

Causal questions are questions that are answered
by a function

David Johnston

September 9, 2021

Researchers in the field of causal inference will often choose a causal framework as one of the first steps of their investigations, or in some cases, one of the first steps of their careers. One could postulate that “causal inference” is what one does when one does work using a causal modelling framework. We argue that “causal inference” is better understood by the kind of questions people ask on rather than the kind of framework people use to answer them. Pearl and Mackenzie (2018) has proposed a three-level hierarchy for classifying causal questions: at the bottom are “seeing” questions followed by “doing” questions with “imagining” questions at the top. We propose an alternative characterisation: “ordinary statistical questions” are questions involving data that are answered by distributions on a given set while “causal statistical questions” are questions involving data that are answered by stochastic functions with given domain and codomain. Potential outcomes and graphical models are features of modelling frameworks, while interventions and counterfactuals are features of causal problems. We show how both potential outcomes and causal graphical models arise in *see-do models*, a generic modelling framework we introduce that addresses causal questions in general, as we define them. We hypothesise that some confusion about interventions and counterfactuals arises from assuming they are given by the modelling framework rather than by the problem under investigation.

or something like that; it could also be functions of a distribution like a maximum likelihood estimate or a p-value

Contents

0.1	Technical prerequisites	3
0.1.1	Cartesian and tensor products	4
0.1.2	Indicator functions, delta measures and function-associated Markov kernels	4
0.1.3	Copy maps and sequences	5
0.1.4	Generalised random variables	5
0.1.5	Disintegration	5
0.1.6	Conditional independence	6
0.2	See-do models	6
0.2.1	See-do models for data-driven decision problems	7

0.1 Technical prerequisites

Many people are familiar with probability theory, but some may be less familiar with *Markov kernels*, which play a central role in the work developed in this paper. Markov kernels are measurable functions that map to probability distributions on some measurable set. Expressions like $\mathbb{P}(Y|X)$ and $x \mapsto \mathbb{P}(Y|do(X = x))$ represent Markov kernels that map from the range of the random variable X to probability distributions on the range of the random variable Y . Conditional probabilities like $\mathbb{P}(Y|X)$ are typically obtained by disintegrating a joint probability $\mathbb{P}(Y, X)$, but Markov kernels can also be things other than conditional probabilities, like “interventional maps” $x \mapsto \mathbb{P}(Y|do(X = x))$.

We will consider only discrete sets in this paper, as uncountable sets raise a number of difficulties we prefer to avoid in this paper. A discrete set is a set X which is at most countably infinite, equipped with the

Footnote?: These difficulties may be a general phenomenon - for example, letting $X := \text{Range}(X)$, if X is real-valued, then for almost every $x \in X$ there are many choices for $\mathbb{P}(Y|do(X = x))$ that all satisfy the definition given by Pearl (2009) because $\mathbb{P}(X = x)$ can be positive for an at most countable subset of X . Also, there are examples of theorems that hold for discrete sets only Heymann et al. (2021)

We will take advantage of the fact that we are working with discrete sets and define probability measures as vectors, measurable functions as covectors

and Markov kernels as matrices.

Given a set X , a probability measure \mathbb{P} on X is a covector in $\mathbb{R}^{|X|}$; $\mathbb{P} := (\mathbb{P}^i)_{i \in X}$. We require that

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

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

An *event* A is a subset of X , and we define $\mathbb{P}(A) := \sum_{i \in A} \mathbb{P}^i$.

A measurable function $f : X \rightarrow Y$ is a vector in $Y^{|X|}$; $f := (f^i)_{i \in X}$ where Y is a vector space.

Given discrete sets X and Y , a Markov kernel $\mathbb{K} : X \rightarrow \Delta(Y)$ is a matrix in $\mathbb{R}^{|X| \times |Y|}$; $\mathbb{K} = (K_i^j)_{i \in X, j \in Y}$ where

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

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

We use subscripts to refer to rows of a Markov kernel $\mathbb{K}_x := (K_x^j)_{j \in Y}$; these are all probability measures.

$\mathbb{P}\mathbb{K}$ refers to a matrix-matrix product in the usual way, and similarly $\mathbb{P}f$ is a covector-vector product and $\mathbb{K}f$ is a matrix-vector product.

Products have the following properties: $\mathbb{P}\mathbb{K}$ is a probability measure, $\mathbb{K}f$ is a measurable function and $\mathbb{P}f := \mathbb{E}_{\mathbb{P}}[f]$ is a scalar which we define as the expectation of f under \mathbb{P} . Given another Markov kernel $\mathbb{L} : Y \rightarrow \Delta(Z)$, the matrix product $\mathbb{K}\mathbb{L}$ is also a Markov kernel.

0.1.1 Cartesian and tensor products

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

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

Given functions $f : W \rightarrow Y$ and $g : X \rightarrow Z$, the tensor product $f \otimes g : W \times X \rightarrow Y \times Z$ is defined by $(f \otimes g)_{(w, x)} = (f_w, g_x)$.

0.1.2 Indicator functions, delta measures and function-associated Markov kernels

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

For any X and any $A \subset X$, $\mathbb{1}[A]$ is the function defined by $\mathbb{1}[A]_x = \llbracket x \in A \rrbracket$. Thus $\mathbb{P}[A] = \mathbb{P}\mathbb{1}[A]$. We use square brackets to highlight the fact that $\mathbb{1}[A]$ is a function rather than a scalar.

For any X and any $x \in X$, $\delta[x]$ is the probability measure defined by $\delta[x]^i = \llbracket x = i \rrbracket$.

We define the Markov kernel $\underline{f} : X \rightarrow \Delta(\mathcal{Y})$ associated with the function $f : X \rightarrow Y$ with the matrix that sends $x \mapsto \delta_{f_x}$. Alternatively, we can define it by its rows: $\underline{f} := (\delta_{f_x})_{x \in X}$.

0.1.3 Copy maps and sequences

The copy map $\bigvee : X \rightarrow \Delta(X \times X)$ is the Markov kernel with $\bigvee_x^{(x', x'')} = \llbracket x = x' \ \& \ x = x'' \rrbracket$.

Given $X : E \rightarrow X$ and $Y : E \rightarrow Y$, the *sequence* random variable $(X, Y) : E \rightarrow X \times Y$ is defined as $\bigvee(X \otimes Y)$. That is, $(X, Y)_i = (X_i, Y_i)$.

0.1.4 Generalised random variables

It is typical to define a probability space as a probability measure along with its underlying set and its σ -algebra: $(\mathbb{P}, (E, \mathcal{E}))$. Here where E is sometimes called the sample space and \mathcal{E} is sometimes called the set of events; as we are considering discrete sets, in this paper we always have \mathcal{E} is the power set of E and we will henceforth only mention the set E .

Given a probability space (\mathbb{P}, E) , we can define *random variables* as measurable functions $X : E \rightarrow X$. The *marginal distribution* of X is given by $\mathbb{P}[X] := \mathbb{P}X$.

Here we want to consider “Markov kernel spaces”, which is a Markov kernel along with its domain and underlying set of its codomain: (\mathbb{K}, D, F) . Given such a triple, a *random variable* is a function $F \rightarrow Y$ for some vector space Y and a *state variable* is a function $D \rightarrow Y'$ for some vector space Y' . The *complete state variable* D is the identity function on D . Probabilities and conditional probabilities that we can define on the space (\mathbb{K}, D, F) usually have to be conditioned on D , but there are some exceptions.

Something that is either a random variable or a state variable is just a *variable*.

For each $d \in D$, any random variable $X : F \rightarrow X$ has a unique marginal distribution $\mathbb{K}[X|D]_d := \mathbb{K}_d X$.

To save space, we say that the marginal distribution of a sequence like (X, Y) is $\mathbb{K}[XY|D]_d$.

0.1.5 Disintegration

Conditional probabilities are *disintegrations* of probability measures. Given a probability space (\mathbb{P}, E) and random variables $X : E \rightarrow X$ and $Y : E \rightarrow Y$, the probability of X given Y is any Markov kernel $\mathbb{P}[Y|X]$ such that $\mathbb{P}[XY]^{ij} = P[X]^i P[Y|X]_i^j$. Note that this is generally non-unique. However, wherever $P_i^X > 0$, $P_{ij}^{Y|X}$ must be equal to $\frac{P_{ij}^{XY}}{P_i^X}$.

We define disintegrations of kernels analogously. Given a Markov kernel space (\mathbb{K}, D, F) , complete state variable D and variables X, Y , $\mathbb{K}[Y|XD]$ is any Markov kernel such that $\mathbb{K}[XY|D]_i^{jk} = \mathbb{K}[X|D]_i^j \mathbb{K}[Y|XD]_{ij}^k$.

As previously mentioned, in the kernel space (\mathbb{K}, D, F) there is in general no unique marginal distribution of (X, Y) and similarly there is generally no unique distribution of X conditioned on Y . However, such a distribution might exist if X and Y are independent of D .

0.1.6 Conditional independence

Given a Markov kernel space (\mathbb{K}, D, F) , and variables X, Y, Z we say X is independent of Y given Z , notated $X \perp\!\!\!\perp_{\mathbb{K}} Y|Z$ iff a version of $\mathbb{K}[X|YZ]$ exists and $\mathbb{K}[X|YZ]_i^j = \mathbb{K}[X|YZ]_{i'}^j$ for all $i, i' \in Y$.

A version of $\mathbb{K}[X|Z]$ exists iff $X \perp\!\!\!\perp_{\mathbb{K}} D|Z$ or $Z = D$, and in the former case is given by any kernel satisfying $\mathbb{K}[X|Z]_i^j = \mathbb{K}[X|DZ]_{ik}^j$ for any version of $\mathbb{K}[X|DZ]$ and all $k \in D$.

0.2 See-do models

We will first introduce *see-do models* as a type of model that functions as the basic kind of thing which we will use to examine questions in the decision theoretic, potential outcomes and graphical models approach.

See-do models can be understood as generalisations of statistical models. Statistical models are a ubiquitous type of model in statistics and machine learning that consist of a set of *states* S , and for each state the model prescribes a single probability distribution on a given set of *outcomes* O .

Definition 0.2.1 (Statistical model). A statistical model is a set of states S , a set of outcomes O and a Markov kernel $\mathbb{T} : S \rightarrow \Delta(O)$.

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

This example actually goes beyond our formal definitions here in that θ is real-valued between 0 and 1. Extending probability theory to real-valued spaces is well understood, see for example Çinlar (2011), but in that setting the existence of disintegrations on kernel spaces (section 0.1.5) is a problem to which we presently only have a partial solution. Discrete sets allow us to discuss see-do models without going into this difficulty. The price we pay is that to properly model the above problem we require θ to take on discrete values, for example restricting it to the rationals.

A see-do model adds the following structure to a statistical model:

- The state is a pair consisting of a *hypothesis* $h \in H$ and a *decision* $d \in D$; $S = H \times D$
- The outcome is a pair consisting of an *observation* $x \in X$ and a consequence $y \in Y$
- The observation is conditionally independent of the decision given the hypothesis

We can use see-do models to model situations where we have some hypotheses and the opportunity to make an observation that takes values in X . Depending on what we see, we can select a decision from a set of possibilities D , and the ultimate consequence depends probabilistically on the decision we selected as well as whichever hypothesis turns out to best describe the world.

Definition 0.2.2. A *see-do model* (\mathbb{T}, H, D, X, Y) is a Markov kernel space $(\mathbb{T}, H \times D, O)$ along with four variables: the *hypothesis* $H : H \times D \times O \rightarrow H$, the *decision* $D : H \times D \times O \rightarrow D$, the *observation* $X : H \times D \times O \rightarrow X$ and the *consequence* $Y : H \times D \times O \rightarrow Y$, all given by the projections onto the respective spaces. In addition, a see-do model must observe the conditional independence:

$$X \perp\!\!\!\perp_{\mathbb{T}} D | H \quad (5)$$

See-do models feature variables D and H that act like Dawid’s “non-stochastic regime indicators” described by Dawid (2002, 2012, 2020). In particular, see-do models induce a collection of probability measures indexed by the elements of $H \times D$, just as regime indicators induce collections of indexed probability measures. Dawid’s regime indicators seem to typically do a similar job to a decision variable D rather than a decision-hypothesis pair.

The hypothesis set is similar to the parameter set described by Lattimore and Rohde (2019) that relates pre- and post-interventional distributions. Lattimore and Rohde consider models with a prior distribution over this parameter set. A similar type of model can be created by taking the product of a prior over the hypothesis set and a see-do model. See-do models are somewhat similar to the models proposed by Savage (1954) for decision problems if we identify *states* with *hypotheses* and *acts* with *decisions*. Savage’s models consider deterministic rather than stochastic functions from acts to outcomes, and did not explicitly distinguish observations from consequences. Savage’s models themselves form the basis of the “decision theoretic” approach to causal inference set out by Heckerman and Shachter (1995) (where I use quotes to indicate that there are several distinct “decision theoretic” approaches in existence).

0.2.1 See-do models for data-driven decision problems

We can be more precise about the type of decision problem see-do models are appropriate for.

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

Then making the following additional assumptions:

1. We have random variables X, Y and D defined by $X_{xdy} = x, Y_{xdy} = y$ and $D_{xdy} = d$
2. Holding the hypothesis h fixed the observations as have the same distribution under any strategy: $\mathbb{P}_{h,\mathbb{S}}[X] = \mathbb{P}_{h,\mathbb{S}'}[X]$ for all $h, \mathbb{S}, \mathbb{S}'$
3. The chosen strategy is a version of the conditional probability of decisions given observations: $\mathbb{S} = \mathbb{P}_{h,\mathbb{S}}[D|X]$
4. There exists some strategy \mathbb{S} that is strictly positive
5. For any $h \in H$ and any two strategies \mathbb{Q} and \mathbb{S} , $\mathbb{P}_{h,\mathbb{Q}}[Y|DX] = \mathbb{P}_{h,\mathbb{S}}[Y|DX]$

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

Proof. Consider some probability $\mathbb{P} \in \Delta(X \times D \times Y)$. By the definition of disintegration (section 0.1.5), we can write

$$\mathbb{P}[XDY]^{ijk} = \mathbb{P}[X]^i \mathbb{P}[D|X]_i^j \mathbb{P}[Y|XD]_{ij}^k \quad (6)$$

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

$$\mathbb{T}_{hj}^{kl} = \mathbb{T}[XY|HD]_{hj}^{kl} \quad (7)$$

$$= \mathbb{P}_{h,\mathbb{S}}[X]^k \mathbb{P}_{h,\mathbb{S}}[Y|XD]_{kj}^l \quad (8)$$

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

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

We then have

$$\mathbb{T}[X'|H'D']_{hj}^k = \mathbb{T}X'_{hj}^k \quad (9)$$

$$= \sum_l \mathbb{T}_{hj}^{kl} \quad (10)$$

$$= \mathbb{P}_{h,\mathbb{S}}[X]^k \quad (11)$$

$$= \mathbb{T}[X'|H'D']_{hj'}^k \quad (12)$$

Thus $X' \perp\!\!\!\perp_{\mathbb{T}} D'|H'$ and so $\mathbb{T}[X'|H']$ exists (section 0.1.6)

Applying Equation 6 to $\mathbb{P}_{h,\mathbb{S}}$:

$$\mathbb{P}_{h,\mathbb{S}}[\mathbf{XDY}]^{ijk} = \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (13)$$

$$= \mathbb{P}_{h,\mathbb{S}}[\mathbf{X}]^i \mathbb{P}_{h,\mathbb{S}}[\mathbf{Y}|\mathbf{XD}]_{ij}^k \quad (14)$$

$$= \mathbb{P}_{h,\mathbb{S}}[\mathbf{D}|\mathbf{X}]_i^j \mathbb{T}[\mathbf{XY}|\mathbf{HD}]_{hj}^{ik} \quad (15)$$

$$= \mathbb{S}_i^j \mathbb{T}[\mathbf{XY}|\mathbf{HD}]_{hj}^{ik} \quad (16)$$

$$= \mathbb{T}[\mathbf{X}|\mathbf{HD}] \mathbb{S}_i^j \quad (17)$$

If \mathbb{V} is strictly positive then \mathbb{T} is unique.

Fix some $h \in H$ and strictly positive \mathbb{S} and let \mathbb{T}_h be the unique 2-comb such that Equation 8 holds for $\mathbb{P}_{h,\mathbb{S}}$. Consider an arbitrary alternative strategy \mathbb{Q} . Then there is some 2-comb \mathbb{U} , not necessarily unique, such that Equation 8 holds for $\mathbb{P}_{h,\mathbb{Q}}$. □

References

- Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- 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.
- Philip Dawid. The Decision-Theoretic Approach to Causal Inference. In *Causality*, pages 25–42. John Wiley & Sons, Ltd, 2012. ISBN 978-1-119-94571-0. doi: 10.1002/9781119945710.ch4. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119945710.ch4>.
- 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>.
- Benjamin Heymann, Michel de Lara, and Jean-Philippe Chancelier. Causal Inference Theory with Information Dependency Models. August 2021. URL <https://arxiv.org/abs/2108.03099v2>.
- 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.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2 edition, 2009.

Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, New York, 1 edition edition, May 2018. ISBN 978-0-465-09760-9.

Leonard J. Savage. *The Foundations of Statistics*. Wiley Publications in Statistics, 1954.

Appendix: