
Causal Statistical Decision Problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We develop the notion of a causal statistical decision problem as an extension
2 of the statistical decision theory of Wald. Suppose we have a dataset and some
3 set of available decisions. Assume we know what state we would like the world
4 to occupy, but we are uncertain about how our decisions affect the state of the
5 world. We introduce the notion of *consequences* that relate decisions to states
6 of the world, and *causal theories* that relate observations to consequences. A
7 strength of this perspective is that it is not motivated by any notion of a “true
8 cause” or “causal effect”. We connect causal statistical decision problems to
9 statistical decision problems and show that two leading approaches to causality -
10 Causal Bayesian Networks and Potential Outcomes - have natural representations
11 as causal theories. We argue that the causal theory associated with a CBN may
12 be considered incomplete and discuss how different extensions can lead to very
13 different properties.

14 1 Introduction

15 The decision theoretic approach to statistics casts statistical problems in terms of learning to output
16 decisions that minimise a loss rather than learning true properties of a data generating distribution.
17 Statistical decision theory plays a role of fundamental importance in modern machine learning;
18 loss functions underpin the development of algorithms, and the analysis of losses is critical to the
19 theoretical treatment of learning algorithms.

20 It is widely accepted that problems of causal inference are different to statistical problems. Causal
21 problems are held to demand causal knowledge that is not in the vocabulary of statistical problems
22 [Pearl, 2009, Cartwright, 1994]. There are two leading approaches to formalising “causal knowledge”
23 and posing data-driven causal problems: one based on Causal Bayesian Networks and the other on
24 Potential Outcomes.

25 Causal Bayesian Networks (CBNs) posit that there are causal relationships among a set of random
26 variables that can be encoded by a directed acyclic graph (DAG). An investigator with access to the
27 true graph and a joint probability distribution over all the variables present in that graph can calculate
28 a wide variety of causal effects, and partial access to these objects will enable to partial knowledge
29 of causal effects. A causal effect in this framework is tied to the intuitive notion of “the result of
30 intervening to set particular variables to particular values”.

31 Potential Outcomes (PO) posits a large joint distribution over observed variables X , Y and partly
32 unobserved “potential outcome” variables X_0 , Y_1 and so forth. A potential outcome variable Y_i is
33 interpreted as “the value of Y that would be observed if the action identified by i were taken”. Under
34 some conditions, an investigator with access to a joint distribution over observed variables may be
35 able to infer certain properties of the distribution over potential outcome variables such as $\mathbb{E}[X_i]$.

36 Queries in the CBN framework may be concerned with identification of causal effects given a graph
37 and a probability distribution [Tian and Pearl, 2002], or with the determination of the true causal
38 graph given just a probability distribution [Spirtes et al., 2000]. Queries in the PO framework usually

concern identification of properties of the distribution of potential outcome variables known as *treatment effects* given a dataset and certain assumptions about this distribution [Rubin, 2005, Robins and Richardson, 2010]. In both cases, these queries fit the paradigm of “determining true properties of nature” rather than “learning to output a decision that minimises a loss”.

The first contribution of this paper is the notion of a *causal statistical decision problem* (CSDP) that proceeds from a natural extension of an ordinary statistical decision problem (SDP) introduced by [Wald, 1950]. We suppose that, in contrast to an ordinary SDP where we have known preferences over (decision, state of nature) pairs, we know only our preferences over the *outcomes* of decisions, which we represent with a utility function. Uncertainty over the consequences of decisions is represented by a *causal theory* that connects observed data with *consequence maps*.

We show by a reduction that results concerning standard SDPs are also true of (at least) a subset of CSDPs. We also show that both Causal Bayesian Networks and joint distributions over potential outcomes have a natural representation as causal theories. Together these results show, for example, that the class of Bayes decision functions is a complete class for CSDPs based on Causal Bayesian Networks provided certain conditions on the utility and size of the available set of decisions are met.

The notion of a causal theory presented here can naturally represent models cast in terms of CBNs or POs, but there are many causal theories that cannot easily be represented by either. We discuss a question motivated by this more general perspective: *given a CBN with observable predictions, what should be assumed when the data doesn’t match these predictions?* We show that different answers to this question yield widely divergent conclusions.

A key strength of our perspective is the possibility of theoretical treatment of causal learning from a viewpoint that is agnostic about the nature of “causal knowledge”. Causal knowledge is a tricky domain from philosophy to practice, and there are many proposals for causal assumptions that do not neatly fit in either the CBN or PO camps [Bongers et al., 2016, Dawid, 2010, Bengio et al., 2019]. The theory presented here is capable of posing questions such as “does a proposed causal learning method work?” without first requiring commitments on the nature of causal knowledge. Substantial progress in machine learning has been the result of developing generic principles and learning techniques that are relevant to many datasets from many domains and are less reliant on the judgement of domain experts. We believe this separation of concerns is crucial to the advancement of generic techniques of causal learning.

Our approach is similar to that of Dawid [2012], but where he takes a “bottom-up” approach of developing a decision theoretic answer to particular causal questions, our approach is “top-down”, proceeding from a general account of a causal problem to the particular objects needed to answer it. It also shares similarities with Causal Decision Theory developed by Lewis [1981], though the connection with statistical decision theory is better understood at this point.

2 Definitions & Notation

We use the following standard notation: $[N]$ refers to the set of natural numbers $\{1, \dots, N\}$. Sets are ordinary capital letters X while σ -algebras are calligraphic capitals \mathcal{X} and random variables are sans serif capitals $X : _ \rightarrow X$. The calligraphic \mathcal{G} refers to a directed acyclic graph rather than a σ -algebra. Sets of probability measures or stochastic maps are script capitals: $\mathcal{H}, \mathcal{T}, \mathcal{J}$.

A measurable space (E, \mathcal{E}) is a set E and a σ -algebra $\mathcal{E} \subset \mathcal{P}(E)$ containing the measurable sets. A probability measure $\mu \in \Delta(\mathcal{E})$ is a nonnegative map $\mathcal{E} \rightarrow [0, 1]$ such that $\mu(\emptyset) = 0$, $\mu(E) = 1$ and for countable $\{E_i\} \in \mathcal{E}$, $\mu(\cup_i E_i) = \sum_i \mu(E_i)$. We assume all measurable spaces discussed are standard. That is, they are isomorphic to either a subset of \mathbb{N} with the discrete σ -algebra, or \mathbb{R} with the Borel σ -algebra.

Given two measurable spaces (E, \mathcal{E}) and (F, \mathcal{F}) , a *Markov kernel* or *stochastic map* $K : E \rightarrow \Delta(\mathcal{F})$ is a map where $x \mapsto K(x; B)$ is \mathcal{E} -measurable for every $B \in \mathcal{F}$ and $B \mapsto K(x; B)$ is a probability measure on (F, \mathcal{F}) for every $x \in E$. Abusing notation somewhat, we will write the set of Markov kernels of type $E \rightarrow \Delta(\mathcal{F})$ as $\Delta(\mathcal{F})^D$.

If we have two random variables $X : _ \rightarrow X$ and $Y : _ \rightarrow Y$, the conditional probability $P(Y|X)$ is a Markov kernel $X \rightarrow \Delta(\mathcal{Y})$. Formally, given $\mu \in \Delta(\mathcal{E})$ and a sub- σ -algebra $\mathcal{E}' \subset \mathcal{E}$, there is a Markov kernel $\mu|_{\mathcal{E}'} : E \rightarrow \Delta(\mathcal{E})$ such that for $A \in \mathcal{E}$ and $B \in \mathcal{E}'$, $\int_B \mu|_{\mathcal{E}'}(y; A) d\mu(y) = \mu(A \cap B)$.

91 $\mu_{|\mathcal{E}'}$ is a *conditional probability distribution* with respect to \mathcal{E}' . This result may not hold if (E, \mathcal{E}) is
 92 not a standard measurable space [Çinlar, 2011].

93 Given a set of random variables $\mathbf{X} = \{X^i\}_{i \in [N]}$ with domain (E, \mathcal{E}) , $\mu_{\mathbf{X}} : E \rightarrow \Delta(\mathcal{E})$ is a
 94 conditional probability distribution with respect to the σ -algebra generated by \mathbf{X} : $\sigma(\cup_{i \in [N]} \sigma(\mathcal{X}^i))$.
 95 We will use this subscript notation rather than the more common bar notation (e.g. $\mu(\cdot | \mathbf{X})$) to express
 96 conditional probability from here onwards.

97 Two Markov kernels $K : E \rightarrow \Delta(\mathcal{F})$ and $K' : E \rightarrow \Delta(\mathcal{F})$ are μ -almost surely equivalent given
 98 $\mu \in \Delta(\mathcal{E})$ if for all $A \in \mathcal{E}, B \in \mathcal{F}$, $\int_A K(x; B) d\mu = \int_A K'(x; B)$.

99 **Kernel products:** Kernel products allow common operations to be written compactly. The notation
 100 here borrows heavily from Çinlar [2011] and Fong [2013]. More details can be found in Appendix
 101 A. For the following, assume $K : E \rightarrow \Delta(\mathcal{F})$, $L : F \rightarrow \Delta(\mathcal{G})$, and $M : G \rightarrow \Delta(\mathcal{H})$ are Markov
 102 kernels, μ is a probability measure on (E, \mathcal{E}) .

103 The *kernel-kernel* product KL is a Markov kernel $E \rightarrow \Delta(\mathcal{G})$ such that $KL(x; B) :=$
 104 $\int_F K(x; dy) L(y; B)$, $x \in E, B \in \mathcal{G}$. Kernel-kernel products are associative: $(KL)M =$
 105 $K(LM)$.

106 The *measure-kernel* product of μ and K , μK is a probability measure on (F, \mathcal{F}) such that $\mu K(B) =$
 107 $\int_E \mu(dx) K(x; B)$, $B \in \mathcal{F}$. Measure-kernel products are also associative: $(\mu K)L = \mu(KL)$.

108 **Special kernels:** $I_{(E)}$ is the identity kernel $E \rightarrow \Delta(\mathcal{E})$ defined by $x \mapsto \delta_x$. It has the properties
 109 $\mu I_{(E)} = \mu$, $K I_{(F)} = K$, $I_{(E)} K = K$.

110 Given some measurable function $g : E \rightarrow F$, the kernel $F_g : E \rightarrow \Delta(\mathcal{F})$ is defined by $x \mapsto \delta_{g(x)}$. It
 111 is easy to check that $F_g F_g = F_g$. For $\mu \in \Delta(\mathcal{E})$, $\mu F_g(A) = \mu(g^{-1}(A))$. This notation allows us to
 112 consistently represent a marginal distribution μF_X and a marginal kernel κF_X .

113 Given $\mu \in \Delta(\mathcal{E})$, $\mu \curlyvee (I_{(E)} \otimes K)$ is a distribution in $\Delta(\mathcal{E} \otimes \mathcal{F})$ given by

$$\mu \curlyvee (I_{(E)} \otimes K)(A \times B) = \int_A K(x; B) d\mu(x) \quad \forall A \in \mathcal{E}, B \in \mathcal{F} \quad (1)$$

114 The symbol \curlyvee is read “splitter”.

115 3 Causal Statistical Decision Problems

	SDPs	CSDPs
State of the world	\mathcal{H} , hypothesis class	\mathcal{T} , causal theory
Observations	\mathbf{X}	\mathbf{X}
Decisions	D	D
Known preferences	ℓ , loss	U , generalised utility
Derived preferences	ℓ , loss	L , causal loss

Table 1: Comparison of SDPs and CSDPs

116 We develop causal statistical decision problems (CSDPs) inspired by statistical decision problems
 117 (SDPs) of Wald [1950]. CSDPs differ from SDPs in that our preferences (i.e. utility or loss) are
 118 known less directly in former case. We show that every SDP can be represented by a CSDP and that
 119 the converse is sometimes but not always possible. We show that an analogue of the fundamental
 120 *complete class theorem* of SDPs applies to the class of CSDPs that can be represented by SDPs, but
 121 whether such a theorem applies more generally is an open question.

122 Following [Ferguson, 1967], we consider SDPs and CSDPs to represent normal form two person
 123 games. At the most abstract level the games represent the options and possible payoffs available to
 124 the decision maker, and this representation allows us to compare the two types of problem. In their
 125 more detailed versions, CSDPs and SDPs differ in their representation of the state of the world and in
 126 the type of function that represents preferences. These differences are summarised in Table 1.

127 **Definition 3.1** (Normal form two person game). A normal form game is a triple $\langle \mathcal{S}, A, L \rangle$ where \mathcal{S}
 128 and A are arbitrary sets and $L : \mathcal{S} \times A \rightarrow [0, \infty)$ is a loss function.

129 The set \mathcal{S} is a set of possible states that the environment may occupy and A is a set of actions
 130 the decision maker may take. The decision maker seeks an action in A that minimises the loss L .
 131 Generally there is no action that minimises the loss for all environment states. A minimax solution is
 132 an action that minimises the worst case loss: $a_{mm}^* = \arg \min_{a \in A} [\sup_{s \in \mathcal{S}} L(s, a)]$.

133 If the set \mathcal{S} is equipped with a σ -algebra \mathcal{S} and a probability measure $\xi \in \Delta(\mathcal{S})$ which we
 134 will call a “prior”, a Bayes solution minimizes the expected risk with respect to ξ : $a_{ba}^* =$
 135 $\arg \min_{a \in A} \int_{\mathcal{S}} L(s, a) \xi(ds)$.

136 **Definition 3.2** (Admissible Action). Given a normal form two person game $\langle \mathcal{S}, A, L \rangle$, an action
 137 $a \in A$ is *strictly better* than $a' \in A$ iff $L(s, a) \leq L(s, a')$ for all $s \in \mathcal{S}$ and $L(s_0, a) < L(s_0, a')$ for
 138 some $s_0 \in \mathcal{S}$. If only the first holds, then a is as good as a' . An *admissible action* is an action $a \in A$
 139 such that there is no action strictly better than a .

140 **Definition 3.3** (Complete Class). A class C of decisions is a *complete class* if for every $a \notin C$ there
 141 is some $a' \in C$ that is strictly better than a .

142 C is an *essentially complete* class if for every $a \notin C$ there is some $a' \in C$ that is as good as a .

143 **Definition 3.4** (Reduction). A normal form two person game $\alpha = \langle \mathcal{S}^\alpha, A, L^\alpha \rangle$ can be reduced
 144 to a different game sharing the same action set $\beta = \langle \mathcal{S}^\beta, A, L^\beta \rangle$ if there is a surjective function
 145 $g : \mathcal{S}^\alpha \rightarrow \mathcal{S}^\beta$ such that for every $a \in A$, $s \in \mathcal{S}^\alpha$, $L^\alpha(s, a) = L^\beta(g(s), a)$.

146 Because CSDPs and SDPs posit states of nature of different types, they cannot represent exactly the
 147 same game. Reduction is a notion that allows us to say tht the game represented by a CSDP and an
 148 SDP are “essentially” the same. Reduction preserves the important properties of admissibility and
 149 completeness as shown in Appendix B.

150 A statistical decision problem represents a normal form two-person game where the available actions
 151 are *decision functions* that output a decision given data, the states of the environment are associated
 152 with probability measures on some measurable space and we assume a loss expressing preferences
 153 over decisions and states is known.

154 **Definition 3.5** (Statistical Decision Problem). A statistical decision problem (SDP) is a triple
 155 $\langle \mathcal{H}, D, \ell \rangle$. $\mathcal{H} \subset \Delta(\mathcal{E})$ is a hypothesis class representing possible states of the environment, D
 156 is the set of available decisions, $X : (E, \mathcal{E}) \rightarrow (X, \mathcal{X})$ is a random variable representing the in-
 157 formation available for the statistician to make a decision and $\ell : \mathcal{H} \times D \rightarrow [0, \infty)$ is a loss
 158 function.

159 Denote by \mathcal{J} the set of decision kernels $X \rightarrow \Delta(D)$. Recall that $\mu F_X(A) = \mu(X^{-1}(A))$. For $J \in \mathcal{J}$
 160 and $\mu \in \mathcal{H}$, the risk $R : \mathcal{J} \times \mathcal{H} \rightarrow [0, \infty)$ is defined as $R(J, \mu) = \int_D \ell(\mu, y) \mu F_X J(dy)$. The triple
 161 $\langle \mathcal{H}, \mathcal{J}, R \rangle$ forms a two player normal form game.

162 The loss function ℓ expresses preferences over (state, decision) pairs. However, it may be the case
 163 that our preferences don’t naturally apply directly to such pairs. For a doctor deciding whether to
 164 prescribe a treatment to a patient, it is clear that this patient being healthy in the future is preferable
 165 to them being sick. This motivates the definition of a causal statistical decision problem, utilising a
 166 preferences defined over outcomes represented by a generalised utility rather than (state, decision)
 167 pairs. A generalised utility, as opposed to an ordinary utility, allows for a simple reduction from
 168 CSDPs to SDPs. In order to compute the loss associated with a decision a map from decisions to
 169 outcomes is required, which we term a *consequence*.

170 **Definition 3.6** (Consequences). Given a measurable outcome space (F, \mathcal{F}) and a measurable decision
 171 space (D, \mathcal{D}) , a Markov kernel $\kappa : D \rightarrow \Delta(\mathcal{F})$ is a *consequence mapping*, or just a consequence for
 172 short.

173 **Definition 3.7** (Causal state). Given a consequence $\kappa : D \rightarrow \Delta(\mathcal{F})$, a measurable observation space
 174 (E, \mathcal{E}) and some distribution $\mu \in \Delta(\mathcal{E})$, the pair $\tau := (\kappa, \mu)$ is a *causal state* on D, E and F . We
 175 refer to κ as the consequence and μ as the observed state.

176 We allow the “observation” space E and the “outcome” space F to differ as it may be desirable
 177 to avoid modelling consequences on variables that are observed but irrelevant to preferences (see
 178 Theorems B.6 and B.5). In practice these spaces often coincide.

179 **Definition 3.8** (Causal Theory). A causal theory \mathcal{T} is a set of causal states sharing the same decision,
 180 observation and outcome spaces. We abuse notation to assign the “type signature” $\mathcal{T} : E \times D \rightarrow F$
 181 for a causal theory with observed distributions in $\Delta(\mathcal{E})$ and consequences $D \rightarrow \Delta(\mathcal{F})$.

Definition 3.9 (Causal Statistical Decision Problem). A causal statistical decision problem (CSDP) is a tuple $\langle (\mathcal{T}, (E, \mathcal{E})\mathbf{X}), (D, \mathcal{D}), (U, (F, \mathcal{F})) \rangle$. \mathcal{T} is a causal theory on D, E and F , D is the decision set, $\mathbf{X} : (E, \mathcal{E}) \rightarrow (X, \mathcal{X})$ is a random variable representing the given information and $U : \Delta(\mathcal{F} \otimes \mathcal{D}) \rightarrow \mathbb{R}$ is a generalised utility expressing preference over joint distributions of decisions and outcomes which we assume is bounded above.

From the generalised utility U we can define a loss $L : \mathcal{T} \times \Delta(\mathcal{D}) \rightarrow [0, \infty]$ by

$$L((\kappa, \mu), \gamma) := \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \vee (I_{(D)} \otimes \kappa)) - U(\gamma \vee (I_{(D)} \otimes \kappa)) \quad (2)$$

For $(\kappa, \mu) \in \mathcal{T}$ and $\gamma \in \Delta(\mathcal{D})$. This is well defined wherever U is bounded above. Note that L does not depend on the data generating distribution μ ; henceforth we will suppress this argument and write $L(\kappa, \gamma) := L((\kappa, \mu), \gamma)$.

Given a decision function $J \in \mathcal{J}$ and $(\kappa, \mu) \in \mathcal{T}$, we define the risk $R : \mathcal{J} \times \mathcal{T} \rightarrow [0, \infty)$ by $R(J, \kappa, \mu) := L(\kappa, \mu F_{\mathbf{X}} J)$. The triple $\langle \mathcal{T}, \mathcal{J}, R \rangle$ is a normal form two person game.

If there exists some measurable $u : F \times D \rightarrow \mathbb{R}$ such that for all $\xi \in \Delta(\mathcal{F} \otimes \mathcal{D})$, $U(\xi) = \mathbb{E}_{\xi}[u]$ then we call U an ordinary utility. An ordinary induces a loss $L(\kappa, \gamma) = \mathbb{E}_{\gamma}[l^{\kappa}]$ where $l^{\kappa} : D \rightarrow [0, \infty)$ is defined by

$$l^{\kappa}(d) := \sup_{\gamma' \in \Delta(\mathcal{D})} \mathbb{E}_{\gamma' \vee (I_{(D)} \otimes \kappa)}[u] - \mathbb{E}_{\kappa(d; \cdot)}[u(\cdot, d)] \quad (3)$$

Reduction from a CSDP to an ordinary SDP is sufficient to import results from statistical decision theory such as Theorem 3.10. The complete class theorem along with Corollary 3.14 show that, at least for CSDPs with countable decision sets, finite causal theories and ordinary utilities, any admissible decision rule is a Bayes rule given *some* prior on \mathcal{T} . Stronger versions of Theorem 3.10 exist for SDPs, and stronger versions are likely to exist for CSDPs as well.

Theorem 3.10 (Complete class theorem (CSDP)). *Given an CSDP $\alpha := \langle (\mathcal{T}, E), D, \mathbf{X}, U \rangle$ with risk R , if there is a reduction to an SDP $\beta := \langle (\mathcal{H}, F), D, \mathbf{Y}, \ell \rangle$ with risk R' such that $|\mathcal{H}| < \infty$ and $\inf_{J \in \mathcal{J}, \mu \in \mathcal{H}} R'(J, \mu) < -\infty$, then the set of all Bayes decision functions is a complete class for β and the set of all admissible Bayes decision functions is a minimal complete class for β .*

Proof. This follows from Lemmas B.2 and B.3. See appendix B. \square

Any statistical decision problem can be reduced to a CSDP featuring a causal theory where decisions have no effect.

Theorem 3.11. *Every SDP $\langle (\mathcal{H}, E, \mathbf{X}), D, \ell \rangle$ can be reduced to a CSDP.*

Proof. Choose a theory that matches every probability measure $\mu \in \Delta(\mathcal{E})$ with a consequence map that itself always yields μ . Then construct a generalised utility U that induces an identical risk. See Appendix B. \square

A CSDP cannot, in general, be reduced to a statistical decision problem - for example, if we choose the utility to be the variance of some random variable we may be able to achieve higher utility through a randomised decision than any nonrandomised decision under conditions where a regular SDP cannot have this property (see Example B.9 in Appendix B). It is an open question whether this reduction is generally possible if the problem features a ordinary utility.

Theorem 3.12 reduces a CSDP to a SDP by associating each pair (μ, κ) in the causal theory with a distribution over $E \times F \times D$. This may be possible by finding a distribution on the product space for which μ is a marginal probability and κ is a conditional distribution, but this is not necessary.

Theorem 3.12. *Given a CSDP $\beta = \langle (\mathcal{T}, E, \mathbf{X}), D, (U, F) \rangle$ where U is an ordinary generalised utility, let $\mathcal{K} = \{\kappa | (\kappa, \mu) \in \mathcal{T}\}$ be the set of consequences. β is reducible to a statistical decision problem on the measurable space $(E \times F \times D, \mathcal{E} \otimes \mathcal{F} \otimes \mathcal{D})$ if there is some surjective map $m : \Delta(\mathcal{F} \otimes \mathcal{D}) \rightarrow \mathcal{K}$.*

Proof. We construct the map h based on the map m and show that, given an ordinary utility U , it is possible to construct a loss ℓ such that the resulting SDP features the same risk assignments as the original CSDP. See Appendix B. \square

226 **Corollary 3.13.** *If the cardinality of $\Delta(\mathcal{F} \otimes \mathcal{D})$ is at least as large as the cardinality of the set of*
 227 *Markov kernels $D \rightarrow \Delta(\mathcal{F})$ then an CSDP with an ordinary utility can always be reduced to a SDP.*

228 A major open question, then, is if (E, \mathcal{E}) and (D, \mathcal{D}) are standard measurable spaces, the conditions
 229 for Corollary 3.13 hold in general. Corollary 3.14 shows that the reduction can be made in general if
 230 D is a denumerable set.

231 **Corollary 3.14.** *A CSDP $\langle (\mathcal{T}, (E, \mathcal{E}), \mathbf{X}), (D, \mathcal{D}), (U, (F, \mathcal{F})) \rangle$ where D is a denumerable set and*
 232 *U is an ordinary generalised utility can be reduced to a statistical decision problem.*

233 *Proof.* Take some probability measure $\pi \in \Delta(\mathcal{D})$ such that $\pi(\{y\}) > 0$ for all $y \in D$. Such a π
 234 exists by the denumerability of \mathcal{D} . The map $m : \Delta(\mathcal{F} \otimes \mathcal{D}) \rightarrow \mathcal{K}$ given by $m(\xi)(y; A) := \frac{\xi(A \times \{y\})}{\pi(\{y\})}$
 235 is surjective. The result follows from Theorem 3.12. \square

236 4 Causal Bayesian Networks

237 A Causal Bayesian Network (CBN) is a directed acyclic graph (DAG) \mathcal{G} containing a set of nodes
 238 $\{X^i\}_{i \in [N]}$ which we identify with random variables on some space (E, \mathcal{E}) . Given a decision $y \in$
 239 D (called a *do-intervention* in other treatments) and a distribution $\mu \in \Delta(\mathcal{E})$ that is *compatible*
 240 (Definition 4.1) with \mathcal{G} , \mathcal{G} induces an *interventional* distribution $\mu^{\mathcal{G}, y}$. The set of pairs $(\mu, y \mapsto \mu^{\mathcal{G}, y})$
 241 for μ compatible with \mathcal{G} is a causal theory $\mathcal{T}_{\mathcal{G}}$.

242 In all following discussion, we assume the observed data represented by \mathbf{X} is a sequence of indepen-
 243 dent and identically distributed random variables $\mathbf{X} = (X_t)_{t \in T}$. We identify distributions over the
 244 sequence \mathbf{X} with distributions over the initial observation X_0 and subsequently drop the subscript.

245 The CBN convention is to denote an interventional distribution with $\mu(\cdot | do(X^i = a))$. Here we
 246 associate every allowable set of *do* statements with an element of the decision space (D, \mathcal{D}) equipped
 247 with random variables $\{D^i\}_{i \in [N]}$ such that for $y \in D$, $\mu^y(\cdot) := P(\cdot | [do(X^j = D^j(y))]_{j \in [N]})$. The
 248 special element $*$ corresponds to a passive intervention which is denoted by the absence of a *do*()
 249 statement in regular CBN notation.

250 **Definition 4.1** (Compatibility). Given a DAG \mathcal{G} , *d-separation* is a ternary relation amongst sets of
 251 nodes the details for which we refer readers to Pearl [2009]. For a set of nodes $\{X^i\}_{i \in [N]}$ we write
 252 $X^i \perp_{\mathcal{G}} X^j | \mathbf{X}$ to say X^i is d-separated in \mathcal{G} from X^j by $\mathbf{X} \subset \{X^i\}_{i \in [N]}$.

253 Given a measurable space (E, \mathcal{E}) , $\mu \in \Delta(\mathcal{E})$ and a set of random variables $\{X^i\}_{i \in [N]}$ on E , X^i is
 254 independent of X^j conditional on \mathbf{X} if $\mu|_{\mathbf{X}} \Upsilon (F_{X^i} \otimes F_{X^j}) = \mu|_{\mathbf{X}} F_{X^i} \mu|_{\mathbf{X}} F_{X^j}$, μ -almost surely. This
 255 is written $X^i \perp_{\mu} X^j | \mathbf{X}$.

256 μ is compatible with \mathcal{G} if $X^i \perp_{\mathcal{G}} X^j | \mathbf{X} \implies X^i \perp_{\mu} X^j | \mathbf{X}$

257 **Definition 4.2** (Causal Bayesian Network).

258 Consider a directed acyclic graph \mathcal{G} with nodes $\mathbf{X} = \{X^i | i \in [N]\}$, a measurable space (E, \mathcal{E}) and
 259 a set of random variables $X^i : E \rightarrow X^i$ and $X = \times_{i \in [N]} X^i$ along with decision space (D, \mathcal{D}) and
 260 random variables $\{D^i\}_{i \in [N]}$ where $D^i : D \rightarrow X^i \cup \{*\}$.

261 Given any $y \in D$ let $S(y) \subset [N]$ be the set of all indices i such that $D^i(y) \neq *$. Let $\mathcal{H}_{\mathcal{G}} \subset \Delta(\mathcal{X})$
 262 be the set of distributions compatible with \mathcal{G} . Given arbitrary $\mu \in \mathcal{H}_{\mathcal{G}}$ and $y \in D$ the \mathcal{G}, μ, y -
 263 interventional distribution denoted $\mu^{\mathcal{G}, y}$ is given by the following three conditions:

- 264 1. $\mu^{\mathcal{G}, y}$ is compatible with \mathcal{G}
- 265 2. For all $i \in S(y)$, $\mu^{\mathcal{G}, y} F_{X^i} = \delta_{D^i(y)} F_{X^i}$
- 266 3. For all $i \notin S(y)$, $\mu^{\mathcal{G}, y}_{\text{Pa}_{\mathcal{G}}(X^i)} F_{X^i} = \mu|_{\text{Pa}_{\mathcal{G}}(X^i)} F_{X^i}$, $\mu^{\mathcal{G}, y}$ -almost surely

267 $\text{Pa}_{\mathcal{G}}(X^i)$ are the parents of X^i with respect to the graph \mathcal{G} and $\mu|_{\text{Pa}_{\mathcal{G}}(X^i)}$ is the conditional probability
 268 with respect to μ and the σ -algebra generated by the set $\text{Pa}_{\mathcal{G}}(X^i)$. Recall that $\mu \Upsilon (\otimes_{i \notin S(y)} F_{X^i})$ is the
 269 joint distribution of $\{X^i | i \in S(y)\}$.

A CBN has a graph \mathcal{G} with edges $\{V^i\}_{i \in [N]}$, random variables $\{X^i\}_{i \in [N]}$ and decision variables $\{D^i\}_{i \in [N]}$ which are all associated with one another in the obvious way. It would be nice to have a simple way of expressing this bundle of things and the corresponding associations

270 To establish that the map $\kappa^{\mathcal{G}, \mu} : D \rightarrow \Delta(\mathcal{X})$ given by $y \mapsto \mu^{\mathcal{G}, y}$ is a consequence map, we must
 271 shown that it is measurable with respect to the σ -algebra generated by the set of variables D^i ; this is
 272 shown by Theorem C.1 provided in Appendix C. Defining $\mathcal{H}_{\mathcal{G}} \subset \Delta(\mathcal{X})$ to be the set of distributions
 273 compatible with \mathcal{G} , the set of pairs $\{(\mu, \kappa^{\mu}) | \mu \in \mathcal{H}_{\mathcal{G}}\}$ is the causal theory $\mathcal{T}_{\mathcal{G}}$.

274 **Extending the theory induced by a CBN** The causal theory $\mathcal{T}_{\mathcal{G}}$ defined above associates a conse-
 275 quence with every probability distribution compatible with \mathcal{G} but not every probability distribution in
 276 $\Delta(\mathcal{X})$. It is arguably not reasonable to assume *a priori* that the conditional independences implied
 277 by \mathcal{G} hold in the observed data. We might therefore regard the theory $\mathcal{T}_{\mathcal{G}}$ to be incomplete, and seek
 278 some extension of the theory for distributions not in $\mathcal{H}_{\mathcal{G}}$.

279 **Example 4.3** (Extension of a CBN). Consider the graph $\mathcal{G} = C \rightarrow A \rightarrow B$, which implies a
 280 single conditional independence: $C \perp\!\!\!\perp B | A$.

281 Suppose the three associated random variables A, B and C each take values in $\{0, 1\}$ and suppose (un-
 282 realistically) we know all μ in the set of possible joint distributions \mathcal{H} share the marginal distribution
 283 $\mu F_B := \zeta$ and the conditional distribution $\mu_{|\{A\}} F_B = \iota$ and C is “almost” independent of B given A:

$$\max_{x \in \{0,1\}^3, y \in \{0,1\}} |\mu_{|\{A,C\}} F_B(x; \{y\}) - \iota(x; \{y\})| < \epsilon \quad (4)$$

284 Suppose that only interventions on A are possible and the problem supplies a generalised utility
 285 such that, overloading B, $U(\xi) = \mathbb{E}_{\xi}[B]$. For convenience, we restrict our attention to the subset of
 286 decisions $D' = \{y | D_B(y) = D_C(y) = *\}$ and consequence maps marginalised over A and C. Define
 287 $\kappa^{\mathcal{G}}$ by

$$\kappa^{\mathcal{G}}(y; Z) := \begin{cases} \iota(D_A(y); Z) & D_A(y) \neq * \\ \zeta(Z) & D_A(y) = * \end{cases} \quad (5)$$

288 It can be verified that the causal theory $\mathcal{T}_{\mathcal{G}}$ induced by \mathcal{G} and the set of compatible distributions
 289 $\mathcal{H}_{\mathcal{G}} \subset \mathcal{H}$ is the set of pairs $\{(\nu, \kappa^{\mathcal{G}}) | \nu \in \mathcal{H}_{\mathcal{G}}\}$.

290 Consider two options for extending this to distributions $\nu \in \mathcal{H}$ but not in $\mathcal{H}_{\mathcal{G}}$, noting that one could
 291 imagine many possibilities: $\mathcal{T}_{\mathcal{G}}^{\subseteq}$ is the union of causal theories given by all graphs \mathcal{G}' on $\{A, B, C\}$

292 such that $\mathcal{G} \subset \mathcal{G}'$ (in this case, just \mathcal{G} and $C \xrightarrow{\quad} A \xrightarrow{\quad} B$), and $\mathcal{T}_{\mathcal{G}}^{\circ}$ is the union of causal theories
 293 given by the all DAGs on the set of nodes $\{A, B, C\}$.

294 The theory $\mathcal{T}_{\mathcal{G}}^{\subseteq}$ is given by $\mathcal{T}_{\mathcal{G}} \cup \{(\nu, \eta^{\nu}) | \nu \in \mathcal{H} \setminus \mathcal{H}_{\mathcal{G}}\}$ where

$$\eta^{\nu} := \begin{cases} (y; Z) \mapsto \sum_{c \in \{0,1\}} \nu F_C(\{c\}) \nu_{|\{A,C\}} F_B(D_A(y), c; Z) & D_A(y) \neq * \\ \zeta(Z) & D_A(y) = * \end{cases} \quad (6)$$

295 $\mathcal{T}_{\mathcal{G}}^{\circ}$ is the set of states associated with three types of graph: those featuring no arrow $A \not\rightarrow B$,

296 those featuring $A \rightarrow B$ but not $C \rightarrow B$ and $C \rightarrow A$ and the graph $C \xrightarrow{\quad} A \xrightarrow{\quad} B$. These
 297 possibilities yield $\mathcal{T}_{\mathcal{G}}^{\circ} = \mathcal{T}_{\mathcal{G}}^{\subseteq} \cup \{(\nu, y \mapsto \zeta) | \nu \in \mathcal{H} \setminus \mathcal{H}_{\mathcal{G}}\}$.

298 By 4, $|\eta(x; \{y\}) - \iota(x; \{y\})| < \epsilon$ for all $x \in A \cup \{*\}$ and $y \in B$ and therefore for $J \in \mathcal{J}$,
 299 $|U(\mu^{J\vee}(I_{(D)} \otimes \eta)) - U(\mu^{J\vee}(I_{(D)} \otimes \iota))| < \epsilon$. Therefore a small ϵ ensures $\mathcal{T}_{\mathcal{G}}^{\subseteq}$ yields a risk set
 300 “close” to the risk given by $\mathcal{T}_{\mathcal{G}}$ for any J . On the other hand, $|\iota(x; \{y\}) - \zeta(\{y\})|$ is independent of ϵ ,
 301 so $\mathcal{T}_{\mathcal{G}}^{\circ}$ yields a risk set that contains points that do not converge to the risk set induced by $\mathcal{T}_{\mathcal{G}}$ with
 302 small ϵ .

303 Extensions of the “base theory” $\mathcal{T}_{\mathcal{G}}$ can yield very different risk sets even when the departure from
 304 compatibility is slight and we limit those extensions to being based on CBNs. This example is
 305 complementary to results indicating that with unknown variable ordering (which may be regarded as
 306 analogous to $\mathcal{T}_{\mathcal{G}}^{\circ}$) or with unmeasured confounders it is not possible to construct a test that uniformly
 307 converges to the true graph equivalence class [Robins et al., 2003, Zhang and Spirtes, 2003]; our
 308 example shows that some misses may be benign and others may not. We will finally note that the
 309 more general theory $\mathcal{T}_{\mathcal{G}}^{\circ}$ still has a nontrivial risk set, and hence (potentially) nontrivial implications
 310 for decision making. We think that the investigation of risk sets for “extended theories” discussed here
 311 or graph learning algorithms considered in the CBN literature presents many interesting questions.

5 Potential Outcomes

Potential Outcomes is an alternative to the approach typified by Causal Bayesian Networks for formulating causal questions and hypotheses. Causal queries in the Potential Outcomes framework concern the distribution of random variables X_0, X_1 representing potential outcomes, or “the value X would have taken if action 0 or 1 were taken respectively” (Hernán and Robins [2018]). This is similar, but not the same, as the question answered by a consequence map which is “what is the distribution of X if I take actions 0 or 1?”

A natural connection between these informal notions of potential outcomes and consequence maps is given by the notion of consequence consistency. Let $\Delta(\mathcal{Y}_o)$ be the space of joint distributions over real and potential outcomes of X . A consequence map $\kappa : D \rightarrow \Delta(\mathcal{Y}_o)$ is consequence consistent if

$$(\delta_i \kappa)_{|X_i} F_X(w; A) = \delta_{X_i(w)}(A) \quad (7)$$

Consequence consistency is similar to the consistency condition [Richardson and Robins, 2013], but the latter does not involve consequences.

A causal theory that is consequence consistent need not have any particular relationship between an “observed” distribution $\mu \in \Delta(\mathcal{Y}_o)$ and an associated consequence κ ; one choice to make this connection is equality of the distributions of potential outcomes $\mu F_{X_i} = \delta_i \kappa F_{X_i}$, $i \in D$. Example D.1 in Appendix D shows that other choices may be preferred.

6 Equivalence of causal problems

Under what conditions could we consider a consequence consistent theory \mathcal{T}^{cc} associated with some distribution over potential outcomes to be “equivalent” to some causal theory \mathcal{T}^G associated with a CBN \mathcal{G} or vice versa?

The question of whether \mathcal{T}^G is consequence consistent with respect to some distribution over potential outcomes is easy to answer in the affirmative as consequence consistency is a trivial requirement if we choose potential outcomes $X_y := X$ for all $y \in D$.

The question of whether a consequence consistent theory \mathcal{T}^{cc} can in general be represented by a Causal Bayesian Network is then also straightforwardly answered in the negative, as conditions 2 and 3 of Definition 4.2 are in general non-trivial (condition 1 is trivial given a fully connected DAG \mathcal{G}).

The trivial potential outcome $X_y = X$ clashes with the informal idea that a potential outcome represents the value X would have taken had action y been taken - we might expect, for example, if $\delta_y \kappa F_X \neq \mu F_X$ then X would at least sometimes take a different value if the action y is taken than if it is not.

We might tentatively propose a more extensive set of assumptions to characterise a “Potential Outcomes” theory, which we will write \mathcal{T}^{po} .

Definition 6.1 (Potential Outcomes Causal Theory). A causal theory \mathcal{T}^{po} is a “Potential Outcomes” theory with respect to random variable $X : E \rightarrow X$ and potential outcome variable $X_i : E \rightarrow X$, $i \in D$ if for every $(\mu, \kappa) \in \mathcal{T}$, κ is consequence consistent (Eq. 7) and

$$\mu F_{X_i} = \delta_i \kappa F_{X_i} \quad (8)$$

If we consider only joint distributions over potential outcomes, a PO causal theory associates a unique consequence with each distribution

Note that the condition of consistency [Richardson and Robins, 2013], which is a very standard condition in the Potential Outcomes literature, is:

$$\mu_{|\{X_i, Z\}} F_X(w; A) = \delta_{X_i(w)}(A) \quad w \in Z^{-1}(i) \quad (9)$$

Where the random variable Z is a variable that is informally understood to be “intervenable” in a similar manner to intervention in Causal Bayesian Networks. A Potential Outcomes Causal Theory invokes a very general notion of Potential Outcomes where such intervenable variables may not exist, and so consistency may not be a sensible notion.

354 We can specify causal theories with a CBN \mathcal{G} that are not potential outcomes causal theories. Consider
 355 the graph X (with a single node and no edges). By condition 2 of Definition 4.2, the consequences in
 356 $\mathcal{T}^{\mathcal{G}}$ will all yield X distributed as a delta function for certain decisions. However, in general $\mathcal{T}^{\mathcal{G}}$ will
 357 contain distributions on the observation space E for which no variable is distributed according to a
 358 delta function. $\mathcal{T}^{\mathcal{G}}$ therefore cannot be a Potential Outcomes Causal Theory. We will outline below
 359 how a Potential Outcomes theory is not, in general, a theory associated with any CBN \mathcal{G} .

360 Rather than demand that we can represent the same theory with a CBN and with PO, we might ask
 361 instead if a problem featuring a PO theory can in general be reduced to a problem featuring a CBN
 362 theory and vice versa. This is in keeping with our approach that a CSDP represents at a high level a
 363 two person game and the latter determines the decision-relevant aspects of the problem.

364 **Definition 6.2** (Potential Outcomes CSDP). A CSDP $\langle (\mathcal{T}, (E, \mathcal{E}), X), (D, \mathcal{D}), (U, (F, \mathcal{F})) \rangle$ is a *Po-*
 365 *tential Outcomes CSDP* (POCSDP) if $E = F$, D is denumerable and there exists a set of potential
 366 outcome variables $X_i : E \rightarrow X$, $i \in D$ with respect to which \mathcal{T} is a Potential Outcomes causal
 367 theory.

368 **Definition 6.3** (CBN CSDP). A CSDP $\langle (\mathcal{T}, (E, \mathcal{E}), X), (D, \mathcal{D}), (U, (F, \mathcal{F})) \rangle$ is a *Causal Bayesian*
 369 *Network CSDP* (CBNCSDP) with respect to some finite DAG $\mathcal{G} = (V, W)$ if $E = F$ and \mathcal{T} is the
 370 theory induced by \mathcal{G}

371 Theorem 6.4 shows that, supposing D is denumerable, every CSDP can be reduced to a PO CSDP.
 372 For denumerable D , then, it suffices to show that conditions 1-3 of Definition 4.2 are nontrivial. Take
 373 some CSDP $\alpha = \langle (\mathcal{T}, E, X), D, (U, E) \rangle$ and suppose there is no $(\kappa, \mu) \in \mathcal{T}$, $y \in D$, $z \in X$ such that
 374 $\delta_y \kappa F_X(A) = \delta_z(A)$. Then it is straightforward to see that α cannot satisfy condition 3 of Definition
 375 4.2. Suppose that there is no $(\kappa, \mu) \in \mathcal{T}$, $y \in D$ such that $\delta_y \kappa = \mu$; it is then straightforward that
 376 conditions 1 and 2 of Definition 4.2 cannot be simultaneously satisfied.

(...and all the
other stuff
you need).

377 In both cases it is straightforward to posit generalised utilities such that α cannot be reduced.

378 Lifting condition 2 from the definition of a CBN yields CBNs with *generalized interventions*.

379 I strongly suspect this corresponds to the class of influence diagrams of [Dawid, 2010]

380 . Because conditions 1+2 are nontrivial, there exist POCSDPs that cannot be reduced to CSDPs based
 381 on CBNs with generalised interventions. Lifting conditions 2 and 3 yields a causal theory where we
 382 require only that the distributions given by every consequence κ are compatible with some DAG \mathcal{G} ,
 383 which we will call an *independence-only CBN*

384 I strongly suspect this is closely related to the notion of Extended Conditional Independence
of [Dawid, 2012]

385 . Condition 1 of Definition 4.2 can always be satisfied by choosing a graph \mathcal{G} that is fully connected,
 386 so lifting conditions 2 and 3 is sufficient to ensure that every POCSDP can be reduced to a CSDP
 387 featuring an independence-only CBN, and in fact an independence-only CBN can represent every PO
 388 causal theory.

389 The single world intervention graphs of Richardson and Robins [2013] are DAGs that rep-
resent independences among distributions over potential outcome variables. They might be
interpretable as POCSDPs.

390 The generalised versions of CBNs yield theories that generally associate multiple conse-
quences with each given distribution. However a generalized CBN still yields a unique causal
theory

391 **Theorem 6.4** (Reduction to PO). A CSDP $\alpha = \langle (\mathcal{T}, E, X), D, (U, E) \rangle$ where D is denumerable can
 392 be reduced to a PO CSDP.

393 *Proof.* Suppose $D = [M]$ or $D = \mathbb{N}^+$. Take $E' = E \times E^D$ and for $i \in D \cup \{0\}$, $x := (x_0, x_1, \dots) \in$
 394 E' define the projection $P_i(x_0, x_1, \dots) := x_i$ and the potential outcome variable $X_i := X \circ P_i$.

Take a map f from \mathcal{T} to causal states on E' such that, letting $(\kappa' F_X, \mu') := f(\kappa, \mu)$, for all $y \in D$ and $A_0, A_1, \dots \in \mathcal{E}$:

$$\mu^{po}(A_1 \times \dots) := \prod_{y' \in D} \delta_{y'} \kappa(A_{y'}) \quad (10)$$

$$\kappa'(y; A_0 \times A_1 \times \dots) := \int_{A_1 \times \dots} \delta_{x_y}(A_0) \mu^{po}(dx) \quad (11)$$

$$= \prod_{y' \in D \setminus \{y\}} \delta_{y'} \kappa(A_{y'}) \int_{A_y} \delta_{x_y}(A_0) \delta_y \kappa(dx_y) \quad (12)$$

$$\mu'(A_0 \times A_1 \times \dots) := \prod_{y' \in D \setminus \{y\}} \delta_{y'} \kappa(A_{y'}) \int_{A_y} \delta_{x_y}(A_0) \mu(dx_y) \quad (13)$$

It can be verified that κ' is a Markov kernel.

Note that by the definition of conditional probability, for $A, B \in \mathcal{X}$, $\int_{X_y^{-1}(A)} (\delta_y \kappa')|_{X_y} F_{X_0}(x; B) \delta_y \kappa'(dx) = \delta_y \kappa' \Upsilon(F_{X_0} \otimes F_{X_y})(A, B)$. Thus by 12, $\delta_x(A)$ is a version of $(\delta_y \kappa')|_{X_y} F_{X_0}(x; A)$, so κ' is consequence consistent.

Furthermore, $\mu' F_{X_y} = \delta_y \kappa F_X = \delta_y \kappa' F_{X_y}$ for $y \geq 1$. Therefore defining \mathcal{T}' to be the image of \mathcal{T} under f , we can see that \mathcal{T}' is a PO causal theory with respect to “observable” X_0 and “potential outcomes” $X_y, y \in D$.

For $A \in \mathcal{E}$:

$$\kappa' F_{P_0}(y; A) = \int_E \delta_z(A) \delta_y \kappa(dz) \quad (14)$$

$$= \int_A \kappa(y; dz) \quad (15)$$

$$= \kappa(y; A) \quad (16)$$

For all $B \in \mathcal{X}$

$$\mu' F_{X_0}(B) = \int_E \delta_z(X^{-1}(B)) \mu(dz) \quad (17)$$

$$= \mu F_X(B) \quad (18)$$

For all $J \in \mathcal{J}$ we have

$$U(\mu F_X J \Upsilon(I_{(D)} \otimes \kappa) = U(\mu' F_{X_0} J \Upsilon(I_{(D)} \otimes \kappa' F_{P_0})) \quad (19)$$

$$(20)$$

Therefore, given the PO CSDP $\beta = \langle (\mathcal{T}', E', X_0), D, (U, E) \rangle$, for all $J \in \mathcal{J}$, $R^\alpha(J, \kappa, \mu) = R^\beta(J, f(\kappa, \mu))$. Thus β is a reduction of α witnessed by f . \square

Corollary 6.5. *A CBN CSDP for which D is a denumerable set can be reduced to a PO CSDP.*

7 Conclusion

We have shown that CSDPs are an intuitive extension of SDPs and that causal theories that play a fundamental role in CSDPs can naturally represent models posed using the language of CBNs or PO. We believe that causal theories are quite general and capable of representing alternative approaches to causality such as IFMOCS [Peters et al., 2011] or approaches based on group invariance [Besserve et al., 2018].

This perspective raises many questions, for example: 1) Under what conditions do versions of the No-Free Lunch theorems hold for CSDPs? 2) Example 4.3 deals with a crude notion of “continuity” of a causal theory - whether a “nearby” distribution induces a similar risk set, which itself has implications for learnability of a causal theory. More generally, what properties may be used to characterise the

learnability of a causal theory? 3) The notation here borrows heavily from [Fong, 2013], whose diagrammatic representation of Markov kernels is closely related to the DAGs associated with CBNs. Can consequence maps be generically and informatively represented using diagrams similar to DAGs? 4) We have proposed consequence maps and causal theories as “relatively minimal” objects to satisfy the need to connect data, decisions and outcomes. Are there strictly more general objects that may be used instead, and if so under what assumptions are consequence maps and causal theories necessary? The general perspective proposed in this paper naturally incorporates the two major causal inference frameworks and, for the first time to our knowledge, allows a range of fundamental questions to be formally posed, such as *what are the characteristics of a causal statistical decision problem that make it “learnable”*? Whilst we don’t have all the answers, at least we have opened the way to ask such foundational questions!

References

- Yoshua Bengio, Tristan Deleu, Nasim Rahaman, Rosemary Ke, Sébastien Lachapelle, Olexa Bilaniuk, Anirudh Goyal, and Christopher Pal. A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms. *arXiv:1901.10912 [cs, stat]*, January 2019. URL <http://arxiv.org/abs/1901.10912>. arXiv: 1901.10912.
- Michel Besserve, Bernhard Schölkopf, Dominik Janzing, et al. Group invariance principles for causal generative models. In *International Conference on Artificial Intelligence and Statistics*, pages 557–565, 2018.
- Stephan Bongers, Jonas Peters, Bernhard Schölkopf, and Joris M. Mooij. Theoretical Aspects of Cyclic Structural Causal Models. *arXiv:1611.06221 [cs, stat]*, November 2016. URL <http://arxiv.org/abs/1611.06221>. arXiv: 1611.06221.
- Nancy Cartwright. *No Causes in, No Causes out*. Oxford University Press, April 1994. ISBN 978-0-19-159716-9. URL <https://www.oxfordscholarship.com/view/10.1093/0198235070.001.0001/acprof-9780198235071-chapter-3>.
- Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- A. Philip Dawid. Beware of the DAG! In *Causality: Objectives and Assessment*, pages 59–86, February 2010. URL <http://proceedings.mlr.press/v6/dawid10a.html>.
- 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>.
- Thomas S. Ferguson. *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press, July 1967. ISBN 978-1-4832-2123-6.
- 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.
- Sara Geneletti and A Philip Dawid. *Defining and identifying the effect of treatment on the treated*. Citeseer, 2007.
- MA Hernán and JM Robins. *Causal Inference*. Chapman & Hall/CRC, 2018.
- David Lewis. Causal decision theory. *Australasian Journal of Philosophy*, 59(1):5–30, March 1981. ISSN 0004-8402. doi: 10.1080/00048408112340011. URL <https://doi.org/10.1080/00048408112340011>.
- Mohammad Ali Mansournia, Julian P. T. Higgins, Jonathan A. C. Sterne, and Miguel A. Hernán. Biases in randomized trials: a conversation between trialists and epidemiologists. *Epidemiology (Cambridge, Mass.)*, 28(1):54–59, January 2017. ISSN 1044-3983. doi: 10.1097/EDE.0000000000000564. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5130591/>.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2 edition, 2009.

467 Jonas Peters, Joris M Mooij, Dominik Janzing, and Bernhard Schölkopf. Identifiability of causal
468 graphs using functional models. In *Proceedings of the Twenty-Seventh Conference on Uncertainty*
469 *in Artificial Intelligence*, pages 589–598. AUAI Press, 2011.

470 Thomas S Richardson and James M Robins. Single world intervention graphs (swigs): A unification
471 of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social*
472 *Sciences, University of Washington Series. Working Paper*, 128(30):2013, 2013.

473 James M Robins and Thomas S Richardson. Alternative graphical causal models and the identification
474 of direct effects. *Causality and psychopathology: Finding the determinants of disorders and their*
475 *cures*, pages 103–158, 2010.

476 James M. Robins, Richard Scheines, Peter Spirtes, and Larry Wasserman. Uniform consistency in
477 causal inference. *Biometrika*, 90(3):491–515, September 2003. ISSN 0006-3444. doi: 10.1093/
478 biomet/90.3.491. URL <https://academic.oup.com/biomet/article/90/3/491/231406>.

479 Donald B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies.
480 *Journal of Educational Psychology*, 66(5):688–701, 1974. ISSN 1939-2176(Electronic),0022-
481 0663(Print). doi: 10.1037/h0037350.

482 Donald B. Rubin. Causal Inference Using Potential Outcomes. *Journal of the American Statistical As-*
483 *sociation*, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880.
484 URL <https://doi.org/10.1198/016214504000001880>.

485 Peter Spirtes, Clark N Glymour, Richard Scheines, David Heckerman, Christopher Meek, Gregory
486 Cooper, and Thomas Richardson. *Causation, prediction, and search*. MIT press, 2000.

487 Jin Tian and Judea Pearl. A general identification condition for causal effects. In *Aaai/iaai*, pages
488 567–573, July 2002.

489 Abraham Wald. *Statistical decision functions*. Statistical decision functions. Wiley, Oxford, England,
490 1950.

491 Jiji Zhang and Peter Spirtes. Strong Faithfulness and Uniform Consistency in Causal Inference.
492 In *Proceedings of the Nineteenth Conference on Uncertainty in Artificial Intelligence*, UAI’03,
493 pages 632–639, San Francisco, CA, USA, 2003. Morgan Kaufmann Publishers Inc. ISBN 978-0-
494 12-705664-7. URL <http://dl.acm.org/citation.cfm?id=2100584.2100661>. event-place:
495 Acapulco, Mexico.

Supplement to: Causal Statistical Decision Problems

A Markov Kernels

This is an expanded version of Section 2 that explains some notation more thoroughly.

A measurable space (E, \mathcal{E}) is a set E and a σ -algebra $\mathcal{E} \subset \mathcal{P}(\mathcal{E})$ containing the measurable sets. A probability measure $\mu \in \Delta(\mathcal{E})$ is a nonnegative map $\mathcal{E} \rightarrow [0, 1]$ such that $\mu(\emptyset) = 0$, $\mu(E) = 1$ and for countable $\{E_i\} \in \mathcal{E}$, $\mu(\cup_i E_i) = \sum_i \mu(E_i)$.

We assume all measurable spaces discussed are standard. That is, they are isomorphic to either a subset of \mathbb{N} with the discrete σ -algebra, or \mathbb{R} with the Borel σ -algebra.

Given two measurable sets (E, \mathcal{E}) and (F, \mathcal{F}) , a *Markov kernel* K is a map $E \times \mathcal{F} \rightarrow [0, 1]$ where

1. The map $x \mapsto K(x; B)$ is \mathcal{E} -measurable for every $B \in \mathcal{F}$
2. The map $B \mapsto K(x; B)$ is a probability measure on (F, \mathcal{F}) for every $x \in E$

Abusing notation somewhat, we will give Markov kernels the alternate type signature $K : E \rightarrow \Delta(\mathcal{F})$, noting that due to part 1 not every map with this type is a Markov kernel. We will sometimes write the set of Markov kernels of type $E \rightarrow \Delta(\mathcal{F})$ as $\Delta(\mathcal{F})^D$, noting again that given part 1, the set of Markov kernels of this type may be smaller than $\Delta(\mathcal{F})^D$.

If we have two random variables $X : _ \rightarrow X$ and $Y : _ \rightarrow Y$, the conditional probability $P(Y|X)$ is a Markov kernel $X \rightarrow \Delta(\mathcal{Y})$. Formally, given $\mu \in \Delta(\mathcal{E})$ and a sub- σ -algebra $\mathcal{E}' \subset \mathcal{E}$, there is a Markov kernel $\mu|_{\mathcal{E}'} : E \rightarrow \Delta(\mathcal{E})$ such that for $A \in \mathcal{E}$ and $B \in \mathcal{E}'$, $\int_B \mu|_{\mathcal{E}'}(y; A) d\mu(y) = \mu(A \cap B)$. $\mu|_{\mathcal{E}'}$ is a *conditional probability distribution* with respect to \mathcal{E}' . This result may not hold if (E, \mathcal{E}) is not a standard measurable space [Çinlar, 2011].

Given a set of random variables $\mathbf{X} = \{X^i\}_{i \in [N]}$ with domain (E, \mathcal{E}) , $\mu|_{\mathbf{X}} : E \rightarrow \Delta(\mathcal{E})$ is a conditional probability distribution with respect to the σ -algebra generated by \mathbf{X} : $\sigma(\cup_{i \in [N]} \sigma(\mathcal{X}^i))$. We will use this subscript notation rather than the more common bar notation (e.g. $\mu(\cdot | \mathbf{X})$) to express conditional probability from here onwards.

Two Markov kernels $K : E \rightarrow \Delta(\mathcal{F})$ and $K' : E \rightarrow \Delta(\mathcal{F})$ are μ -almost surely equivalent given $\mu \in \Delta(\mathcal{E})$ if

$$\int_A K(x; B) d\mu = \int_A K'(x; B) d\mu \quad \forall A \in \mathcal{E}, B \in \mathcal{F} \quad (21)$$

A.1 Operations with Markov kernels

For the following, assume K is a Markov kernel from $E \rightarrow \Delta(\mathcal{F})$, K' a kernel $E \rightarrow \Delta(\mathcal{H})$, L is a Markov kernel $F \rightarrow \Delta(\mathcal{G})$, μ is a probability measure on (E, \mathcal{E}) , ν is a probability measure on (F, \mathcal{F}) and f is a nonnegative measurable function $F \rightarrow \mathbb{R}$.

The notation here borrows heavily from Çinlar [2011] and Fong [2013].

A.1.1 Kernel products

The kernel-kernel product KL is a Markov kernel $E \rightarrow \Delta(\mathcal{G})$ such that $KL(x; B) := \int_F K(x; dy) L(y; B)$, $x \in E, B \in \mathcal{G}$.

The measure-kernel product of μ and K , μK is a probability measure on (F, \mathcal{F}) such that $\mu K(B) = \int_E \mu(dx) K(x; B)$, $B \in \mathcal{F}$.

The kernel-function product Kf is a nonnegative measurable function $E \rightarrow \mathbb{R}$ such that $Kf(x) := \int_F K(x; dy) f(y)$, $x \in E$.

Kernel products are in general associative: $(KL)M = K(LM)$.

A.1.2 Special kernels

$I_{(E)}$ is a kernel $E \rightarrow \Delta(\mathcal{E})$ defined by $x \mapsto \delta_x$. It has the properties $\mu I_{(E)} = \mu$, $K I_{(F)} = K$, $I_{(E)} K = K$, $I_{(F)} f = f$.

538 Υ_E is a kernel $E \rightarrow \Delta(\mathcal{E} \otimes \mathcal{E})$ defined by $x \mapsto \delta_{(x,x)}$. We will subsequently leave the space implicit.
 539 The symbol Υ is pronounced “splitter”.

540 Given $M : H \rightarrow \Delta(\mathcal{I})$, $K \otimes M$ is a Markov kernel $E \times H \rightarrow \Delta(\mathcal{F} \otimes \mathcal{I})$ where

$$K \otimes M(x, y; A \times B) := K(x; A)M(y; B) \quad (22)$$

541 Given $N : I \rightarrow \Delta(\mathcal{J})$, it can be verified that $(K \otimes M)(L \otimes N) = KL \otimes MN$.

542 $\Upsilon(K \otimes K')$ is a Markov kernel $E \rightarrow \Delta(\mathcal{F} \otimes \mathcal{H})$ and

$$\Upsilon(K \otimes K')(x; A \times B) = \int_E K(x'; A)K'(x''; B)\delta_{(x,x)}(dx' \times dx'') \quad (23)$$

$$= K(x; A)K'(x; B) \quad (24)$$

543 We can overload notation to use $\Upsilon(K \otimes K' \otimes K'')$ for the nested construction $\Upsilon(K \otimes \Upsilon(K' \otimes K''))$.

544 Let $(*, \{\emptyset, *\})$ be an indiscrete measurable set. \uparrow_E is a kernel $E \rightarrow \Delta(\{\emptyset, *\})$ defined by $x \mapsto \mathbb{1}_*$.

545 We have $\Upsilon(I \otimes \uparrow) = I$. The symbol \uparrow is pronounced “stopper”.

546 Given some measurable function $g : E \rightarrow F$, the kernel $F_g : E \rightarrow \Delta(\mathcal{F})$ is defined by $x \mapsto \delta_{g(x)}$. It
 547 is easy to check that $F_g F_g = F_g$. For $\mu \in \Delta(\mathcal{E})$, the product μF_g is the push forward measure $g_*\mu$.

$$\mu F_g(A) = \int_E \delta_{g(x)}(A) d\mu \quad (25)$$

$$= \mu(g^{-1}(A)) \quad (26)$$

$$= g_*\mu(A) \quad (27)$$

548 Given two random variables $X : (E, \mathcal{E}) \rightarrow (X, \mathcal{X})$ and $Y : (E, \mathcal{E}) \rightarrow (Y, \mathcal{Y})$, the product $\mu^\Upsilon(F_X \otimes$
 549 $F_Y)$ is the joint distribution of X and Y .

$$\mu^\Upsilon(F_X \otimes F_Y)(A, B) = \int_E \delta_{X(x)}(A) \delta_{Y(x)}(B) d\mu \quad (28)$$

$$= \mu(X^{-1}(A) \cap Y^{-1}(B)) \quad (29)$$

550 B Appendix: Causal Statistical Decision Problems

551 **Lemma B.1** (Reduction preserves admissibility). *If a CSDP β with induced game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be*
 552 *reduced to a statistical decision problem α with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$ then a decision function*
 553 *$J \in \mathcal{J}$ is admissible in β iff it is admissible in α .*

554 *Proof.* Suppose $J \in \mathcal{J}$ is inadmissible in α . Then there is some $J' \in \mathcal{J}, \mu \in \mathcal{H}$ such that $R'(J', \mu) <$
 555 $R'(J, \mu)$ and $R'(J', \nu) \leq R'(J, \nu)$ for all $\nu \in \mathcal{H}$. Let h be the function that witnesses the reduction.
 556 Then we have for all $\tau \in h^{-1}(\mu)$, $R(J', \tau) = R'(J', \mu) < R(J, \tau) = R'(J, \nu)$ and for all $\nu \in \mathcal{H}$,
 557 $\chi \in h^{-1}(\nu)$, $R(J', \chi) = R'(J', \nu) \leq R(J, \chi) = R'(J, \nu)$. The set $\bigcup_{\nu \in \mathcal{H}} h^{-1}(\nu) = \mathcal{T}$, so J is
 558 inadmissible in β .

559 Suppose $J \in \mathcal{J}$ is admissible in β . Then there is some $J' \in \mathcal{J}, \tau \in \mathcal{T}$ such that $R(J', \tau) < R(J, \tau)$ and
 560 $R(J', \chi) \leq R(J, \chi)$ for all $\chi \in \mathcal{T}$. Then we have $R'(J', h(\tau)) = R(J', \tau) < R(J, \tau) = R'(J, h(\tau))$
 561 and $R'(J', h(\chi)) = R(J', \chi) \leq R(J, \chi) = R'(J, h(\chi))$. Because h is surjective, J is admissible in
 562 α . \square

563 **Corollary B.2** (Reduction preserves completeness). *If a causal decision problem β with induced*
 564 *game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be reduced to a statistical decision problem α with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$, then*
 565 *an (essentially) complete class with respect to α is (essentially) complete with respect to β .*

566 **Lemma B.3** (Induced Bayes rule). *If a CSDP β with induced game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be reduced to a*
 567 *statistical decision problem α with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$ witnessed by $h : \mathcal{T} \rightarrow \mathcal{H}$ and $J_{ba}^\xi \in \mathcal{J}$ is*
 568 *a Bayes rule with respect to the problem α and the prior ξ then J_{ba}^ξ is a Bayes rule with respect to the*
 569 *problem β and the induced prior ξ_h .*

570 *Proof.* For any $J \in \mathcal{J}, \tau \in \mathcal{T}$, by the properties of the push-forward measure

$$\int_{\mathcal{T}} R(J, \tau) d\xi_h = \int_{\mathcal{H}} R'(J, h(\tau)) d\xi \quad (30)$$

571 And therefore, if a Bayes rule exists,

$$\arg \min_{J \in \mathcal{J}} \int_{\mathcal{T}} R(J, \tau) d\xi_h = \arg \min_{J \in \mathcal{J}} \int_{\mathcal{H}} R'(J, h(\tau)) d\xi \quad (31)$$

572 \square

573 **Theorem B.4** (Complete class theorem (CSDP)). *Given an CSDP $\alpha := \langle \langle \mathcal{T}, E \rangle, D, \mathbf{X}, U \rangle$ with risk*
 574 *R , if there is a reduction to an SDP $\beta := \langle \langle \mathcal{H}, F \rangle, D, \mathbf{Y}, \ell \rangle$ with risk R' such that $|\mathcal{H}| < \infty$ and*
 575 *$\inf_{J \in \mathcal{J}, \mu \in \mathcal{H}} R'(J, \mu) < -\infty$ then the set of all Bayes decision functions is a complete class and the*
 576 *set of all admissible Bayes decision functions is a minimal complete class.*

577 *Proof.* Given the conditions, the Bayes decision functions in β form a complete class and admissible
 578 Bayes rules a minimal complete class [Ferguson, 1967].

579 By Corollary B.2 the Bayes rules for β are complete in α , and the admissible Bayes rules for β are
 580 essentially complete in α .

581 Every (admissible) Bayes rule for β is a(n admissible) Bayes rule for α , so the set of (admissible)
 582 Bayes rules for α is also (essentially) complete in α . \square

583 **Theorem B.5** (Reduction of a CSDP on observations). *A CSDP $\alpha =$*
 584 *$\langle \langle \mathcal{T}^\alpha, (E, \mathcal{E}), \mathbf{X} \rangle, D, (U, (F, \mathcal{F})) \rangle$ where, for $\zeta \in \Delta(\mathcal{E} \otimes \mathcal{D})$ can be reduced to a problem*
 585 *$\beta = \langle \langle \mathcal{T}^\beta, (X, \mathcal{X}), \text{id}_X \rangle, D, (U, (F, \mathcal{F})) \rangle$ by marginalization.*

586 *Proof.* Consider the mapping $g : \mathcal{T}^\alpha \rightarrow \mathcal{T}^\beta$ given by $(\kappa, \mu) \mapsto (\kappa, \mu F_X)$.

587 For $J \in \mathcal{J}$, $(\kappa, \mu) \in \mathcal{T}^\alpha$

$$R^\alpha(J, \kappa, \mu) = \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \Upsilon(I_{(D)} \otimes \kappa)) - U(\mu F_X J \Upsilon(I_{(D)} \otimes \kappa)) \quad (32)$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \Upsilon(I_{(D)} \otimes \kappa)) - U(\mu F_X F_X J \Upsilon(I_{(D)} \otimes \kappa)) \quad (33)$$

$$= R^\beta(J, g(\kappa, \mu)) \quad (34)$$

588 \square

589 **Theorem B.6** (Reduction of a CSDP on the utility). *Given a CSDP $\alpha = \langle (\mathcal{T}^\alpha, (E, \mathcal{E}), \mathbf{X}), D, (U, (F, \mathcal{F})) \rangle$ where, for $\zeta \in \Delta(\mathcal{E} \otimes \mathcal{D})$, if $U(\zeta) = U'(\zeta(I_{(D)} \otimes F_Y))$*
 590 *for some $Y : F \rightarrow Y$ and $U' : \Delta(\mathcal{Y}) \rightarrow \mathbb{R}$ then α has Y -observable utility. Such a problem can be*
 591 *reduced to a problem $\beta = \langle (\mathcal{T}^\beta, (E, \mathcal{E}), \mathbf{X}), D, (U', (Y, \mathcal{Y})) \rangle$ by marginalization.*
 592

593 *Proof.* Consider the mapping $g : \mathcal{T}^\alpha \rightarrow \mathcal{T}^\beta$ given by $(\kappa, \mu) \mapsto (\kappa F_Y, \mu)$.

594 We have for $J \in \mathcal{J}$, $(\kappa, \mu) \in \mathcal{T}^\alpha$

$$R^\alpha(J, \kappa, \mu) = \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \Upsilon(I_{(D)} \otimes \kappa)) - U(\mu F_X J \Upsilon(I_{(D)} \otimes \kappa)) \quad (35)$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U'(\gamma' \Upsilon(I_{(D)} \otimes \kappa)(I_{(D)} \otimes F_Y)) - U'(\mu F_X J \Upsilon(I_{(D)} \otimes \kappa)(I_{(D)} \otimes F_Y)) \quad (36)$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U'(\gamma' \Upsilon(I_{(D)} \otimes \kappa F_Y)) - U'(\mu F_X J \Upsilon(I_{(D)} \otimes \kappa F_Y)) \quad (37)$$

$$= R^\beta(J, g(\kappa, \mu)) \quad (38)$$

595 \square

596 **Theorem B.7.** *Every SDP $\langle (\mathcal{H}, E, \mathbf{X}), D, \ell \rangle$ can be reduced to a CSDP.*

597 *Proof.* Take D to be the projection from $D \times E$ to D . For each $\mu \in \mathcal{H}$ define the consequence
 598 $\kappa_\mu : d \mapsto \mu$ for all $d \in D$. Take the causal theory $\mathcal{T} = \{(\kappa_\mu, \mu) | \mu \in \mathcal{H}\}$ for some $\pi \in \Delta(\mathcal{D})$ and
 599 the pseudo-utility $U(\nu) = -\mathbb{E}_\nu[\ell(P_E^\nu, D)]$ to construct the CSDP $\langle (\mathcal{T}, E, \mathbf{X}), D, (U, E) \rangle$. We will
 600 show that the original problem can be reduced to this.

601 For $\gamma \in \Delta(\mathcal{D})$ the induced loss L is

$$L(\kappa_\mu, \gamma) = - \sup_{\gamma' \in \Delta(\mathcal{D})} \mathbb{E}_{\gamma' \Upsilon(I_{(D)} \otimes \kappa_\mu) | E} [\ell(\gamma' \Upsilon(I_{(D)} \otimes \kappa_\mu) | E, D)] + \mathbb{E}_{\gamma \Upsilon(I_{(D)} \otimes \kappa_\mu)} [\ell(\gamma \Upsilon(I_{(D)} \otimes \kappa_\mu) | E, D)] \quad (39)$$

$$= \mathbb{E}_\gamma[\ell(\mu, D)] \quad (40)$$

602 For the surjective map, take $g : \mathcal{H} \rightarrow \mathcal{T}$ defined by $g(\mu) = \kappa_\mu$.

603 Denote by R the risk associated with the SDP $\langle (\mathcal{H}, E), D, \mathbf{X}, \ell \rangle$ and by R' the risk associated with
 604 the CSDP $\langle (\mathcal{T}, E), D, \mathbf{X}, U \rangle$. Then

$$R'(J, \kappa, \mu) = \int_D \ell(\mu, y) \mu F_X J(dy) \quad (41)$$

$$= R(J, g(\kappa, \mu)) \quad (42)$$

605 \square

606 **Theorem B.8.** *Given a CSDP $\beta = \langle (\mathcal{T}, E, \mathbf{X}), D, (U, F) \rangle$ where U is an ordinary pseudo-utility, let*
 607 $\mathcal{K} = \{\kappa | (\kappa, \mu) \in \mathcal{T}\}$ *be the set of consequences. β is reducible to a statistical decision problem on*
 608 *the measurable space $(E \times F \times D, \mathcal{E} \otimes \mathcal{F} \otimes \mathcal{D})$ if there is some surjective map $m : \Delta(\mathcal{F} \otimes \mathcal{D}) \rightarrow \mathcal{K}$.*

609 *Proof.* Let $\mathcal{H} \subset \Delta(\mathcal{E} \otimes \mathcal{F} \otimes \mathcal{D})$ be some hypothesis class and let m^\dagger be a right inverse of m . Define
 610 $h : \mathcal{T} \rightarrow \mathcal{H}$ by $(\kappa, \mu) \mapsto \mu \otimes m^\dagger(\kappa)$.

611 Let $k : \Delta(\mathcal{F})^D \times D \rightarrow \mathbb{R}$ be the differential loss induced by the ordinary pseudo-utility U (see
 612 Equation 3).

613 Given the projections $F : E \times F \times D \rightarrow F$ and $D : E \times F \times D \rightarrow D$ and arbitrary $\xi \in \Delta(\mathcal{E} \otimes \mathcal{F} \otimes \mathcal{D})$
 614 define $\ell : \mathcal{H} \times D \rightarrow [0, \infty)$ by

$$\ell(\xi, y) = k(m(\xi F_{\bigcup_{(F \otimes D)} }), y) \quad (43)$$

615 Note that

$$\ell(h(\kappa, \mu), y) = k(\kappa, y) \quad (44)$$

616 Define $X' : E \times F \times D \rightarrow X$ by $(a, b, c) \mapsto X(a)$.

617 Then, given the statistical decision problem $\langle (\mathcal{H}, E \times F \times D, X'), D, \ell \rangle$, we have for all $J \in \mathcal{J}$,
 618 $(\kappa, \mu) \in \mathcal{T}$ the risk

$$R'(J, h(\kappa, \mu)) = \int_D \ell(h(\kappa, \mu), y) h(\kappa, \mu) F_{X'} J(dy) \quad (45)$$

$$= \int_D \ell(h(\kappa, \mu), y) (\mu \otimes m^\dagger(\kappa)) F_{X'} J(dy) \quad (46)$$

$$= \int_D k(\kappa, y) \mu F_X J(dy) \quad (47)$$

$$= R(J, \kappa, \mu) \quad (48)$$

619

□

620 **Example B.9** (Irreducible CSDP). The choice of decision function in an SDP does not affect the
 621 state, while this choice does affect the outcome in an CSDP. For an SDP, then, the risk of a mixed
 622 decision function is equal to the mixture of risks of each atomic decision function but this is not true
 623 in general for an CSDP.

624 Take the CSDP $\langle (\mathcal{T}, E), D, X, U \rangle$ where $E = D = \{0, 1\}$, $Y : E \rightarrow \{0, 1\}$ is the identity function,
 625 $U : \mu \mapsto -\text{Var}_\mu[Y]$ and $\mathcal{T} = \{(d \mapsto \delta_d, \nu) | \nu \in \Delta(\mathcal{E})\}$.

626 For any $(\kappa, \mu) \in \mathcal{T}$ and $J \in \mathcal{J}$ we have

$$R(J, \kappa, \mu) = 0.25 - \text{Var}_{\mu F_X J}(Y) \quad (49)$$

627 Consider the forgetful decision functions $J_0 : x \mapsto \text{Bernoulli}(0)$ and $J_{1/2} : x \mapsto \text{Bernoulli}(\frac{1}{2})$ and
 628 $J_1 : x \mapsto \text{Bernoulli}(1)$ for all $x \in X$. Note that $J_{1/2}(x; A) = \frac{1}{2}(J_0(x; A) + J_1(x; A))$ for all
 629 $x \in X, A \in \mathcal{D}$. For any statistical decision problem with risk R' ,

$$R'(J_{1/2}, \mu) = \int_D \ell(\mu, y) \mu F_X J_{1/2}(dy) \quad (50)$$

$$= \frac{1}{2} \left(\int_D \ell(\mu, y) \mu F_X J_0(dy) + \int_D \ell(\mu, y) \mu F_X J_1(dy) \right) = \frac{1}{2} (R'(J_0, \mu) + R'(J_1, \mu)) \quad (51)$$

630 But

$$R(J_{1/2}, \kappa, \mu) = 0 \quad (52)$$

$$\neq \frac{1}{2} (R(J_0, \kappa, \mu) + R(J_1, \kappa, \mu)) \quad (53)$$

631 **Corollary B.10.** The class of nonrandomized decision functions is not essentially complete for
 632 CSDPs. The stochastic decision function $J_{1/2}$ is strictly better than any deterministic function in the
 633 above example.

C Appendix: CBN is a causal theory

Theorem C.1. Given a measurable set (E, \mathcal{E}) and a graph \mathcal{G} over a set of random variables $\{X^i\}_{i \in [N]}$ where $X^i : E \rightarrow X^i$, a decision set (D, \mathcal{D}) and random variables $\{D^i\}_{i \in [N]}$ with $D^i : (D, \mathcal{D}) \rightarrow (X^i \cup \{*\}, \sigma(X^i \cup \{*\}))$. Given $\mu \in \mathcal{G}_{\mathcal{G}}$, let μ^y be the \mathcal{G}, μ, y -interventional distribution (Definition 4.2).

Then the map $\kappa^{\mu, \mathcal{G}} : D \rightarrow \Delta(\mathcal{E})$ given by $y \mapsto \mu^y$ is a Markov kernel with respect to (D, \mathcal{D}) and (E, \mathcal{E}) .

Proof. The DAG \mathcal{G} induces a partial ordering on the RV's X^i by $X^i < X^j$ if $X^i \rightarrow X^j$ is in \mathcal{G} . Without loss of generality, suppose the total ordering X^0, \dots, X^N is consistent with the partial ordering induced by \mathcal{G} .

Let $\kappa^i : \mathcal{E} \rightarrow \Delta(\mathcal{X}^i)$ be defined by $\kappa^i(x; A) := \mu_{|X^{<i}} F_{X^i}$. Note that by the compatibility of μ , for all $x \in \mathcal{E}$, $A \in \mathcal{X}^i$ we also have

$$\kappa^i(x; A) = \mu_{|Pa_{\mathcal{G}}(X^i)} F_{X^i}(x; A) \quad (54)$$

Consider $\kappa^{i,*} : D \times E \rightarrow \Delta(\mathcal{X}^i)$ given by

$$\kappa^{i,*}(y, pa^i; A) := \begin{cases} \kappa^i(pa^i; A) & D^i(y) = * \\ \delta_{D^i(y)}(A) & D^i(y) \neq * \end{cases} \quad (55)$$

Clearly for every $(d, pa^i) \in D \times E$ the map $A \mapsto \kappa^{i,*}(d, pa^i; A)$ is a probability distribution on \mathcal{X}^i . Fix $B \in \mathcal{X}_i$ and let $\kappa_B^{i,*} = \kappa_i^{i,*}(\cdot; B)$.

Then for any $A \in \mathcal{B}([0, 1])$

$$[\kappa_B^{i,*}]^{-1}(A) = [D^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \quad \text{if } 0, 1 \notin A \quad (56)$$

$$= [D^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [D^i]^{-1}(B) \times X^{Pa_{\mathcal{G}}(i)} \quad \text{if } 1 \in A \wedge 0 \notin A \quad (57)$$

$$= [D^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [D^i]^{-1}(B^C) \times X^{Pa_{\mathcal{G}}(i)} \quad \text{if } 0 \in A \wedge 1 \notin A \quad (58)$$

$$= [D^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [D^i]^{-1}(X^i) \times X^{Pa_{\mathcal{G}}(i)} \quad \text{if } 0 \in A \wedge 1 \in A \quad (59)$$

Note that $\sigma(Pa_{\mathcal{G}}(X^i)) \subset \mathcal{E}$ and $[\kappa_i^B]^{-1}(A) \in \sigma(Pa_{\mathcal{G}}(X^i))$. Further note that $\{*\}$, B and B^C are in $\sigma(X^i \cup \{*\})$. Therefore, in every case the result is an element of $\mathcal{E} \otimes \mathcal{D}$ and $\kappa^{i,*}$ is a Markov kernel.

Then $\iota^{\mathcal{G}} : D \rightarrow \Delta(\mathcal{X})$ defined below is a Markov kernel.

$$\iota^{\mathcal{G}} : (y; A) \mapsto \int_{A^0} \kappa^{0,*}(y; dx^0) \dots \int_{A^{N-1}} \kappa^{N-1,*}(y, x^{n-2}; dx^{n-1}) \kappa^{N,*}(y, x^{n-1}; A^N) \quad (60)$$

for $y \in D$, $A \in E$ and $A^i = [X^i]^{-1}(A)$.

From Equations 54, 55 and 60 we can verify that, given some $i \in N$, if $D^i(y) = \{*\}$ then $[\delta_y \iota^{\mathcal{G}}]_{Pa_{\mathcal{G}}(X^i)} = \kappa_i = \mu_{Pa_{\mathcal{G}}(X^i)} F_{X^i}$ and if $D^i(y) \neq \{*\}$ then $\delta_y \iota^{\mathcal{G}} = \delta_{D^i(y)} F_{X^i}$. From Equation 60 and the compatibility of μ with \mathcal{G} it further follows that $\delta_y \iota^{\mathcal{G}}$ is compatible with \mathcal{G} . Therefore $\delta_y \iota^{\mathcal{G}} = \mu^y$ and so $\iota^{\mathcal{G}} = \kappa^{\mu, \mathcal{G}}$. \square

D Appendix: Counterfactuals

A causal theory for Potential Outcomes is associated with a much larger hypothesis class than any causal theory that works only with distributions over observable variables. Theorems B.5 and B.6 show that given any SCDP based on Potential Outcomes, provided that the potential outcome variables are unobserved and the utility does not depend on them, a reduced SCDP can be constructed by marginalising over potential outcomes. Potential outcomes are not universally excluded by this; there are some examples of problems where one does care about the values of potential outcome variables. The *effect of treatment on the treated* (ETT) that depends on counterfactual quantities and has some relevance to decision preferences Rubin [1974], though it is controversial whether this dependence is necessary Geneletti and Dawid [2007]. More straightforwardly, the legal standard of “no harm but for the defendant’s negligence” does seem to invoke fundamentally counterfactual considerations Pearl [2009].

Example D.1 (Performance bias). Suppose we have a CSDP $\langle (\mathcal{T}, E), D, X, (U, E) \rangle$ where the observed data X is from a randomised controlled trial (RCT), $Y_0 : E \rightarrow Y$ and $Y_1 : E \rightarrow Y$ are random variables representing a particular outcome of interest under no treatment and treatment respectively and $Y : E \rightarrow Y$ represents the “realised” outcome of interest and for $\xi \in \Delta(\mathcal{E})$, $U(\xi) = \mathbb{E}_\xi[Y]$.

Under usual assumptions about RCTs, if we suppose the observed data are distributed according to $\mu \in \Delta(\mathcal{E})$ it is possible (given infinite data X) to determine $\mathbb{E}_\mu[Y_0]$ and $\mathbb{E}_\mu[Y_1]$ [Rubin, 2005].

Consequence consistency is assumed, but performance bias is suspected, which can lead to $\delta_i \kappa Y_i$ differing from $\mathbb{E}_\mu[Y_i]$ [Mansournia et al., 2017].

1. Assume performance bias is absent, so the theory must satisfy $\delta_i \kappa Y_i = \mathbb{E}_\mu[Y_i]$
2. Assume performance bias has a uniform additive effect: the theory satisfies $\delta_i \kappa Y_i = \mathbb{E}_\mu[Y_i] + k$. In this case the average treatment effect can still be estimated from the data: $\delta_1 \kappa Y_1 - \delta_0 \kappa Y_0 = \mathbb{E}[Y_1] - \mathbb{E}[Y_0]$ which may be sufficient to find a decision function minimising the risk
3. Avoid assumptions about the effect of performance bias; the theory satisfies no particular relationship between $\mathbb{E}_\mu[Y_i]$ and $\delta_i \kappa Y_i$ and we may therefore expect preferred decision function to ignore the data

The question of specifying this relationship arises naturally when we consider connecting Potential Outcomes to CSDPs. Nonetheless, the possibility of deviations from option 1 above are often treated as “external to the causal problem”. For example, Mansournia et al. [2017] states:

In this case, it might be more appropriate to say that the intention-to-treat effect from the trial is not generalizable or transportable to other settings rather than saying that it is “biased”