
Answering Complex Causal Queries With the Maximum Causal Set Effect

Anonymous Author(s)

Affiliation

Address

email

Abstract

The standard tools of causal inference have been developed to answer simple causal queries which can be easily formalized as a small number of statistical estimands in the context of a particular structural causal model (SCM); however, scientific theories often make diffuse predictions about a large number of causal variables. This article proposes a framework for parameterizing such complex causal queries as the maximum difference in causal effects associated with two sets of causal variables of a researcher specified size. We term this estimand the *Maximum Causal Set Effect* (MCSE) and develop an estimator for it that is asymptotically consistent and conservative in finite samples under assumptions that are standard in the causal inference literature. This estimator is also asymptotically normal and amenable to the non-parametric bootstrap, facilitating classical statistical inference about this novel estimand. We compare this estimator to more common latent variable approaches and find that it can uncover larger causal effects in both real world and simulated data.

1 Introduction

Recent advances in machine learning technology have made it possible to non-parametrically estimate many parameters present in complex structural causal models (SCMs). Specifically, such estimating technology has rapidly advanced for three major causal inference settings: the many causes setting, the many moderators setting, and the many mediators setting. All three settings represent a situation in which a particular causal query can be stated in terms of a large number of combinations of different variables. Specifically, a researcher could estimate a different treatment effect associated with each of the many different possible combinations of causes [Imbens, 2000, Wang and Blei, 2019, Li et al., 2019, Wang et al., 2018, Zheng et al., Forthcoming], a different conditional treatment effect for each of the many different combinations of moderators [Green and Kern, 2012, Athey and Imbens, 2016, Grimmer et al., 2017, Wager and Athey, 2018, Künzel et al., 2019], and a different mediated effect for each of the many different combinations of mediators [Zhou and Yamamoto, 2020, Daniel et al., 2015]. Such causal queries are complex in the sense that they require summarizing the combined influence of a large number of causal variables.

The main challenge for applied researchers in such settings is that standard causal inference algorithms are designed to provide a different estimate associated with each of the many causal variables rather than a single number summarizing the combined influence of all the causal variables together. Consider, for example, the setting of inferring the causal effect of actors on a film’s box office performance. Wang and Blei [2019] provide a framework for estimating the average treatment effect associated with every actor on a film’s performance. While certainly useful for making predictions

about which actors a director should cast, an economist studying the film industry might prefer a single number which summarizes the general importance of actors in general for a film’s box office success. As discussed in the next section, such settings are common in scientific research, suggesting the need for novel causal estimands to parameterize the predictions of such theories in the context of a particular SCM.

Contribution The contribution of this paper is threefold. First, it introduces the notion of a complex causal query and argues that existing causal estimands are of limited utility to applied researchers in the face of such queries. Second, it defines a novel estimand – the *Maximum Causal Set Effect* (MCSE) – which can be used to provide an interpretable answer to such complex queries. Finally, the paper introduces an estimator for this estimand. The estimator is based on techniques proposed in the double Q-learning literature [Hasselt, 2010] and is asymptotically consistent and conservative in finite samples under assumptions that are standard in the causal inference literature. It is also asymptotically normal and amenable to non-parametric bootstrap techniques, facilitating classical statistical inference about the MCSE.

2 Setting and Previous Work

2.1 Problem Overview

Standard approaches to causal inference [Pearl, 2009] typically begin with the researcher specifying an SCM and then defining a causal query which can be answered based on the assumed SCM. Under certain assumptions about the SCM, it may be possible to estimate the answer to that causal query using the conventional tools of statistical inference. The standard tools of causal inference are designed with settings in mind where the predictions of a scientific theory take the form of a *simple causal query*. Such queries are stated in terms of some low dimensional *causal variable* t and some outcome Y . For example, a question like how much does a medical procedure reduce the risk of disease, represents a simple causal query because it is defined in terms of a single unidimensional treatment. Such queries can be easily quantified using conventional statistical estimands because they are directly formulated in terms of a small number of theoretically motivated variables.

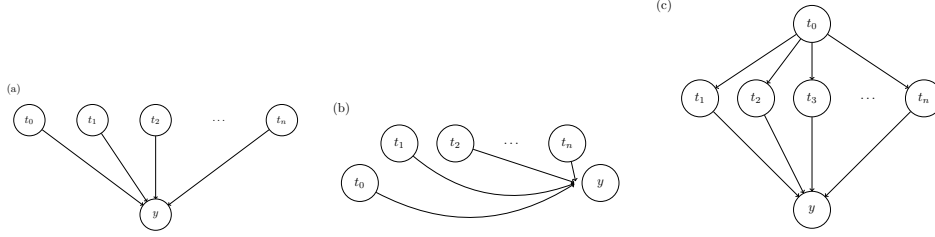
This paper instead focuses on situations where a scientific theory makes diffuse predictions about the importance of a large number of causal variables, defying the stylization of simple causal queries. Such queries are common in scientific research. For example:

- **Genome Wide Association Studies (GWAS)** – GWAS attempt to quantify the causal effect of a huge number of individual genotypes on the likelihood that some trait is expressed [Stephens and Balding, 2009, Visscher et al., 2017]
- **Personality** – psychologists are often interested in the effect certain personality traits (such as extraversion or neuroticism) might have on life outcomes [Pervin, 2003], but such traits are only observed by the researcher as responses to a large number of survey questions.
- **Text** – language is complex and multi-faceted and the causal effect of the wording of a document on a user’s response requires an assessment of the contribution of many different topics or words together [Fong and Grimmer, 2016, Egami et al., 2018, Fong and Grimmer].
- **Complex medical treatments** – many medical treatments cannot be reduced to a single low dimensional representation. For example, radiation exposure is observed as a high dimensional vector [Nabi et al., 2017] and medical researchers might also wish to understand the combined importance of many procedures using electronic medical records [Gottesman et al., 2013].

Such causal queries are *complex* because they require estimating the joint influence of many causal variables.

The SCM undergirding such complex queries can take many forms. Three major examples are: (a) the many causes setting where the researcher wishes to understand the joint influence of many

Figure 1: Visualization of Causal Graphs With Complex Queries



Note: Figure visualizes SCMs corresponding to complex causal queries in the case of (a) many causes, (b) many moderators, and (c) many mediators. (a) visualizes the case where treatment types $t_0 \dots t_n$ each influence y described in Wang and Blei [2019]. (b) visualizes the case where the causal effect of t_0 on y is directly modified by $t_1 \dots t_n$ as described in VanderWeele and Robins [2007]. (c) visualizes the case where the effect of t_0 on y is mediated by $t_1 \dots t_n$ as described by Zhou and Yamamoto [2020].

82 treatments (b) many moderators setting where the researcher wishes to understand how effect of a
 83 binary treatment varies based on many variables (c) the many mediators setting where the researcher
 84 wishes to model how a causal effect can be decomposed into many different pieces. These SCM's are
 85 visualized in Figure 1 in the form of directed acyclical graphs (DAGs). The unifying trait of a complex
 86 causal query is that it asks about the importance of many arrows present in each DAG.

87 Techniques developed in the context of simple causal queries cannot be readily used to answer
 88 complex ones. While the standard tools of causal inference can be used to estimate causal effects
 89 corresponding to every combination of causal variables in SCM's like those visualized in Figure 1,
 90 they do not provide applied researchers with a single unambiguous estimate with which to summarize
 91 the joint causal effect of many such variables.

92 2.2 Previous Work

93 The only existent proposal for addressing the challenge presented by complex causal queries in the
 94 machine learning literature is to dimension reduce the relevant causal variables and then focus on a
 95 simple causal query defined in terms of that latent trait [Fong and Grimmer, 2016, Fong and Grimmer,
 96 Nabi et al., 2017]. This strategy has only been proposed in the many causes setting, but could also
 97 be extended to the many moderators or many mediators cases as well. Such a strategy is inherently
 98 reductive and risks understating the magnitude of causal effects because it disregards all variation in
 99 the treatment types that is not accounted for in the latent trait. Additionally such latent traits are often
 100 scale invariant and so may lack a scientifically meaningful interpretation.

101 2.3 Assumptions and Notation

102 We assume that the researcher observes a set of N independent (t_i, Y_i, \mathbf{x}_i) triplets where Y_i is the
 103 outcome, and t_i is a length K vector indicating the *treatment type* received by unit i , and \mathbf{x}_i is a
 104 length J vector representing a set of background covariates that causal effects should be adjusted for.
 105 Additionally, let \mathcal{T} denote the support of the distribution of t_i .

106 We also assume that the researcher has knowledge of the population distribution of t_i : $g(t)$. In many
 107 settings, the empirical distribution of t_i will be the most logical choice, but other choices may be
 108 reasonable as well if the population distribution is known to the researcher, as might be the case when
 109 conducting survey research or if the treatment types were experimentally randomized.

110 Finally, we assume that the researcher has specified some SCM and has specified a simple causal
 111 query, $\tau(\mathcal{T}', \mathcal{T}'')$, which is defined in terms of two subsets: $\mathcal{T}', \mathcal{T}'' \subseteq \mathcal{T}$. In the many causes case,
 112 $\tau(\mathcal{T}', \mathcal{T}'')$ might take the form:

$$\tau(\mathcal{T}', \mathcal{T}'') \equiv \mathbb{E}(\mathbb{E}(Y_i | \text{do}(t)) | t \in \mathcal{T}') - \mathbb{E}(\mathbb{E}(Y_i | \text{do}(t)) | t \in \mathcal{T}'')$$

where $\text{do}(\cdot)$ represents some causal intervention [Pearl, 2009]. This estimand represents the average effect of receiving a set of treatments contained in \mathcal{T}' rather than \mathcal{T}'' .

In the many moderators or many mediators case on the other hand let the zeroth element of the treatment types vector received by unit i , $\mathbf{t}_{i,0} \in \{0, 1\}$, denote the level of some binary treatment received by unit i . Similarly, let the remaining elements of \mathbf{t}_i be denoted $\mathbf{t}_{i,-0}$ and indicate the level received by unit i on the many moderators or mediators. Then a possible choice for $\tau(\mathcal{T}', \mathcal{T}'')$ might be:

$$\begin{aligned} \tau(\mathcal{T}', \mathcal{T}'') &\equiv \mathbb{E}(\mathbb{E}(Y_i | \text{do}(\mathbf{t}_{i,0} = 1)) - \mathbb{E}(Y_i | \text{do}(\mathbf{t}_{i,0} = 0)) | \mathbf{t}_{i,-0} \in \mathcal{T}') \\ &\quad - \mathbb{E}(\mathbb{E}(Y_i | \text{do}(\mathbf{t}_{i,0} = 1)) - \mathbb{E}(Y_i | \text{do}(\mathbf{t}_{i,0} = 0)) | \mathbf{t}_{i,-0} \in \mathcal{T}'') \end{aligned}$$

which represents the difference in average treatment effects between units with moderators or mediators contained in \mathcal{T}' rather than \mathcal{T}'' .

While these two choices of $\tau(\mathcal{T}', \mathcal{T}'')$ are likely to be useful in a number of situations, the framework could easily be generalized to a much wider range of causal quantities of interest. For example, $\tau(\mathcal{T}', \mathcal{T}'')$ could easily be defined in terms of ratios of different average outcomes, outcome quantiles, or instrumental variables approaches, etc.

3 The Maximum Causal Set Effect

The challenge for applied researchers in the presence of such complex causal queries is that a different value of $\tau(\mathcal{T}', \mathcal{T}'')$ can be defined for every distinct pair of sets $\mathcal{T}', \mathcal{T}'' \subseteq \mathcal{T}$, leaving the analyst without a single unambiguous causal estimand to summarize their findings. In this section, we define a causal quantity of interest which overcomes this challenge by focusing on the contrast between two sets $\mathcal{T}_q^{\text{Max}}$ and $\mathcal{T}_q^{\text{Min}}$ which maximize $\tau(\mathcal{T}', \mathcal{T}'')$. To avoid choosing sets $\mathcal{T}_q^{\text{Max}}$ and $\mathcal{T}_q^{\text{Min}}$ which correspond to unrepresentative edge cases, we require that the sets be of a researcher specified size: q . Formally, let the set of subsets of \mathcal{T} such that the probability that \mathbf{t}_i is in \mathcal{T} is at least q be defined as: $\mathcal{T}_q \equiv \{\mathcal{T}' \subseteq \mathcal{T} : P(\mathbf{t}_i \in \mathcal{T}') \geq q\}$ where $P(\mathbf{t}_i \in \mathcal{T}') = \int_{\mathcal{T}} g(\mathbf{t}) \mathbb{1}\{\mathbf{t} \in \mathcal{T}'\} d\mathbf{t}$.

We then define MCSE_q as:

$$\text{MCSE}_q = \max_{\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q} \tau(\mathcal{T}', \mathcal{T}'') = \tau(\mathcal{T}_q^{\text{Max}}, \mathcal{T}_q^{\text{Min}})$$

We refer to $\mathcal{T}_q^{\text{Max}}$ as the *maximum causal set* and $\mathcal{T}_q^{\text{Min}}$ as the *minimum causal set*. For many applications, the MCSE will have an intuitive and scientifically meaningful interpretation. In the actors example, it might be used to answer a question like what is the expected difference in box office performance between a film cast with one of the 10% best performing casts rather than one of the bottom 10% worst performing casts? Similarly, in the genetics example, it might answer the question, what is the difference in the efficacy of some drug for patients with one of the top 10% most treatment enhancing sets of genes rather than one of the bottom 10% most treatment diminishing sets of genes?

4 Estimation

This section outlines an algorithm for estimating MCSE_q . Sample splitting is a major part of this algorithm and this section develops the procedure in the context of a single data split. The efficiency of this estimator can also easily be improved by rotating the roles that each subset of the data plays and then averaging the results, a procedure known as crossfitting [Chernozhukov et al., 2017], which we discuss in Appendix A.

4.1 Algorithm Overview

A basic result in the Q-learning literature is that a single sample estimator for the maximum expected value will have an upward bias. Since conservative estimators are easier to interpret and necessary for valid hypothesis testing, we follow the lead of Hasselt [2010] in using a split sample estimator for this estimation task. This approach is also useful in demonstrating the asymptotic normality of the resulting estimator as well.

Specifically, we begin by assuming that the analyst has randomly split the observations into two equally sized sets, \mathcal{S}^{Est} and $\mathcal{S}^{\text{Prob}}$. We further assume that the analyst has specified two models. The first uses the elements of the splitting set to make predictions about the probability that any $\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$ are the true maximum and minimum causal sets and we denote its predictions: $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$. The second model makes a prediction about $\tau(\mathcal{T}', \mathcal{T}'')$ for any two $\mathcal{T}', \mathcal{T}'' \subseteq \mathcal{T}$, and we denote its predictions $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$. Note $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ should make use only of outcomes that are included in $\mathcal{S}^{\text{Prob}}$ while $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ should only use the outcomes in \mathcal{S}^{Est} so that, $\forall \mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$, $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}}) \perp \hat{\tau}(\mathcal{T}', \mathcal{T}'')$ conditional on observing the sample values of \mathbf{t}_i and x_i for all units. After specifying models for $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ and $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$, estimation proceeds as a weighted average of the estimates for $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ for every $\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$:

$$\widehat{\text{MCSE}}_q = \sum_{\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q} \hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}}) \hat{\tau}(\mathcal{T}', \mathcal{T}'')$$

4.2 Point Estimation Properties

A major requirement for the good behavior of this estimator is that $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ obey the basic probability axioms and that it assign zero probability to sets of treatment types which are too small to be plausible candidates for $\mathcal{T}_q^{\text{Max}}$ and $\mathcal{T}_q^{\text{Min}}$. These requirements are entirely verifiable by the analyst through the careful construction of $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ and are formalized in the following assumption:

Assumption 1. $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ satisfies the following conditions:

- $\sum_{\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q} \hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}}) = 1$
- $\forall \mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q, 0 \leq \hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}}) \leq 1$
- $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}}) = 0$ for all $\mathcal{T}', \mathcal{T}'' \notin \mathcal{T}_q$

Under Assumption 1, $\widehat{\text{MCSE}}_q$ can be interpreted as a weighted average of estimators for the causal effect of being treated with a treatment type in one set rather than another. Because MCSE_q is defined as the maximum of such causal effects for any two subsets of \mathcal{T} of the required size, it will always be greater than the expectation of this average, leading to the following proposition:

Proposition 1. If $\forall \mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q, \mathbb{E}(\hat{\tau}(\mathcal{T}', \mathcal{T}'')) \leq \tau(\mathcal{T}', \mathcal{T}'')$ and the conditions of Assumption 1 hold, then:

$$\mathbb{E}(\widehat{\text{MCSE}}_q) \leq \text{MCSE}_q$$

Proof in appendix C.1

The conditions for finite sample conservatism are relatively mild (for example, $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ could be misspecified or inconsistent); however, as formalized in the next proposition, the conditions for the consistency of MCSE_q are a bit stronger and require that $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ converge to a binary indicator identifying $\mathcal{T}_q^{\text{Min}}$ and $\mathcal{T}_q^{\text{Max}}$:

Proposition 2. If $\forall \mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$,

$$\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}} | \mathcal{S}^{\text{Prob}}) \xrightarrow[n \rightarrow \infty]{p} \mathbb{1}\{\mathcal{T}' = \mathcal{T}_q^{\text{Max}}\} \mathbb{1}\{\mathcal{T}'' = \mathcal{T}_q^{\text{Max}}\}$$

and

$$\hat{\tau}(\mathcal{T}', \mathcal{T}'') \xrightarrow[n \rightarrow \infty]{p} \tau(\mathcal{T}', \mathcal{T}'')$$

then

$$\widehat{\text{MCSE}}_q \xrightarrow[n \rightarrow \infty]{p} \text{MCSE}_q$$

187 This result will also hold if convergence in probability is replaced with almost sure convergence.

188 Proof in Appendix C.2.

189 Many machine learning techniques (e.g. support vector machines, regression trees, etc.) will not
 190 readily produce probabilistic estimates for $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$, instead generating only a
 191 binary prediction for the two sets $\mathcal{T}_q^{\text{Min}}$ and $\mathcal{T}_q^{\text{Max}}$. The following proposition shows that such binary
 192 estimators will perform at best as well as probabilistic estimators as long as the two estimators have
 193 the same expectation:

Proposition 3. Let, $d(\mathcal{T}', \mathcal{T}'') \in \{0, 1\}$ and $w(\mathcal{T}', \mathcal{T}'') \in [0, 1]$ represent two choices for $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$. Let $\widehat{\text{MCSE}}_q^d$ and $\widehat{\text{MCSE}}_q^w$ represent the corresponding estimators for MCSE_q . Then if $\forall \mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$, $\mathbb{E}(d(\mathcal{T}', \mathcal{T}'')) = \mathbb{E}(w(\mathcal{T}', \mathcal{T}''))$,

$$\mathbb{E} \left(\left(\text{MCSE}_q - \widehat{\text{MCSE}}_q^w \right)^2 \right) \leq \mathbb{E} \left(\left(\text{MCSE}_q - \widehat{\text{MCSE}}_q^d \right)^2 \right)$$

194 Proof in Appendix C.3

195 A direct implication of this result is that bootstrap aggregation can be used to improve the performance
 196 of any binary predictor for $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ to create a probabilistic estimator without
 197 changing the expected value of the predictions.

198 4.3 Interval Estimation

199 While the previous section establishes the properties of the point estimator for MCSE_q , such results
 200 will be of little utility for applied researchers without a corresponding framework for measuring
 201 the uncertainty of those estimates. In this section, we begin the process of providing such results
 202 by introducing the assumption that $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ can be represented as a linear combination of the
 203 estimation set outcomes:

Assumption 2. Let $Z = \{\mathbf{t}_i, x_i : i \in \mathcal{S}^{\text{Est}}\}$. For any $\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q$ there exists a set of transformations $\{f_i(Z, \mathcal{T}', \mathcal{T}'') : i \in \mathcal{S}^{\text{Est}}\}$ such that:

$$\hat{\tau}(\mathcal{T}', \mathcal{T}'') = \sum_{i \in \mathcal{S}^{\text{Est}}} f_i(Z, \mathcal{T}', \mathcal{T}'') Y_i$$

204 Many common estimators for causal effects (e.g. matching, weighting, regression techniques, etc) fit
 205 this form, so such an assumption will not be unduly restrictive in many settings.

206 This assumption eases the derivation of asymptotic normality because it shows that $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ can be
 207 represented as the sum of independent random variables. The following proposition uses the central
 208 limit theorem derived by Neumann [2013] to show that multiplication by $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$
 209 will not impact this convergence so that asymptotic normality of $\widehat{\text{MCSE}}_q$ can be preserved under
 210 some mild regularity conditions:

211 **Proposition 4.** If $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ satisfies assumption 1; $\forall i, \mathbb{E}(Y_i^2) < \infty$; and
 212 $\forall \epsilon > 0$,

$$\sum_{i \in \mathcal{S}^{\text{Est}}} \frac{1}{|\mathcal{S}^{\text{Est}}|} \mathbb{E} \left(f_i(Z, \mathcal{T}', \mathcal{T}'')^2 Y_i^2 \mathbb{1}\{|f_i(Z, \mathcal{T}', \mathcal{T}'')| > \epsilon\} \right) \xrightarrow{|\mathcal{S}^{\text{Est}}| \rightarrow \infty} 0$$

Then, conditional on observing the estimation set values of \mathbf{t}_i and \mathbf{x}_i ,

$$\frac{(\widehat{MCSE}_q - \mathbb{E}(\widehat{MCSE}_q))}{\sqrt{\widehat{Var}(\widehat{MCSE}_q)}} \xrightarrow{D} \mathcal{N}(0, 1)$$

213 Proof in Appendix C.4

214 The final result necessary for conducting classical statistical inference is a corresponding variance
 215 estimator. This can be most easily accomplished via the non-parametric bootstrap. Specifically,
 216 Mammen [1992] shows that the non-parametric bootstrap is consistent for an asymptotically normal
 217 estimator that can be represented as a linear transformation of some set of independent observations.
 218 The following lemma uses assumption 2 to provide just such a result:

Lemma 1.

$$\widehat{MCSE}_q = \sum_{i \in \mathcal{S}^{Est}} Y_i w_i$$

219 where $w_i = \sum_{\mathcal{T}', \mathcal{T}'' \in \mathcal{T}_q} \hat{P}(\mathcal{T}' = \mathcal{T}_q^{Max} \cap \mathcal{T}'' = \mathcal{T}_q^{Min}) f_i(Z, \mathcal{T}', \mathcal{T}'')$

220 *Proof.* The proof follows trivially by using assumption 2 to substitute $\sum_{i \in \mathcal{S}^{Est}} f_i(Z, \mathcal{T}', \mathcal{T}'')$ for
 221 $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ in the definition of \widehat{MCSE}_q and then changing the order of summation. \square

222 So the variance and confidence intervals of \widehat{MCSE}_q can be consistently estimated by bootstrap
 223 resampling from the set $\{Y_i w_i : i \in \mathcal{S}^{Est}\}$.¹

224 5 Experiments

225 5.1 Benchmarks on Synthetic Data

226 We first consider the performance of this estimation procedure using synthetic data. Specifically, to
 227 assess the performance of this estimator, we implemented it on synthetic version of the many causes
 228 setting. First, we generated a set of N length K vectors of causes for each unit i as $\mathbf{t}_i \sim \mathcal{N}(0, \Sigma)$
 229 where Σ is some matrix with ones on the diagonal elements and some value $\rho \in [0, 1]$ in the off
 230 diagonal elements. We then generated the outcome as $\mu_i = \mathbf{t}_i' \beta$ where β is a length K vector
 231 composed of i.i.d draws from the standard normal distribution. Finally, we normalized μ_i so that the
 232 corresponding value of \widehat{MCSE}_q was always 1 and generated the outcome variables as $Y_i = \mu_i + \epsilon_i$
 233 where $\epsilon_i \sim \mathcal{N}(0, 1)$.

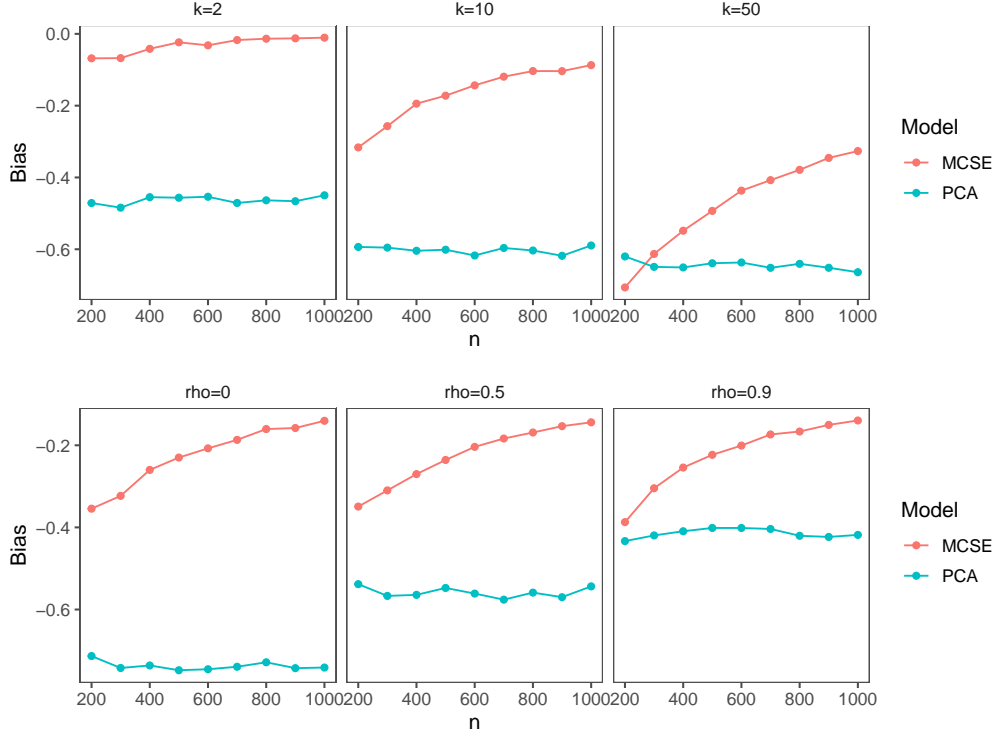
234 We implemented two estimators on this dataset. The first is the split sample \widehat{MCSE}_q estimator
 235 described in this paper². Note that under this simulation set up, all the assumptions needed for the
 236 theoretical results presented in Section 4 to hold are known to be true, so \widehat{MCSE}_q should be unbiased
 237 and consistent. We compared the performance of \widehat{MCSE}_q with an estimate for \widehat{MCSE}_q generated
 238 using a linear regression of Y_i on the first principal component of \mathbf{t}_i .³ This estimator corresponds to
 239 the current state of the art for drawing causal inferences in the face of a complex causal query, which
 240 involves using dimension reduction techniques to simplify the complex causal query into a simple
 241 one. We repeated this procedure 100 times for each combination of $K = 2, 10$, and 50 ; $\rho = 0, .5$ and
 242 $.9$; and values of N between 100 and 1,000.

¹Note, clustered standard errors can also be easily generated using the block bootstrap.

²Specifically, one using monte carlo sampling from the asymptotic distribution of linear regression of Y_i on \mathbf{t}_i as $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{Max} \cap \mathcal{T}'' = \mathcal{T}_q^{Min})$ and a linear regression for $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$. See Appendix B.1 for more details on the implementation of \widehat{MCSE}_q

³See appendix B.2 for details on the implementation of this estimator.

Figure 2: Simulation Results



Note: Red dots identify the bias of the method for quantifying the combined effect of many causes proposed in this paper while blue dots show the bias of dimension reduction techniques that represent the current state of the art for this same task.

Figure 2 visualizes the results of this analysis. Each point in the figure represents the average of all 300 iterations of the simulation procedure with the same values of n and K or n and ρ .⁴ Because the bias of both estimators is large relative to their variance in this setting, Figure 2 focuses on the bias of the estimators.⁵ These estimates show that \widehat{MCSE}_q is a large improvement over the latent trait model, generating significantly less biased estimates even when ρ is large and the principal components analysis (PCA) should perform well. Importantly, the bias of \widehat{MCSE}_q appears to vanish asymptotically while the PCA estimator shows little convergence as the sample size increases.

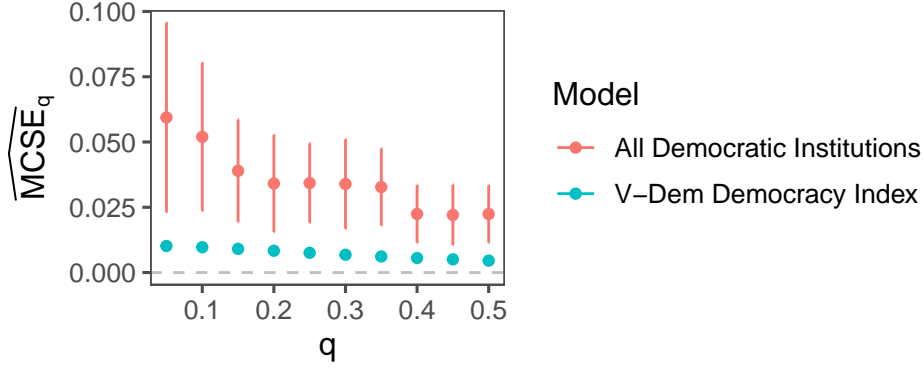
5.2 An Application to Real World Data

Our second application focuses on the role of democratic political institutions in reducing the likelihood of civil war onset. Democracy is a fundamental concept when modeling the quality of governance, but drawing inferences about its effect represents a straightforward example of the multiple causes setting. In particular, democracy cannot be measured as a single unambiguous feature – instead it is a confluence of many conceptually related by empirically distinct features describing different aspects of a system of governance. The causal effect of democracy on outcomes like conflict initiation is typically measured using a dimension reduction of the features representing the individual institutions [Treier and Jackman, 2008]; however, such strategies have led to conflicting results about the importance of democracy for political stability [Vreeland, 2008, Fearon and Laitin, 2003].

⁴Note, the monte carlo error in these estimates is quite low. The standard error associated with these average is never higher than .019 for any of the points.

⁵Appendix 4 presents estimates for the root mean squared error, which show a similar pattern

Figure 3: The Causal Effect of Democratic Political Institutions on the Probability of Civil War Onset



Note: The red dots identify estimates for the $MCSE_q$ made using the methodology outlined in this paper and represent the combined influence of many different democratic institutions together. The blue dots instead represent the influence of just a univariate latent trait produced by the maintainers of the V-Dem Dataset that is frequently used to model democracy.

Note 2: Confidence intervals adjusted for clustering by country.

260 Consequently, the role of democratic political institutions in reducing civil war onset represents a
 261 useful case for comparing latent trait models with with the MCSE.

262 Specifically, we used a linear model with 4 lagged outcomes and fixed effects for the country and year
 263 for both $\hat{P}(\mathcal{T}' = \mathcal{T}_q^{\text{Max}} \cap \mathcal{T}'' = \mathcal{T}_q^{\text{Min}})$ and $\hat{\tau}(\mathcal{T}', \mathcal{T}'')$ and measured democratic political institutions
 264 using the 128 features describing the system of governance present in a country in the Varieties of
 265 Democracy Dataset (V-Dem).⁶ The red dots and confidence intervals in Figure 3 show the estimates
 266 for $MCSE_q$ quantifying the effect of these political institutions on civil war onset for many different
 267 values of q . In particular, they suggest that countries with one of the 10% most conflict reducing
 268 institutions have roughly a 5% lower risk of civil war than countries with some of the 10% most
 269 conflict inducing institutions. The blue dots instead represent predictions for the $MCSE_q$ made using
 270 the predictions of a linear regression of just the V-Dem democracy variable on the probability of civil
 271 war onset.⁷ The estimates for MCSE are significantly larger than those generated using the more
 272 typical univariate model, suggesting that the the MCSE can successfully recover causal effects that
 273 standard latent variable approaches cannot.⁸

274 Conclusion

275 Non-parametric estimation techniques and high dimensional datasets increasingly confront re-
 276 searchers with estimates for a huge number of distinct causal estimands. While the capacity to
 277 fit such models represents tremendous progress for the estimation and computational techniques
 278 that support them, scientific theories rarely make predictions about such a large number of distinct
 279 parameters. In this article, we propose a framework for making sense of such model outputs by
 280 focusing on the maximum causal contrast between two sets of a researcher specified size q . We
 281 also develop an estimator for this estimand that is consistent, conservative in finite samples, and
 282 asymptotically normal. While the estimator is developed with the many causes and treatment effect
 283 heterogeneity settings in mind, the framework is extremely flexible and could be extended to a
 284 myriad of other causal qauntities of interest, speaking to its wide applicability and utility for applied
 285 researchers.

⁶See appendix B.1 for more details on these models.

⁷Specifically, we used the linear regression to impute the conditional expectation function, and then estimated the corresponding value of $MCSE_q$ using that imputed conditional expectation.

⁸While there is no straightforward way to generate confidence intervals for the univaraiate MCSE estimates, the coefficient from regressing civil war occurrence on the democracy variable is not statistically significant.

286 **Broader Impacts** Complex causal queries are ubiquitous in scientific research. While statistical
 287 analyses typically begin with the researcher specifying a small number of causal variables to focus
 288 on, scientific theories often make diffuse predictions about many variables working together. Such
 289 settings are particularly common in the social sciences, where causal variables often correspond to
 290 latent constructs that are only observed by the researcher as a set of proxies. For example, concepts
 291 like ideology, intelligence, or good public policy are not observed directly by the researcher, instead
 292 they are only revealed indirectly through a large number proxies such as votes cast in a legislature,
 293 answers to questions on an IQ test, or a large number of policies that may or may not be present in
 294 a particular municipality. Such complex causal queries also emerge in the natural sciences. Most
 295 prominently, genetics research is directly concerned with assessing the influence of a large number of
 296 genes on some outcome. Biological systems more generally often involve the complex interaction of
 297 a large number of distinct processes and could be understood from a similarly framework. These
 298 wide ranging examples speak to the value of the MCSE as an interpretable causal estimand for a wide
 299 range of applied researchers.

300 References

- 301 Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings*
 302 *of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- 303 Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, and Whitney
 304 Newey. Double/debiased/neyman machine learning of treatment effects. *American Economic*
 305 *Review*, 107(5):261–65, 2017.
- 306 Rhian M Daniel, Bianca L De Stavola, SN Cousens, and Stijn Vansteelandt. Causal mediation
 307 analysis with multiple mediators. *Biometrics*, 71(1):1–14, 2015.
- 308 Naoki Egami, Christian J Fong, Justin Grimmer, Margaret E Roberts, and Brandon M Stewart. How
 309 to make causal inferences using texts. *arXiv preprint arXiv:1802.02163*, 2018.
- 310 James D Fearon and David D Laitin. Ethnicity, insurgency, and civil war. *American political science*
 311 *review*, pages 75–90, 2003.
- 312 Christian Fong and Justin Grimmer. Causal inference with latent treatments. *American Journal of*
 313 *Political Science*. URL [https://www.dropbox.com/s/hxxyn5vtpjuyrw4/dexp_rev.pdf?](https://www.dropbox.com/s/hxxyn5vtpjuyrw4/dexp_rev.pdf?dl=0)
 314 [dl=0](https://www.dropbox.com/s/hxxyn5vtpjuyrw4/dexp_rev.pdf?dl=0).
- 315 Christian Fong and Justin Grimmer. Discovery of treatments from text corpora. In *Proceedings of*
 316 *the 54th annual meeting of the association for computational linguistics (volume 1: long papers)*,
 317 pages 1600–1609, 2016.
- 318 Omri Gottesman, Helena Kuivaniemi, Gerard Tromp, W Andrew Faucett, Rongling Li, Teri A
 319 Manolio, Saskia C Sanderson, Joseph Kannry, Randi Zinberg, Melissa A Basford, et al. The
 320 electronic medical records and genomics (emerge) network: past, present, and future. *Genetics in*
 321 *Medicine*, 15(10):761–771, 2013.
- 322 Donald P Green and Holger L Kern. Modeling heterogeneous treatment effects in survey experiments
 323 with bayesian additive regression trees. *Public opinion quarterly*, 76(3):491–511, 2012.
- 324 Justin Grimmer, Solomon Messing, and Sean J Westwood. Estimating heterogeneous treatment
 325 effects and the effects of heterogeneous treatments with ensemble methods. *Political Analysis*, 25
 326 (4):413–434, 2017.
- 327 Hado Hasselt. Double q-learning. *Advances in neural information processing systems*, 23:2613–2621,
 328 2010.
- 329 Guido W Imbens. The role of the propensity score in estimating dose-response functions. *Biometrika*,
 330 87(3):706–710, 2000.

331 Sören R Künzle, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heteroge-
332 neous treatment effects using machine learning. *Proceedings of the national academy of sciences*,
333 116(10):4156–4165, 2019.

334 Fan Li et al. Propensity score weighting for causal inference with multiple treatments. *The Annals of*
335 *Applied Statistics*, 13(4):2389–2415, 2019.

336 Enno Mammen. Bootstrap, wild bootstrap, and asymptotic normality. *Probability Theory and Related*
337 *Fields*, 93(4):439–455, 1992.

338 Razieh Nabi, Todd McNutt, and Ilya Shpitser. Semiparametric causal sufficient dimension reduction
339 of high dimensional treatments. *arXiv preprint arXiv:1710.06727*, 2017.

340 Michael H Neumann. A central limit theorem for triangular arrays of weakly dependent random
341 variables, with applications in statistics. *ESAIM: Probability and Statistics*, 17:120–134, 2013.

342 Judea Pearl. *Causality*. Cambridge university press, 2009.

343 Lawrence A Pervin. *The science of personality*. Oxford university press, 2003.

344 Matthew Stephens and David J Balding. Bayesian statistical methods for genetic association studies.
345 *Nature Reviews Genetics*, 10(10):681–690, 2009.

346 Shawn Treier and Simon Jackman. Democracy as a latent variable. *American Journal of Political*
347 *Science*, 52(1):201–217, 2008.

348 Tyler J VanderWeele and James M Robins. Four types of effect modification: a classification based
349 on directed acyclic graphs. *Epidemiology*, 18(5):561–568, 2007.

350 Peter M Visscher, Naomi R Wray, Qian Zhang, Pamela Sklar, Mark I McCarthy, Matthew A Brown,
351 and Jian Yang. 10 years of gwas discovery: biology, function, and translation. *The American*
352 *Journal of Human Genetics*, 101(1):5–22, 2017.

353 James Raymond Vreeland. The effect of political regime on civil war: Unpacking anocracy. *Journal*
354 *of conflict Resolution*, 52(3):401–425, 2008.

355 Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using
356 random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.

357 Yixin Wang and David M Blei. The blessings of multiple causes. *Journal of the American Statistical*
358 *Association*, pages 1–71, 2019.

359 Yixin Wang, Dawen Liang, Laurent Charlin, and David M Blei. The deconfounded recommender: A
360 causal inference approach to recommendation. *arXiv preprint arXiv:1808.06581*, 2018.

361 Jeffrey M Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.

362 Jiaping Zheng, Alexander D’Amour, and Alexander Franks. Copula-based sensitivity analysis for
363 multi-treatment causal inference with unobserved confounding. *Journal of the American Statistical*
364 *Association*, Forthcoming.

365 Xiang Zhou and Teppei Yamamoto. Tracing causal paths from experimental and observational data.
366 *SocArXiv. January*, 11, 2020.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **Yes**
- (b) Did you describe the limitations of your work? **Yes**. Limitations and requirements for all main results to hold are clearly stated in each theorem
- (c) Did you discuss any potential negative societal impacts of your work? **NA**. This work essentially generalizes causal inference techniques that have been developed in other contexts. While harms can certainly result from the improper use of such techniques, we do not perceive any additional problems emerging from the use of this estimand.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **Yes**

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **Yes**
- (b) Did you include complete proofs of all theoretical results? **Yes**

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **Yes**
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **Yes**
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **Yes**. The standard error of the bias estimates is discussed in footnote 4. They are quite low relative to the magnitude of performance improvement that our approach brings and should be un concerning.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **No**. These analyses are not computationally intensive and can be run locally on most computers.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **Yes**
- (b) Did you mention the license of the assets? **NA**
- (c) Did you include any new assets either in the supplemental material or as a URL? **NA**
- (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **NA**
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **NA**

- 413 5. If you used crowdsourcing or conducted research with human subjects...
- 414 (a) Did you include the full text of instructions given to participants and screenshots, if
- 415 applicable? NA
- 416 (b) Did you describe any potential participant risks, with links to Institutional Review
- 417 Board (IRB) approvals, if applicable? NA
- 418 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 419 spent on participant compensation? NA