- Alan Hájek. Interpretations of Probability. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, fall 2019 edition, 2019. URL <https://plato.stanford.edu/archives/fall2019/entries/probability-interpret/>.
- Chayakrit Krittanawong, Bharat Narasimhan, Zhen Wang, Joshua Hahn, Hafeez Ul Hassan Virk, Ann M. Farrell, HongJu Zhang, and WH Wilson Tang. Association between chocolate consumption and risk of coronary artery disease: a systematic review and meta-analysis. *European Journal of Preventive Cardiology*, July 2020. doi: 10.1177/2047487320936787. URL <http://journals.sagepub.com/doi/10.1177/2047487320936787>. Publisher: SAGE PublicationsSage UK: London, England.

- Finnian Lattimore and David Rohde. Replacing the do-calculus with Bayes rule. *arXiv:1906.07125 [cs, stat]*, December 2019. URL <http://arxiv.org/abs/1906.07125>. arXiv: 1906.07125.
- David K Lewis. Causation. *Journal of Philosophy*, 1986.
- Stephen L. Morgan and Christopher Winship. *Counterfactuals and Causal Inference: Methods and Principles for Social Research*. Cambridge University Press, New York, NY, 2 edition edition, November 2014. ISBN 978-1-107-69416-3.
- 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. Challenging the hegemony of randomized controlled trials: A commentary on Deaton and Cartwright. *Social Science & Medicine*, 2018.
- 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 cigarette-cancer link: evidentiary traditions, corporate denial, global toll. *Tobacco Control*, 21(2):87–91, March 2012. ISSN 0964-4563, 1468-3318. doi: 10.1136/tobaccocontrol-2011-050338. URL <https://tobaccocontrol.bmj.com/content/21/2/87>. Publisher: BMJ Publishing Group Ltd Section: The shameful past.
- Thomas S Richardson and James M Robins. Single world intervention graphs (swigs): A unification of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social Sciences, University of Washington Series. Working Paper*, 128(30):2013, 2013.
- Donald B. Rubin. Causal Inference Using Potential Outcomes. *Journal of the American Statistical Association*, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880. URL <https://doi.org/10.1198/016214504000001880>.
- Leonard J. Savage. *Foundations of Statistics*. Dover Publications, New York, revised edition edition, June 1972. ISBN 978-0-486-62349-8.

- Peter Selinger. A survey of graphical languages for monoidal categories. *arXiv:0908.3347 [math]*, 813:289–355, 2010. doi: 10.1007/978-3-642-12821-9_4. URL <http://arxiv.org/abs/0908.3347>. arXiv: 0908.3347.
- Ilya Shpitser and Judea Pearl. Complete Identification Methods for the Causal Hierarchy. *Journal of Machine Learning Research*, 9(Sep):1941–1979, 2008. ISSN ISSN 1533-7928. URL <https://www.jmlr.org/papers/v9/shpitser08a.html>.
- Statista. Cigarettes - worldwide | Statista Market Forecast, 2020. URL <https://www.statista.com/outlook/50010000/100/cigarettes/worldwide>.
- Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England, 1950.
- 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: