

CSDT

October 30, 2019

1 Invariance and Capital-C Causality

CSDT features *consequences* - that is, probabilistic relations between decisions and results - but it does not feature *causal effects*, which seem to be probabilistic relations between random variables on the observation space E that are not necessarily disintegrations of a joint distribution. Here I propose the following notion of a causal effect in CSDT: if I have prior knowledge about the consequences of my decisions D on some random variable B via $C_0 B$ and I can extend this to the consequences on Y using some Markov kernel $G_\theta : B \rightarrow \Delta(\mathcal{Y})$ via composition $C_\theta Y = C_0 B G_\theta$ then we'll say that G_θ is the “causal effect” of B on Y in state θ .

This notion accords with intuitions about do-interventions - if I have the option to $do(B)$ then I definitely know the effect of my decision on B . Furthermore, using do interventions I assume that the effect of $do(B = x)$ on Y is given by fixing B to x and then computing $P(Y|do(B) = x)$, a special case of the composition above if we let $G_\theta = P(Y|do(B))$.

However, the notion of causal effect proposed here doesn't make additional commitments that $do(B)$ does - namely, that I also know $do(B)$ has no direct effect on any other variable. Unlike the assumption of prior knowledge, the “no direct effects” assumption is cannot even be formulated within CSDT.

The notion of “causal effect” given here permits a (to my knowledge) novel set of assumptions under which this type of causal effect of B on Y can be inferred. Importantly, it avoids any assumptions of “conditional independence of unobservables”, and is to my knowledge the only known case of “causal identifiability” that achieves this.

Suppose $C : D \times \Theta \rightarrow \Delta(\mathcal{E})$ is the consequence of interest, and furthermore given some $\theta \in \Theta$ we have $A_\theta : D \rightarrow \Delta(\mathcal{A})$ and $B : E \rightarrow \Delta(\mathcal{B})$ such that the observations are distributed according to

$$H_\theta := \begin{array}{c} \begin{array}{c} \boxed{A_\theta} \text{---} A \\ \boxed{C_\theta} \text{---} Y \\ \boxed{B} \text{---} B \end{array} \\ \swarrow \gamma_\theta \end{array} \quad (1)$$

That is, the observations \mathbf{H}_θ are the inputs and outputs of \mathbf{C}_θ where the inputs are masked by \mathbf{A}_θ and we run one copy of the outputs through a fixed \mathbf{B} . Without the mask, we could recover \mathbf{C}_θ from \mathbf{H}_θ via disintegration $D \dashrightarrow Y$.

For all θ , assume \mathbf{A}_θ is γ_θ -almost surely right invertible and that γ_θ is strictly positive. Together these are strong assumptions. Note that as θ is unknown, despite the fact that \mathbf{A}_θ is right invertible, we can't recover it's inputs, so we can't simply disintegrate.

This setup is something like an instrumental variables (IV) setup where γ_θ is a source of "exogenous" variation. It is stronger than a typical IV setup in that γ_θ is strictly positive and \mathbf{A}_θ is right invertible, but it is weaker than a typical IV setup in that we make no assumptions at all about model class and fewer assumptions about the relationships between \mathbf{A} , \mathbf{B} and \mathbf{Y} .

It is also somewhat similar to Arjovsky et al. (2019), itself based on Peters et al. (2016), if \mathbf{A} is understood as an environment indicator. Under this interpretation, the assumption that γ_θ is strictly positive is much stronger than Arjovsky et al. (2019), but on the other hand we do not assume any particular relationship between \mathbf{B} and \mathbf{Y} where Arjovsky et al. (2019) assumes a particular form of SEM. An interpretational difference is that while we regard D as a set of feasible decisions, Arjovsky et al. (2019) considers the set of environments thus:

Here, the set of all environments contains all possible experimental conditions concerning our system of variables, both observable and hypothetical.

As an aside, we could potentially interpret such a set of environments as an rich set of decisions which may be coarsened to a realistic set of decisions.

Under certain conditions, this setup allows for the extension of causal knowledge via observed data. In particular, if we find the conditional probability of \mathbf{Y} on \mathbf{A} and \mathbf{B} (written $[Y|\mathbf{A} \otimes \mathbf{B}]_\theta$) is independent of \mathbf{A} , then given prior knowledge for the effect of a decision on \mathbf{B} we can deduce the full consequence map \mathbf{C}_θ from our prior knowledge and the disintegration $[Y|\mathbf{B}]_\theta$ (see Theorem 1.2). Note that no assumptions have been made about "causal" relationships between \mathbf{B} and \mathbf{Y} - $[Y|\mathbf{B}]_\theta$ is an ordinary disintegration, not a platonic Markov kernel/structural equation/FFRCISTG/whatever.

A kernel \mathbf{B} that throws away more information is advantageous in the sense that less prior knowledge is needed to determine \mathbf{C}_θ .

Example 1.1 (Waste collection). A council is deciding on how to implement a waste collection service to reduce littering \mathbf{Y} . For every possible service $d \in D$, the collection schedule \mathbf{S} that will be achieved is known prior to any investigation (services may differ in other ways - e.g. the bin types may differ, and these differences may or may not be known in advance). In addition, the council has obtained weekly collection and littering data from a set A of other councils with their own waste collection services that are known to have faced the same unknown consequence map \mathbf{C}_θ . Each other council has implemented exactly one possible service $d \in D$ and enough councils were surveyed that all possible choices of service have been sampled, though they do not know which councils

have implemented which systems. These conditions ensure that for each service d the set of councils $A_d \subset A$ implementing that service is disjoint from the set of councils implementing any other service, and hence the unknown map from services to councils $\mathbf{A}_\theta : D \rightarrow \Delta(\mathcal{A})$ is right invertible.

It is found by the council's statisticians that the disintegration $[Y|A \otimes S]_\theta$ is independent of A . Thus by theorem 1.2 the impact of any service d on the rate of littering Y can be found via the collection schedule S that will be achieved by that service and the disintegration $[Y|S]_\theta$.

This is a surprisingly strong conclusion from the assumptions made. We appear to have the ability to deduce a “causal effect” from observational data in a context that looks remarkably similar to standard examples of when this *can't* be done. In fact, the conclusion is stronger than a *do*-style causal relationship, as *do* interventions assume we know a decision has no “direct effect” on any variable other than the target, whereas here we only assume the consequence of d on S is known.

An objection might be that other councils' waste collection service choices are determined, in part, by some unobserved background factors. Suppose that these background factors take values in some space K . One means of formalising this is that \mathbf{C}_θ factorises:

$$D - \boxed{\mathbf{C}_\theta} - E = \begin{array}{c} D \\ \swarrow \kappa \\ \boxed{\mathbf{F}_\theta} - E \end{array} \quad (2)$$

Where $\kappa \in \Delta(K)$ is a distribution on background factors and $\mathbf{F}_\theta : D \times K \rightarrow \Delta(\mathcal{E})$ maps decisions and background factors to consequences. For $d \neq d'$ we may have:

$$\begin{array}{c} \triangle \delta_d \\ \swarrow \\ \boxed{\mathbf{F}_\theta} - \boxed{\mathbf{S}} - S \\ \swarrow \\ K \end{array} \neq \begin{array}{c} \triangle \delta_{d'} \\ \swarrow \\ \boxed{\mathbf{F}_\theta} - \boxed{\mathbf{S}} - S \\ \swarrow \\ K \end{array} \quad (3)$$

That is, different decisions may induce different relationships between background characteristics K and collection schedules S . Formally, introducing this extra complication has not violated any of our assumptions - the causal inference is still valid! Practically, we would typically need a much larger set of decisions D in this case to account for the number of plausible relationships between K and S , and this may give us more reason to question the assumption that γ_θ is strictly positive - i.e. that the set of councils implementing each decision has positive measure. Nonetheless, this appears to be quite different to existing conditions for observational causal inference: we have allowed for K to affect both S and Y but in contrast to existing approaches **we do not need to see K in order to – sometimes – infer the “causal effect” of S on Y .**

Theorem 1.2. *Suppose we have a causal theory $\mathbf{T} : \Theta \times D \rightarrow \Delta(\mathcal{E}^2)$ where for $\theta \in \Theta$ we have consequence \mathbf{C}_θ and experiment \mathbf{H}_θ given by 1. Suppose for some $\mathbf{C}_0 : D \rightarrow \Delta(\mathcal{E})$ we have $\mathbf{C}_\theta \mathbf{B} = \mathbf{C}_0 \mathbf{B}$ for all $\theta \in \Theta$.*

For all $\theta \in \Theta$ such that $[Y|A \otimes B]_\theta$ is independent of A we have $C_\theta = C_0 B[Y|B]_\theta$

Proof. By the definition of disintegration

$$\begin{array}{c} \text{Diagram 1: } \gamma_\theta \text{ (triangle) with input } A \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } A_\theta \text{ is connected to } A \text{ and } Y. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \end{array} = \begin{array}{c} \text{Diagram 2: } \gamma_\theta \text{ (triangle) with input } A \text{ and } B. \text{ It has two outputs: } A \text{ and } B. \text{ A box } A_\theta \text{ is connected to } A \text{ and } Y. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \text{ A box } [Y|A \otimes B]_\theta \text{ is connected to } A \text{ and } Y. \end{array} \quad (4)$$

$$= \begin{array}{c} \text{Diagram 3: } \gamma_\theta \text{ (triangle) with input } A \text{ and } B. \text{ It has two outputs: } A \text{ and } B. \text{ A box } A_\theta \text{ is connected to } A \text{ and } Y. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \text{ A box } [Y|B]_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (5)$$

Where 5 follows from the independence of $[Y|A \otimes B]_\theta$ from A .

In addition, we have

$$\begin{array}{c} \text{Diagram 4: } \gamma_\theta \text{ (triangle) with input } D \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \end{array} = \begin{array}{c} \text{Diagram 5: } \gamma_\theta \text{ (triangle) with input } D \text{ and } B. \text{ It has two outputs: } D \text{ and } B. \text{ A box } A_\theta \text{ is connected to } D \text{ and } Y. \text{ A box } A_\theta^{-1} \text{ is connected to } D \text{ and } Y. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (6)$$

Thus

$$\begin{array}{c} \text{Diagram 6: } \gamma_\theta \text{ (triangle) with input } D \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \end{array} = \begin{array}{c} \text{Diagram 7: } \gamma_\theta \text{ (triangle) with input } D \text{ and } B. \text{ It has two outputs: } D \text{ and } B. \text{ A box } C_\theta \text{ is connected to } Y \text{ and } B. \text{ A box } [Y|B]_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (7)$$

From Lemma 1.3 and by the assumption of strict positivity on γ_θ , we therefore have

$$\begin{array}{c} \text{Diagram 8: } C_\theta \text{ (box) with input } Y \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \end{array} = \begin{array}{c} \text{Diagram 9: } C_\theta \text{ (box) with input } Y \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } [Y|B]_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (8)$$

$$\begin{array}{c} \text{Diagram 10: } C_\theta \text{ (box) with input } Y \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \end{array} = \begin{array}{c} \text{Diagram 11: } C_\theta \text{ (box) with input } Y \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } [Y|B]_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (9)$$

$$= \begin{array}{c} \text{Diagram 12: } C_0 \text{ (box) with input } Y \text{ and } B. \text{ It has two outputs: } Y \text{ and } B. \text{ A box } [Y|B]_\theta \text{ is connected to } Y \text{ and } B. \end{array} \quad (10)$$

Where 9 follows from marginalisation of 8.

□

A non-invertible A_θ means 6 doesn't hold. Given that A_θ is arbitrary apart from the fact that it is invertible, I wonder if a channel capacity lower boundon A_θ could give an upper bound on the “distance” between the kernels on the left and right of 6 using an appropriate replacement for A_θ^{-1} .

Lemma 1.3. *Given strictly positive probability measure $\gamma \in \Delta(\mathcal{D})$ and Markov kernels $\mathbf{X} : \mathcal{D} \rightarrow \Delta(\mathcal{E})$ and $\mathbf{Y} : \mathcal{D} \rightarrow \Delta(\mathcal{E})$ if*

$$\begin{array}{c} \text{D} \\ \curvearrowright \\ \triangleleft \gamma \text{---} \boxed{\mathbf{X}} \text{---} \mathcal{E} \end{array} = \begin{array}{c} \text{D} \\ \curvearrowright \\ \triangleleft \gamma \text{---} \boxed{\mathbf{Y}} \text{---} \mathcal{E} \end{array} \quad (11)$$

Then $\mathbf{X} = \mathbf{Y}$.

Proof. We note that both \mathbf{X} and \mathbf{Y} are $\mathcal{D} \dashrightarrow \mathcal{E}$ disintegrations of 11, and so they must be γ -almost surely equal. Strict positivity means they must in fact be equal. \square

References

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant Risk Minimization. *arXiv:1907.02893 [cs, stat]*, July 2019. URL <http://arxiv.org/abs/1907.02893>. arXiv: 1907.02893.

Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5): 947–1012, 2016. ISSN 1467-9868. doi: 10.1111/rssb.12167. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12167>.

Appendix: