# Causal Statistical Decision Problems

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

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

## 1 Introduction

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

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

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

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

Queries in the CBN framework may be concerned with identification of causal effects given a graph and a probability distribution [Tian and Pearl, 2002], or with the determination of the true causal 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 $\{1, ..., N\}$. Sets are ordinary capital letters $X$ while $\sigma$-algebras are calligraphic capitals $\mathcal{X}$ and random variables are sans serif capitals $\mathsf{X} : \_ \to X$. The calligraphic $\mathcal{G}$ refers to a directed acyclic graph rather than a $\sigma$-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 $\sigma$-algebra $\mathcal{E} \subset \mathcal{P}(\mathcal{E})$ containing the measurable sets. A probability measure $\mu \in \Delta(\mathcal{E})$ is a nonnegative map $\mathcal{E} \to [0, 1]$ 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 $\sigma$-algebra, or $\mathbb{R}$ with the Borel $\sigma$-algebra.

Given two measureable spaces $(E, \mathcal{E})$ and $(F, \mathcal{F})$, a *Markov kernel* or *stochastic map* $K : E \to \Delta(\mathcal{F})$ 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 $E \to \Delta(\mathcal{F})$ as $\Delta(\mathcal{F})^D$.

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

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 measureable space [Çinlar, 2011].

93 Given a set of random variables $\mathbf{X} = \{\mathsf{X}^i\}_{i \in [N]}$ with domain $(E, \mathcal{E})$, $\mu_{\mathbf{X}} : E \to \Delta(\mathcal{E})$ is a
94 conditional probability distribution with respect to the $\sigma$-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 $K : E \to \Delta(\mathcal{F})$ and $K' : E \to \Delta(\mathcal{F})$ are $\mu$-almost surely equivalent given
98 $\mu \in \Delta(\mathcal{E})$ if for all $A \in \mathcal{E}, B \in \mathcal{F}, \int_A K(x; B) d\mu = \int_A K'(x; B)$.

**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 $K : E \to \Delta(\mathcal{F}), L : F \to \Delta(\mathcal{G})$, and $M : G \to \Delta(\mathcal{H})$ are Markov
102 kernels, $\mu$ is a probability measure on $(E, \mathcal{E})$.

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

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

**Special kernels:** $I_{(E)}$ is the identity kernel $E \to \Delta(\mathcal{E})$ 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 $g : E \to F$, the kernel $F_g : E \to \Delta(\mathcal{F})$ 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 $\mu F_{\mathsf{X}}$ and a marginal kernel $\kappa F_{\mathsf{X}}$.

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

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

114 The symbol $\curlyvee$ is read "splitter".

# 3 Causal Statistical Decision Problems

|  | SDPs | CSDPs |
|---|---|---|
| State of the world | $\Theta$ | $\mathcal{T}$, causal theory |
| Observation space | $E$ | $E$ |
| Result space | - | $F$ |
| Decisions | $D$ | $D$ |
| Given preferences | $\ell : \Theta \times D \to \mathbb{R}$ | $u : F \to \mathbb{R}$ |
| Loss in a given state | $\ell(\theta, \cdot), \theta \in \Theta$ | $\kappa u(\cdot), (\kappa, \mu) \in \mathcal{T}$ |

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.

3

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

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

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

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

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

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

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

**Definition 3.4** (Statistical Experiment). A *statistical experiment* relative to a set $\Theta$, a measurable space $(E, \mathcal{E})$ and a map $m : \Theta \to \Delta(\mathcal{E})$ is a set $\mathcal{H} = \{\mu_\theta | \theta \in \Theta\}$ where $\mu_\theta := m(\theta)$. The set $\Theta$ indexes the "state of nature".

**Definition 3.5** (Statistical Decision Problem). A statistical decision problem (SDP) is a tuple $\langle \Theta, (\mathcal{H}, m), D, \ell \rangle$. $\mathcal{H} \subset \Delta(\mathcal{E})$ is a statistical experiment relative to states $\Theta$, space $(E, \mathcal{E})$ and map $m : \Theta \to \Delta(\mathcal{E})$, $D$ is the set of available decisions with some $\sigma$-algebra $\mathcal{D}$ and $\ell : \Theta \times D \to \mathbb{R}$ is a loss function where $\ell(\theta, \cdot)$ is measurable with respect to $\mathcal{D}$ and $\mathcal{B}(\mathbb{R})$.

Denote by $\mathcal{J}$ the set of stochastic decision functions $E \to \Delta(\mathcal{D})$. For $J \in \mathcal{J}$ and $\mu_\theta \in \mathcal{H}$, the risk $R : \Theta \times \mathcal{J} \to [0, \infty)$ is defined as $R(J, \theta) = \int_D \ell(\theta, y)\mu_\theta J(dy)$. The triple $\langle \Theta, \mathcal{J}, R \rangle$ forms a two player normal form game.

The loss function $\ell$ expresses preferences over general (state, decision) pairs. It may be the case that our preferences are most directly known over future states of the world - we know which results of our decisions are desirable and which are undesirable, which we represent with a *utility function*. In this case, if we are to induce preferences over the possible decisions, that we have a model that is more informative than a statistical experiment. In particular, we require each state of nature to be associated with both a distribution over the given information and a map from decisions to distributions over results - we call this map a *consequence*, and the object that pairs a distribution and a consequence with each state of the world a *causal theory*.

**Definition 3.6** (Consequences). Given a measurable result space $(F, \mathcal{F})$ and a measurable decision space $(D, \mathcal{D})$, a Markov kernel $\kappa : D \to \Delta(\mathcal{F})$ is a *consequence mapping*, or just a *consequence*.

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

In many cases the observation space $E$ and the results space $F$ might coincide. However, these spaces are defined by different aspects of the given information: the former is fixed by what observations are available and the latter by which parts of the world are relevant to the investigator's preferences (see Theorems B.7 and B.6), and there is not a clear reason to insist that these spaces should always be the same.

**Definition 3.8** (Causal Theory). A causal theory $\mathcal{T}$ is a set of causal states sharing the same decision, observation and outcome spaces. We abuse notation to assign the "type signature" $\mathcal{T} : E \times D \rightharpoonup F$ for a causal theory with observed distributions in $\Delta(\mathcal{E})$ and consequences of type $D \to \Delta(\mathcal{F})$. The causal states of a theory $\mathcal{T}$ may be associated with a master set of states $\Theta$, but in contrast to a statistical experiment this is not necessary to define the basic associated decision problem.

4

**Definition 3.9** (Causal Statistical Decision Problem). A causal statistical decision problem (CSDP) is a triple $\langle \mathcal{T}, D, u \rangle$. $\mathcal{T}$ is a causal theory on $D \times E \to F$, $D$ is the decision set with $\sigma$-algebra $\mathcal{D}$ and $u : F \to \mathbb{R}$ is a measurable utility function expressing preference over the results of decisions.

Define the canonical loss $L : \mathcal{T} \times D \to \mathbb{R}$ by $L : (\kappa, \mu), y \mapsto -\mathbb{E}_{\gamma\kappa}[u]$. This change conforms with the conventions that utilities are maximised while losses are minimised.

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

The loss and the utility differ in that the loss expresses per-state preferences while the utility expresses state independent preferences. While we choose the loss to be a particular function of the utility here, it is possible to allow losses to be a more general class of functions of the utility and state without altering the preference ordering of a CSDP under minimax or Bayes decision rules. Given arbitrary $f : \mathcal{T} \to \mathbb{R}$, define $l : \mathcal{T} \times D \to \mathbb{R}$ by $l : (\kappa, \mu, y) \mapsto af(\kappa, \mu) + b\mathbb{E}_{\delta_y\kappa}[u]$. We can define a loss (relative to $f$) $L : \mathcal{T} \times \Delta(\mathcal{D}) \to [0, \infty]$ by

$$L((\kappa, \mu), \gamma) := \mathbb{E}_\gamma[l(\kappa, \mu, \cdot)] \tag{2}$$

$$= af(\kappa, \mu) - b\mathbb{E}_{\gamma\kappa}[u] \tag{3}$$

$$\tag{4}$$

For $(\kappa, \mu) \in \mathcal{T}$, $\gamma \in \Delta(\mathcal{D})$ and $a \in \mathbb{R}$, $b \in \mathbb{R}^+$.

A common example of a loss of the type above is the *regret*, which takes $a = b = 1$ and $f(\kappa, \mu) = \sup_{\gamma' \in \Delta(\mathcal{D})} \mathbb{E}_{\gamma'\kappa}[u]$. Because expected utility preserves preference orderings under positive affine transformations, the ordering of preferences given a particular state is not affected by the choices of $a, b$ and $f$, nor is the Bayes ordering of preferences given some prior $\xi$ over $\mathcal{T}$. While it may be possible to formulate decision rules for which the choices of $a, b$ and $f$ do matter, we will take these properties as sufficient to allow us to choose $a = 0$ and $b = 1$. More general classes of loss are of interest. *Regret theory*, for example, is a straightforward generalisation of the losses discussed here and is a prominent alternative to expected utility theory [Loomes and Sugden, 1982].

There are obvious similarities between SDPs and CSDPs: both have the same high level representation as a two person game which is arrived at by taking the expectation of a loss with respect to a decision function. In fact, if we consider two decision problems to be the same if they have the same representation as a two player game, we find that CSDPs are a special case of SDPs.

**Theorem 3.10** (CSDPs are a special case of SDPs). *Given any CSDP $\alpha = \langle \mathcal{T}, D, u \rangle$ with two player game representation $\langle \mathcal{T}, \mathcal{J}, R \rangle$, there exists an SDP $\langle \mathcal{T}, (\mathcal{H}, m), D, \ell \rangle$ with the same representation as a two player game.*

*Proof.* Let $m : \mathcal{T} \to \mathcal{H}$ be defined such that $m : (\kappa, \mu) \mapsto \mu$ for $(\kappa, \mu) \in \mathcal{H}$. Define $\ell : \mathcal{T} \times D \to \mathbb{R}$ by $\ell : ((\kappa, \mu), y) \mapsto -\mathbb{E}_{\delta_y\kappa}[u]$. Let $R'((\kappa, \mu), J) = \mathbb{E}_{\mu J}[\ell(\theta, \cdot)]$. Then

$$R'((\kappa, \mu), J) = -\int_D \mathbb{E}_{\delta_y\kappa}[u]\mu J(dy) \tag{5}$$

$$= -\int_D \int_F u(x)\kappa(y; dx)\mu J(dy) \tag{6}$$

$$= -\int_F u(x)\mu J\kappa(dx) \tag{7}$$

$$= R((\kappa, \mu), J) \tag{8}$$

$\square$

The converse is not true, as the set $\Theta$ in an SDP is of an arbitrary type and may not be a causal theory. However, it is possible for any SDP with environmental states $\Theta$ to find a CSDP with causal theory $\mathcal{T}$ such that the games represented by each decision problem are related by a surjective map $f : \Theta \to \mathcal{T}$ which associates each state of nature with a causal state. We call such a map a *reduction* from an SDP to a CSDP.

**Definition 3.11** (Reduction). Given normal form two person games $\alpha = \langle \mathcal{S}^\alpha, A, L^\alpha \rangle$ and $\beta = \langle \mathcal{S}^\beta, A, L^\beta \rangle$, $f : \mathcal{S}^\alpha \to \mathcal{S}^\beta$ is a *reduction* from $\alpha$ to $\beta$ if, defining the image $f(\mathcal{S}^\alpha) = \{f(\theta)|\theta \in \mathcal{S}^\alpha\}$, we have $\langle \mathcal{S}^\beta, A, L^\beta \rangle = \langle f(\mathcal{S}^\alpha), A, L^\alpha \circ (f \otimes I_A) \rangle$.

**Theorem 3.12** (SDP can be reduced to a CSDP). *Given any SDP $\langle \Theta, (\mathcal{H}, m), D, \ell \rangle$ represented as the game $\alpha = \langle \Theta, \mathcal{J}, R \rangle$, there exists a CSDP $\langle \mathcal{T}, D, u \rangle$ represented as the game $\beta = \langle \mathcal{T}, \mathcal{J}, R' \rangle$ such that there is some reduction $f : \Theta \to \mathcal{T}$ from $\alpha$ to $\beta$.*

*Proof.* Take $\mathcal{H} \subset \Delta(\mathcal{E})$ and define $f : \Theta \to \Delta(\mathcal{E}) \times \Delta(\mathcal{B}(\mathbb{R}))^D$ by $f : \theta \mapsto (y \mapsto \delta_{l(\theta,y)}, \mu_\theta)$. Noting that $y \mapsto \delta_{l(\theta,y)}$ is a Markov kernel $D \to \Delta(\mathcal{B}(\mathbb{R}))$, the image $f(\Theta)$ is a causal theory $E \times D \twoheadrightarrow \mathbb{R}$. Consider the CSDP $\langle f(\Theta), D, -I_{(\mathbb{R})} \rangle$. Then, letting $R'$ denote the risk associated with this theory

$$R'((\kappa, \mu), J) = -\int_{\mathbb{R}} \int_D (-x) \delta_{l(\theta,y)}(dx) \mu_\theta J(dy) \tag{9}$$

$$= \int_D l(\theta, y) \mu_\theta J(dy) \tag{10}$$

$$= R(\Theta, J) \tag{11}$$

$\square$

The fundamental *complete class theorem* of SDPs establishes that there are no decision rules that dominate the set of all Bayes rules under some regularity assumptions. By theorem 3.10, this must also be true of CSDPs.

**Theorem 3.13** (Complete class theorem (CSDP)). *Given any CSDP $\alpha := \langle \mathcal{T}, D, u \rangle$ with two player game representation $\langle \mathcal{T}, \mathcal{J}, R \rangle$, if $|\mathcal{T}| < \infty$ and $\inf_{J \in \mathcal{J}, (\kappa,\mu) \in \mathcal{H}} R((\kappa, \mu), J) > -\infty$, then the set of all Bayes decision functions is a complete class for $\alpha$ and the set of all admissible Bayes decision functions is a minimal complete class for $\alpha$.*

*Proof.* By theorem 3.10, there exists an SDP $\beta$ such that $\alpha$ and $\beta$ have the same representation as a two player game. By assumption, $\beta$ has a finite set of states and a risk function that is bounded below. Therefore the Bayes rules on $\alpha$ are a complete class and admissible Bayes rules are a minimal complete class for the problem $\langle \mathcal{T}, \mathcal{J}, R \rangle$ [Ferguson, 1967]. $\square$

# 4 Causal Bayesian Networks

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

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

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

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

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 independent of $X^j$ conditional on $\mathbf{X}$ if $\mu_{|\mathbf{X}}^\curlyvee (F_{X^i} \otimes F_{X^j}) = \mu_{|\mathbf{X}} F_{X^i} \mu_{|\mathbf{X}} F_{X^j}$, $\mu$-almost surely. This is written $X^i \perp\!\!\!\perp_\mu X^j | \mathbf{X}$.

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

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

A CBN has a graph $\mathcal{G}$ with edges $\{V^i\}_{[N]}$, random variables $\{X^i\}_{[N]}$ and deci-

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

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

1. $\mu^{\mathcal{G}, y}$ is compatible with $\mathcal{G}$

2. For all $i \in S(y)$, $\mu^{\mathcal{G}, y} F_{X^i} = \delta_{D^i(y)} F_{X^i}$

3. For all $i \notin S(y)$, $\mu^{\mathcal{G}, y}_{\mathrm{Pa}_{\mathcal{G}}(X^i)} F_{X^i} = \mu_{|\mathrm{Pa}_{\mathcal{G}}(X^i)} F_{X^i}$, $\mu^{\mathcal{G}, y}$-almost surely

$\mathrm{Pa}_{\mathcal{G}}(X^i)$ are the parents of $X^i$ with respect to the graph $\mathcal{G}$ and $\mu_{|\mathrm{Pa}_{\mathcal{G}}(X^i)}$ is the conditional probability with respect to $\mu$ and the $\sigma$-algebra generated by the set $\mathrm{Pa}_{\mathcal{G}}(X^i)$. Recall that $\mu^{\curlyvee}(\otimes_{i \notin S(y)} F_{X^i})$ is the joint distribution of $\{X^i | i \in S(y)\}$.

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

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

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

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

$$\max_{x \in \{0,1\}^3, y \in \{0,1\}} \left| \mu_{|\{A,C\}} F_B(x; \{y\}) - \iota(x; \{y\}) \right| < \epsilon \tag{12}$$

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

$$\kappa^{\mathcal{G}}(y; Z) := \begin{cases} \iota(D_A(y); Z) & D_A(y) \neq * \\ \zeta(Z) & D_A(y) = * \end{cases} \tag{13}$$

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

Consider two options for extending this to distributions $\nu \in \mathcal{H}$ but not in $\mathcal{H}_{\mathcal{G}}$, noting that one could imagine many possibilities: $\mathcal{T}_{\mathcal{G}}^{\subseteq}$ is the union of causal theories given by all graphs $\mathcal{G}'$ on $\{A, B, C\}$ such that $\mathcal{G} \subset \mathcal{G}'$ (in this case, just $\mathcal{G}$ and $C \overset{\frown}{\longrightarrow} A \longrightarrow B$ ), and $\mathcal{T}_{\mathcal{G}}^{\circ}$ is the union of causal theories given by the all DAGs on the set of nodes $\{A, B, C\}$.

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

$$\eta^{\nu} := \begin{cases} (y; Z) \mapsto \sum_{c \in \{0,1\}} \nu F_C(\{c\}) \nu_{|\{A,C\}} F_B(D_A(y), c; Z) & D_A(y) \neq * \\ \zeta(Z) & D_A(y) = * \end{cases} \tag{14}$$

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

299    those featuring $A \longrightarrow B$ but not $C \longrightarrow B$ and $C \longrightarrow A$ and the graph $C \overset{\frown}{\longrightarrow} A \longrightarrow B$. These
300    possibilities yield $\mathcal{T}_{\mathcal{G}}^{\circ} = \mathcal{T}_{\mathcal{G}}^{\subseteq} \cup \{(\nu, y \mapsto \zeta) | \nu \in \mathcal{H} \setminus \mathcal{H}_{\mathcal{G}}\}$.

301    By 12, $|\eta(x; \{y\}) - \iota(x; \{y\})| < \epsilon$ for all $x \in A \cup \{*\}$ and $y \in B$ and therefore for $J \in \mathcal{J}$,
302    $|U(\mu J \curlyvee (I_{(D)} \otimes \eta)) - U(\mu J \curlyvee (I_{(D)} \otimes \iota))| < \epsilon$. Therefore a small $\epsilon$ ensures $\mathcal{T}_{\mathcal{G}}^{\subseteq}$ yields a risk set
303    "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 $\epsilon$,
304    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
305    small $\epsilon$.

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

315 # 5    Potential Outcomes

316    Potential Outcomes is an alternative to the approach typified by Causal Bayesian Networks for
317    formulating causal questions and hypotheses. Causal queries in the Potential Outcomes framework
318    concern the distribution of random variables $\mathsf{X}_0, \mathsf{X}_1$ representing potential outcomes, or "the value
319    $\mathsf{X}$ would have taken if action 0 or 1 were taken respectively" (Hernán and Robins [2018]). This is
320    similar, but not the same, as the question answered by a consequence map which is "what is the
321    distribution of $\mathsf{X}$ if I take actions 0 or 1?"

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

$$(\delta_i \kappa)_{|\mathsf{X}_i} F_{\mathsf{X}}(w; A) = \delta_{\mathsf{X}_i(w)}(A) \tag{15}$$

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

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

331 # 6    Equivalence of causal problems

332    Under what conditions could we consider a consequence consistent theory $\mathcal{T}^{cc}$ associated with some
333    distribution over potential outcomes to be "equivalent" to some causal theory $\mathcal{T}^{\mathcal{G}}$ associated with a
334    CBN $\mathcal{G}$ or vise versa?

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

338    The question of whether a consequence consistent theory $\mathcal{T}^{cc}$ can in general be represented by a
339    Causal Bayesian Network is then also straightforwardly answered in the negative, as conditions 2 and
340    3 of Definition 4.2 are in general non-trivial (condition 1 is trivial given a fully connected DAG $\mathcal{G}$).

341    The trivial potential outcome $\mathsf{X}_y = \mathsf{X}$ clashes with the informal idea that a potential outcome
342    represents the value $\mathsf{X}$ would have taken had action $y$ been taken - we might expect, for example, if
343    $\delta_y \kappa F_{\mathsf{X}} \neq \mu F_{\mathsf{X}}$ then $\mathsf{X}$ would at least sometimes take a different value if the action $y$ is taken than if it
344    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 $\mathsf{X} : E \to X$ and potential outcome variable $\mathsf{X}_i : E \to X$, $i \in D$ if for every $(\mu, \kappa) \in \mathcal{T}$, $\kappa$ is consequence consistent (Eq. 15) and

$$\mu F_{\mathsf{X}_i} = \delta_i \kappa F_{\mathsf{X}_i} \tag{16}$$

> 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:

$$\mu_{|\{\mathsf{X}_i, \mathsf{Z}\}} F_{\mathsf{X}}(w; A) = \delta_{\mathsf{X}_i(w)}(A) \qquad w \in \mathsf{Z}^{-1}(i) \tag{17}$$

Where the random variable $\mathsf{Z}$ 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.

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

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

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

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

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

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

> In both cases it is straightforward to posit generalised utilities such that $\alpha$ cannot be reduced.

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

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

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

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

9

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

392
> The single world intervention graphs of Richardson and Robins [2013] are DAGs that represent independences among distributions over potential outcome variables. They might be interpretable as POCSDPs.

393
> The generalised versions of CBNs yield theories that generally associate multiple consequences with each given distribution. However a generalized CBN still yields a unique causal theory

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

396 *Proof.* Suppose $D = [M]$ or $D = \mathbb{N}^+$. Take $E' = E \times E^D$ and for $i \in D \cup \{0\}$, $x := (x_0, x_1, ...) \in$
397 $E'$ define the projection $\mathsf{P}_i(x_0, x_1, ...) := x_i$ and the potential outcome variable $\mathsf{X}_i := \mathsf{X} \circ \mathsf{P}_i$.

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

$$\mu^{po}(A_1 \times ...) := \prod_{y' \in D} \delta_{y'} \kappa(A_{y'}) \tag{18}$$

$$\kappa'(y; A_0 \times A_1 \times ...) := \int_{A_1 \times ...} \delta_{x_y}(A_0) \mu^{po}(dx) \tag{19}$$

$$= \prod_{y' \in D \setminus \{y\}} \delta_{y'} \kappa(A_{y'}) \int_{A_y} \delta_{x_y}(A_0) \delta_y \kappa(dx_y) \tag{20}$$

$$\mu'(A_0 \times A_1 \times ...) := \prod_{y' \in D \setminus \{y\}} \delta_{y'} \kappa(A_{y'}) \int_{A_y} \delta_{x_y}(A_0) \mu(dx_y) \tag{21}$$

400 It can be verified that $\kappa'$ is a Markov kernel.

401 Note that by the definition of conditional probability, for $A, B \in \mathcal{X}$,
402 $\int_{\mathsf{X}_y^{-1}(A)} (\delta_y \kappa')_{|\mathsf{X}_y} F_{\mathsf{X}_0}(x; B) \delta_y \kappa'(dx) = \delta_y \kappa' \curlyvee (F_{\mathsf{X}_0} \otimes F_{\mathsf{X}_y})(A, B)$. Thus by 20, $\delta_x(A)$ is a
403 version of $(\delta_y \kappa')_{|\mathsf{X}_y} F_{\mathsf{X}_0}(x; A)$, so $\kappa'$ is consequence consistent.

404 Furthermore, $\mu' F_{\mathsf{X}_y} = \delta_y \kappa F_{\mathsf{X}} = \delta_y \kappa' F_{\mathsf{X}_y}$ for $y \geq 1$. Therefore defining $\mathcal{T}'$ to be the image of $\mathcal{T}$
405 under $f$, we can see that $\mathcal{T}'$ is a PO causal theory with respect to "observable" $\mathsf{X}_0$ and "potential
406 outcomes" $\mathsf{X}_y, y \in D$.

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

$$\kappa' F_{\mathsf{P}_0}(y; A) = \int_E \delta_z(A) \delta_y \kappa(dz) \tag{22}$$

$$= \int_A \kappa(y; dz) \tag{23}$$

$$= \kappa(y; A) \tag{24}$$

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

$$\mu' F_{\mathsf{X}_0}(B) = \int_E \delta_z(\mathsf{X}^{-1}(B)) \mu(dz) \tag{25}$$

$$= \mu F_{\mathsf{X}}(B) \tag{26}$$

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

$$U(\mu F_{\mathsf{X}} J \curlyvee (I_{(D)} \otimes \kappa)) = U(\mu' F_{\mathsf{X}_0} J \curlyvee (I_{(D)} \otimes \kappa' F_{\mathsf{P}_0})) \tag{27}$$

$$\tag{28}$$

10

Therefore, given the PO CSDP $\beta = \langle (\mathcal{T}', E', \mathsf{X}_0), D, (U, E) \rangle$, for all $J \in \mathcal{J}$, $R^\alpha(J, \kappa, \mu) = R^\beta(J, f(\kappa, \mu))$. Thus $\beta$ is a reduction of $\alpha$ witnessed by $f$. $\qquad\square$

**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 itelf 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 Schoelkopf, 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.

Graham Loomes and Robert Sugden. Regret Theory: An Alternative Theory of Rational Choice Under Uncertainty. *The Economic Journal*, 92(368):805–824, 1982. ISSN 0013-0133. doi: 10.2307/2232669. URL https://www.jstor.org/stable/2232669.

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.

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

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

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

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

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

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

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

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

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

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

**Supplement to: Causal Statistical Decision Problems**

## A  Markov Kernels

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

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

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

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

511        1. The map $x \mapsto K(x; B)$ is $\mathcal{E}$-measurable for every $B \in \mathcal{F}$

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

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

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

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

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

$$\int_A K(x; B) d\mu = \int_A K'(x; B) d\mu \qquad \forall A \in \mathcal{E}, B \in \mathcal{F} \tag{29}$$

### A.1  Operations with Markov kernels

529  For the following, assume $K$ is a Markov kernel from $E \to \Delta(\mathcal{F})$, $K'$ a kernel $E \to \Delta(\mathcal{H})$, $\mathsf{L}$ is
530  a Markov kernel $F \to \Delta(\mathcal{G})$, $\mu$ is a probability measure on $(E, \mathcal{E})$, $\nu$ is a probability measure on
531  $(F, \mathcal{F})$ and $f$ is a nonnegative measurable function $F \to \mathbb{F}$.

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

### A.1.1  Kernel products

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

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

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

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

### A.1.2  Special kernels

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

544 $\curlyvee_E$ is a kernel $E \to \Delta(\mathcal{E} \otimes \mathcal{E})$ defined by $x \mapsto \delta_{(x,x)}$. We will subsequently leave the space implicit.
545 The symbol $\curlyvee$ is pronounced "splitter".

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

$$K \otimes M(x, y; A \times B) := K(x; A)M(y; B) \tag{30}$$

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

548 $\curlyvee(K \otimes K')$ is a Markov kernel $E \to \Delta(\mathcal{F} \otimes \mathcal{H})$ and

$$\curlyvee(K \otimes K')(x; A \times B) = \int_E K(x'; A)K'(x''; B)\delta_{(x,x)}(dx' \times dx'') \tag{31}$$

$$= K(x; A)K'(x; B) \tag{32}$$

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

550 Let $(*, \{\emptyset, *\})$ be an indiscrete measurable set. $\upharpoonleft_E$ is a kernel $E \to \Delta(\{\emptyset, *\})$ defined by $x \mapsto \mathbb{1}_*$.
551 We have $\curlyvee(I \otimes \upharpoonleft) = I$. The symbol $\upharpoonleft$ is pronounced "stopper".

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

$$\mu F_g(A) = \int_E \delta_{g(x)}(A)d\mu \tag{33}$$

$$= \mu(g^{-1}(A)) \tag{34}$$

$$= g_*\mu(A) \tag{35}$$

554 Given two random variables $\mathsf{X} : (E, \mathcal{E}) \to (X, \mathcal{X})$ and $\mathsf{Y} : (E, \mathcal{E}) \to (Y, \mathcal{Y})$, the product $\mu\curlyvee(F_\mathsf{X} \otimes$
555 $F_\mathsf{Y})$ is the joint distribution of $\mathsf{X}$ and $\mathsf{Y}$.

$$\mu\curlyvee(F_\mathsf{X} \otimes F_\mathsf{Y})(A, B) = \int_E \delta_{\mathsf{X}(x)}(A)\delta_{\mathsf{Y}(x)}(B)d\mu \tag{36}$$

$$= \mu(\mathsf{X}^{-1}(A) \cap \mathsf{Y}^{-1}(B)) \tag{37}$$

## B  Appendix: Causal Statistical Decision Problems

It is possible to define a generalised CSDP where preferences may not obey the Von Neumann-Morgenstern axioms, but this generalisation presents some difficulties.

**Definition B.1** (Generalised Causal Statistical Decision Problem)**.** A causal statistical decision problem (CSDP) is a tuple $\langle(\mathcal{T}, (E, \mathcal{E})\mathsf{X}), (D, \mathcal{D}), (U, (F, \mathcal{F}))\rangle$. $\mathcal{T}$ is a causal theory on $D, E$ and $F$, $D$ is the decision set, $\mathsf{X} : (E, \mathcal{E}) \to (X, \mathcal{X})$ is a random variable representing the given information and $U : \Delta(\mathcal{F} \otimes \mathcal{D}) \to \mathbb{R}$ 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 $L : \mathcal{T} \times \Delta(\mathcal{D}) \to [0, \infty]$ by

$$L((\kappa, \mu), \gamma) := \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \curlyvee (I_{(D)} \otimes \kappa)) - U(\gamma \curlyvee (I_{(D)} \otimes \kappa)) \tag{38}$$

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 $\mu$; 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 $R : \mathcal{J} \times \mathcal{T} \to [0, \infty)$ by $R(J, \kappa, \mu) := L(\kappa, \mu F_{\mathsf{X}} J)$. The triple $\langle \mathcal{T}, \mathcal{J}, R \rangle$ is a normal form two person game.

If there exists some measurable $u : F \times D \to \mathbb{R}$ 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 $l^\kappa : D \to [0, \infty)$ is defined by

$$l^\kappa(d) := \sup_{\gamma' \in \Delta(\mathcal{D})} \mathbb{E}_{\gamma' \curlyvee (I_{(D)} \otimes \kappa)}[u] - \mathbb{E}_{\kappa(d; \cdot)}[u(\cdot, d)] \tag{39}$$

**Lemma B.2** (Reduction preserves admissibility)**.** *If a CSDP $\beta$ with induced game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be reduced to a statistical decision problem $\alpha$ with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$ then a decision function $J \in \mathcal{J}$ is admissible in $\beta$ iff it is admissible in $\alpha$.*

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

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

**Corollary B.3** (Reduction preserves completeness)**.** *If a causal decision problem $\beta$ with induced game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be reduced to a statistical decision problem $\alpha$ with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$, then an (essentially) complete class with respect to $\alpha$ is (essentially) complete with respect to $\beta$.*

**Lemma B.4** (Induced Bayes rule)**.** *If a CSDP $\beta$ with induced game $\langle \mathcal{T}, \mathcal{J}, R \rangle$ can be reduced to a statistical decision problem $\alpha$ with induced game $\langle \mathcal{H}, \mathcal{J}, R' \rangle$ witnessed by $h : \mathcal{T} \to \mathcal{H}$ and $J_{ba}^\xi \in \mathcal{J}$ is a Bayes rule with respect to the problem $\alpha$ and the prior $\xi$ then $J_{ba}^\xi$ is a Bayes rule with respect to the problem $\beta$ and the induced prior $\xi_h$.*

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

$$\int_{\mathcal{T}} R(J, \tau) d\xi_h = \int_{\mathcal{H}} R'(J, h(\tau)) d\xi \tag{40}$$

And therefore, if a Bayes rule exists,

$$\arg\min_{J \in \mathcal{J}} \int_{\mathcal{T}} R(J, \tau) d\xi_h = \arg\min_{J \in \mathcal{J}} \int_{\mathcal{H}} R'(J, h(\tau) d\xi \tag{41}$$

$\qquad\square$

15

**Theorem B.5** (Complete class theorem (CSDP)). *Given an CSDP $\alpha := \langle (\mathfrak{T}, E), D, \mathsf{X}, U \rangle$ with risk $R$, if there is a reduction to an SDP $\beta := \langle (\mathcal{H}, F), D, \mathsf{Y}, \ell \rangle$ 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 and the set of all admissible Bayes decision functions is a minimal complete class.*

*Proof.* Given the conditions, the Bayes decision functions in $\beta$ form a complete class and admissible Bayes rules a minimal complete class [Ferguson, 1967].

By Corollary B.3 the Bayes rules for $\beta$ are complete in $\alpha$, and the admissible Bayes rules for $\beta$ are essentially complete in $\alpha$.

Every (admissible) Bayes rule for $\beta$ is a(n admissible) Bayes rule for $\alpha$, so the set of (admissible) Bayes rules for $\alpha$ is also (essentially) complete in $\alpha$. $\qquad\square$

**Theorem B.6** (Reduction of a CSDP on observations). *A CSDP $\alpha = \langle (\mathfrak{T}^\alpha, (E, \mathcal{E}), \mathsf{X}), D, (U, (F, \mathcal{F})) \rangle$ where, for $\zeta \in \Delta(\mathcal{E} \otimes \mathcal{D})$ can be reduced to a problem $\beta = \langle (\mathfrak{T}^\beta, (X, \mathcal{X}), \mathrm{id}_X), D, (U, (F, \mathcal{F}) \rangle$ by marginalization.*

*Proof.* Consider the mapping $g : \mathfrak{T}^\alpha \to \mathfrak{T}^\beta$ given by $(\kappa, \mu) \mapsto (\kappa, \mu F_\mathsf{X})$.

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

$$R^\alpha(J, \kappa, \mu) = \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \curlyvee (I_{(D)} \otimes \kappa)) - U(\mu F_\mathsf{X} J \curlyvee (I_{(D)} \otimes \kappa)) \tag{42}$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \curlyvee (I_{(D)} \otimes \kappa) - U(\mu F_\mathsf{X} F_\mathsf{X} J \curlyvee (I_{(D)} \otimes \kappa))) \tag{43}$$

$$= R^\beta(J, g(\kappa, \mu)) \tag{44}$$

$\qquad\square$

**Theorem B.7** (Reduction of a CSDP on the utilty). *Given a CSDP $\alpha = \langle (\mathfrak{T}^\alpha, (E, \mathcal{E}), \mathsf{X}), D, (U, (F, \mathcal{F})) \rangle$ where, for $\zeta \in \Delta(\mathcal{E} \otimes \mathcal{D})$, if $U(\zeta) = U'(\zeta(I_{(D)} \otimes F_\mathsf{Y}))$ for some $\mathsf{Y} : F \to Y$ and $U' : \Delta(\mathcal{Y}) \to \mathbb{R}$ then $\alpha$ has $\mathsf{Y}$-observable utility. Such a problem can be reduced to a problem $\beta = \langle (\mathfrak{T}^\beta, (E, \mathcal{E}), \mathsf{X}), D, (U', (Y, \mathcal{Y}) \rangle$ by marginalization.*

*Proof.* Consider the mapping $g : \mathfrak{T}^\alpha \to \mathfrak{T}^\beta$ given by $(\kappa, \mu) \mapsto (\kappa F_\mathsf{Y}, \mu)$.

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

$$R^\alpha(J, \kappa, \mu) = \sup_{\gamma' \in \Delta(\mathcal{D})} U(\gamma' \curlyvee (I_{(D)} \otimes \kappa)) - U(\mu F_\mathsf{X} J \curlyvee (I_{(D)} \otimes \kappa)) \tag{45}$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U'(\gamma' \curlyvee (I_{(D)} \otimes \kappa)(I_{(D)} \otimes F_\mathsf{Y})) - U'(\mu F_\mathsf{X} J \curlyvee (I_{(D)} \otimes \kappa))(I_\mathsf{D} \otimes F_\mathsf{Y})) \tag{46}$$

$$= \sup_{\gamma' \in \Delta(\mathcal{D})} U'(\gamma' \curlyvee (I_{(D)} \otimes \kappa F_\mathsf{Y})) - U'(\mu F_\mathsf{X} J \curlyvee (I_{(D)} \otimes \kappa F_\mathsf{Y})) \tag{47}$$

$$= R^\beta(J, g(\kappa, \mu)) \tag{48}$$

$\qquad\square$

**Theorem B.8.** *Every SDP $\langle (\mathcal{H}, E, \mathsf{X}), D, \ell \rangle$ can be reduced to a CSDP.*

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

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

$$L(\kappa_\mu, \gamma) = - \sup_{\gamma' \in \Delta(\mathcal{D})} \mathbb{E}_{\gamma' \curlyvee (I_{(D)} \otimes \kappa_\mu)_{|\mathsf{E}}} [\ell(\gamma' \curlyvee (I_{(D)} \otimes \kappa_\mu)_{|\mathsf{E}}, \mathsf{D})] + \mathbb{E}_{\gamma \curlyvee (I_{(D)} \otimes \kappa_\mu)} [\ell(\gamma \curlyvee (I_{(D)} \otimes \kappa_\mu)_{|\mathsf{E}}, \mathsf{D})] \tag{49}$$

$$= \mathbb{E}_\gamma [\ell(\mu, \mathsf{D})] \tag{50}$$

16

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

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

$$R'(J, \kappa, \mu) = \int_D \ell(\mu, y) \mu F_\mathsf{X} J(dy) \tag{51}$$

$$= R(J, g(\kappa, \mu)) \tag{52}$$

627  $\square$

628  **Theorem B.9.** *Given a CSDP $\beta = \langle (\mathcal{T}, E, \mathsf{X}), D, (U, F) \rangle$ where $U$ is an ordinary pseudo-utility, let*
629  $\mathcal{K} = \{\kappa | (\kappa, \mu) \in \mathcal{T}\}$ *be the set of consequences. $\beta$ is reducible to a statistical decision problem on*
630  *the measurable space $(E \times F \times D, \mathcal{E} \otimes \mathcal{F} \otimes \mathcal{D})$ if there is some surjective map $m : \Delta(\mathcal{F} \otimes \mathcal{D}) \to \mathcal{K}$.*

631  *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
632  $h : \mathcal{T} \to \mathcal{H}$ by $(\kappa, \mu) \mapsto \mu \otimes m^\dagger(\kappa)$.

633  Let $k : \Delta(\mathcal{F})^D \times D \to \mathbb{R}$ be the differential loss induced by the ordinary pseudo-utility $U$ (see
634  Equation 39).

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

$$\ell(\xi, y) = k(m(\xi F_{\curlyvee_{(\mathsf{F} \otimes \mathsf{D})}}), y) \tag{53}$$

637  Note that

$$\ell(h(\kappa, \mu), y) = k(\kappa, y) \tag{54}$$

638  Define $\mathsf{X}' : E \times F \times D \to X$ by $(a, b, c) \mapsto \mathsf{X}(a)$.

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

$$R'(J, h(\kappa, \mu)) = \int_D \ell(h(\kappa, \mu), y) h(\kappa, \mu) F_{\mathsf{X}'} J(dy) \tag{55}$$

$$= \int_D \ell(h(\kappa, \mu), y) (\mu \otimes m^\dagger(\kappa)) F_{\mathsf{X}'} J(dy) \tag{56}$$

$$= \int_D k(\kappa, y) \mu F_\mathsf{X} J(dy) \tag{57}$$

$$= R(J, \kappa, \mu) \tag{58}$$

641  $\square$

642  **Example B.10** (Irreducible CSDP). The choice of decision function in an SDP does not affect the
643  state, while this choice does affect the outcome in an CSDP. For an SDP, then, the risk of a mixed
644  decision function is equal to the mixture of risks of each atomic decision function but this is not true
645  in general for an CSDP.

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

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

$$R(J, \kappa, \mu) = 0.25 - \mathrm{Var}_{\mu F_\mathsf{X} J}(\mathsf{Y}) \tag{59}$$

649  Consider the forgetful decision functions $J_0 : x \mapsto \mathrm{Bernoulli}(0)$ and $J_{1/2} : x \mapsto \mathrm{Bernoulli}(\frac{1}{2})$ and
650  $J_1 : x \mapsto \mathrm{Bernoulli}(1)$ 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
651  $x \in X, A \in \mathcal{D}$. For any statistical decision problem with risk $R'$,

$$R'(J_{1/2}, \mu) = \int_D \ell(\mu, y) \mu F_\mathsf{X} J_{1/2}(dy) \tag{60}$$

$$= \frac{1}{2} \left( \int_D \ell(\mu, y) \mu F_\mathsf{X} J_0(dy) + \int_D \ell(\mu, y) \mu F_\mathsf{X} J_1(dy) \right) = \frac{1}{2} \left( R'(J_0, \mu) + R'(J_1, \mu) \right) \tag{61}$$

But

$$R(J_{1/2}, \kappa, \mu) = 0 \tag{62}$$

$$\neq \frac{1}{2} \left( R(J_0, \kappa, \mu) + R(J_1, \kappa, \mu) \right) \tag{63}$$

**Corollary B.11.** *The class of nonrandomized decision functions is not essentially complete for CSDPs. The stochastic decision function $J_{1/2}$ is strictly better than any deterministic function in the 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 $\{\mathsf{X}^i\}_{i \in [N]}$ where $\mathsf{X}^i : E \to X^i$, a decision set $(D, \mathcal{D})$ and random variables $\{\mathsf{D}^i\}_{i \in [N]}$ with $\mathsf{D}^i : (D, \mathcal{D}) \to (X^i \cup \{*\}, \sigma(\mathcal{X}^i \cup \{*\}))$. Given $\mu \in \mathcal{G}_\mathcal{G}$, let $\mu^y$ be the $\mathcal{G}, \mu, y$-interventional distribution (Definition 4.2).*

*Then the map $\kappa^{\mu, \mathcal{G}} : D \to \Delta(\mathcal{E})$ 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 $\mathsf{X}^i$ by $\mathsf{X}^i < \mathsf{X}^j$ if $\mathsf{X}^i \to \mathsf{X}^j$ is in $\mathcal{G}$. Without loss of generality, suppose the total ordering $X^0, ..., X^N$ is consistent with the partial ordering induced by $\mathcal{G}$.

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

$$\kappa^i(x; A) = \mu_{|\mathrm{Pa}_\mathcal{G}(\mathsf{X}^i)} F_{\mathsf{X}^i}(x; A) \tag{64}$$

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

$$\kappa^{i,*}(y, pa^i; A) := \begin{cases} \kappa^i(pa^i; A) & \mathsf{D}^i(y) = * \\ \delta_{\mathsf{D}^i(y)}(A) & \mathsf{D}^i(y) \neq * \end{cases} \tag{65}$$

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 $\kappa_B^{i,*} = \kappa_i'(\cdot; B)$.

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

$$[\kappa_B^{i,*}]^{-1}(A) = [\mathsf{D}^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \qquad \text{if } 0, 1 \notin A \tag{66}$$

$$= [\mathsf{D}^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [\mathsf{D}^i]^{-1}(B) \times X^{\mathrm{Pa}_\mathcal{G}(i)} \qquad \text{if } 1 \in A \wedge 0 \notin A \tag{67}$$

$$= [\mathsf{D}^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [\mathsf{D}^i]^{-1}(B^C) \times X^{\mathrm{Pa}_\mathcal{G}(i)} \qquad \text{if } 0 \in A \wedge 1 \notin A \tag{68}$$

$$= [\mathsf{D}^i]^{-1}(\{*\}) \times [\kappa_i^B]^{-1}(A) \cup [\mathsf{D}^i]^{-1}(X^i) \times X^{\mathrm{Pa}_\mathcal{G}(i)} \qquad \text{if } 0 \in A \wedge 1 \in A \tag{69}$$

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

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

$$\iota^\mathcal{G} : (y; A) \mapsto \int_{A^0} \kappa^{0,*}(y; dx^0) ... \int_{A^{N-1}} \kappa^{N-1,*}(y, x^{n-2}; dx^{n-1}) \kappa^{N,*}(y, x^{n-1}; A^N) \tag{70}$$

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

From Equations 64, 65 and 70 we can verify that, given some $i \in N$, if $\mathsf{D}^i(y) = \{*\}$ then $[\delta_y \iota^\mathcal{G}]_{\mathrm{Pa}_\mathcal{G}(\mathsf{X}^i)} = \kappa_i = \mu_{\mathrm{Pa}_\mathcal{G}(\mathsf{X}^i)} F_{\mathsf{X}^i}$ and if $\mathsf{D}^i(y) \neq \{*\}$ then $\delta_y \iota^\mathcal{G} = \delta_{\mathsf{D}^i(y)} F_{\mathsf{X}^i}$. From Equation 70 and the compatibility of $\mu$ 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 $\iota^\mathcal{G} = \kappa^{\mu, \mathcal{G}}$. $\qquad \square$

# 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.6 and B.7 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 $\langle (\mathcal{T}, E), D, \mathsf{X}, (U, E) \rangle$ where the observed data $\mathsf{X}$ is from a randomised controlled trial (RCT), $\mathsf{Y}_0 : E \to Y$ and $\mathsf{Y}_1 : E \to Y$ are random variables representing a particular outcome of interest under no treatment and treatment respectively and $\mathsf{Y} : E \to Y$ represents the "realised" outcome of interest and for $\xi \in \Delta(\mathcal{E})$, $U(\xi) = \mathbb{E}_\xi[\mathsf{Y}]$.

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 $\mathsf{X}$) to determine $\mathbb{E}_\mu[\mathsf{Y}_0]$ and $\mathbb{E}_\mu[\mathsf{Y}_1]$ [Rubin, 2005].

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

1. Assume performance bias is absent, so the theory must satisfy $\delta_i \kappa \mathsf{Y}_i = \mathbb{E}_\mu[\mathsf{Y}_i]$

2. Assume performance bias has a uniform additive effect: the theory satisfies $\delta_i \kappa \mathsf{Y}_i = \mathbb{E}_\mu[\mathsf{Y}_i] + k$. In this case the average treatment effect can still be estimated from the data: $\delta_1 \kappa \mathsf{Y}_1 - \delta_0 \kappa \mathsf{Y}_0 = \mathbb{E}[\mathsf{Y}_1] - \mathbb{E}[\mathsf{Y}_0]$ 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 $\mathbb{E}_\mu[\mathsf{Y}_i]$ and $\delta_i \kappa \mathsf{Y}_i$ 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"