# Disentangling Causal Effects from Sets of Interventions in the Presence of Unobserved Confounders

**Olivier Jeunen**[*]
University of Antwerp
Antwerp, BE

**Ciarán M. Gilligan-Lee**
Spotify Research
London, UK

**Rishabh Mehrotra**
Spotify Research
London, UK

**Mounia Lalmas**
Spotify Research
London, UK

## Abstract

The ability to answer causal questions is crucial in many domains, as causal inference allows one to understand the impact of interventions. In many applications, only a single intervention is possible at a given time. However, in certain important areas, multiple interventions are concurrently applied. Disentangling the effects of single interventions from jointly applied interventions is a challenging task—especially as simultaneously applied interventions can interact. This problem is made harder still by unobserved confounders, which influence both interventions and outcome. We address this challenge by aiming to learn the effect of a single-intervention from both observational data and sets of interventions. We prove that this is not generally possible, but provide identification proofs demonstrating that it can be achieved in certain classes of additive noise models—even in the presence of unobserved confounders. Importantly, we show how to incorporate observed covariates and learn heterogeneous treatment effects conditioned on them for single-interventions.

## 1 Introduction

The ability to answer causal questions is crucial in science, medicine, economics, and beyond, see Gilligan-Lee (2020) for a high-level overview. This is because causal inference allows one to understand the impact of interventions. In many applications, only a single intervention is possible at a given time, or interventions are applied one after another in a sequential manner. However, in some important areas, multiple interventions are concurrently applied. For instance, in medicine, patients that possess many commodities may have to be simultaneously treated with multiple prescriptions; in computational advertising, people may be targeted by multiple concurrent campaigns; and in dietetics, the nutritional content of meals can be considered a joint intervention from which we wish to learn the effects of individual nutritional components.

Disentangling the effects of single interventions from jointly applied interventions is a challenging task—especially as simultaneously applied interventions can interact, leading to consequences not seen when considering single interventions separately. This problem is made harder still by the possible presence of unobserved confounders, which influence both interventions and outcome. This paper addresses this challenge, by aiming to learn the effect of a single-intervention from both observational data and sets of interventions. We prove that this is not generally possible, but provide identification proofs demonstrating it can be achieved in certain classes of non-linear causal models with additive Gaussian noise—even in the presence of unobserved confounders. Importantly, we show

---

[*]Work was partially done during an internship at Spotify. Correspondence to `olivierjeunen@gmail.com`

how to incorporate observed covariates, which can be high-dimensional, by learning heterogeneous treatment effects conditioned on them for single-interventions.

Our main contributions are:

1. A proof that without restrictions on the causal model, single-intervention effects cannot be identified from observations and joint-interventions. (§3.1,3.2)

2. Proofs that single-interventions *can* be identified from observations and joint-interventions when the causal model belongs to certain classes of additive noise models. (§3.2,3.3)

3. An algorithm that learns the parameters of the proposed causal model and disentangles single interventions from joint interventions. (§4)

4. An empirical validation of our method on synthetic data. (§5)

## 2 Related Work

**Disentangling multiple concurrent interventions:** Parbhoo et al. (2021) study the question of disentangling multiple, simultaneously applied interventions from observational data. They propose a specially designed neural network for the problem and show good empirical performance on some datasets. However, they do not address the formal identification problem, nor do they address possible presence of unobserved confounders. By contrast our work derives the conditions under which identifiability holds. We moreover propose an algorithm that can disentangle multiple interventions even in the presence of unobserved confounders—as long as both observational and interventional data is available. Related work by Parbhoo et al. (2020) investigated the intervention-disentanglement problem from a reinforcement learning perspective, where each intervention combination constitutes a different action that a reinforcement learning agent can take. Unlike this approach, our work explicitly focuses on modelling the interactions between multiple interventions to learn their individual effects.

Closer to our work is Saengkyongam and Silva (2020), who investigate identifiability of joint effects from observations and single-intervention data. They prove this is not generally possible, but provide an identification proof for non-linear causal models with additive Gaussian noise. Our work addresses a complementary question; we want to learn the effect of a single-intervention from observational data and sets of interventions. Additionally, another difference between our work and that of Saengkyongam and Silva (2020) is that they do not consider identification of individual-level causal effects given observed covariates.

In a precursor to the work by Saengkyongam and Silva (2020), Nandy et al. (2017) developed a method to estimate the effect of joint interventions from observational data when the causal structure is unknown. This approach assumed linear causal models with Gaussian noise, and only proved identifiability in this case under a sparsity constraint. However, like Saengkyongam and Silva (2020), our result does not need the linearity assumption, and no sparsity constraints are required in our identification proof.

Finally, others including Schwab et al. (2020); Egami and Imai (2018); Lopez and Gutman (2017); Ghassami et al. (2021) explored how to estimate causal effects of a single categorical-, or continuous-valued treatment, where different intervention values can produce different outcomes. Unlike our work, they do not consider multiple concurrent interventions that can interact.

**Combining observations and interventions:** Bareinboim and Pearl (2016) have investigated non-parametric identifiability of causal effects using both observational and interventional data, in a paradigm they call "data fusion." More general results were studied by Lee et al. (2020), who provided necessary and sufficient graphical conditions for identifying causal effects from arbitrary combinations of observations and interventions. Recent work in Correa et al. (2021) explored identification of counterfactual—as opposed to interventional—distributions from combinations of observational and interventional data. Finally, Ilse et al. (2021) investigated the most efficient way to combine observational and interventional data to estimate causal effects. They demonstrated they could significantly reduce the number of interventional samples required to achieve a certain fit when adding sufficient observational training samples. However, they only prove their method theoretically in the linear-Gaussian case. In the non-linear case, they parameterise their model using normalising flows and demonstrate their method empirically. They only consider estimating single-interventions, and do not deal with multiple, interacting interventions.

**Additive noise models:** While certain causal quantities may not be generally identifiable from observational and interventional data, by imposing restrictions on the structural functions underlying causal models, one can obtain *semi-parametric* identifiability results. One of the most common restrictions are additive noise models (ANMs), first studied in the context of causal discovery by Hoyer et al. (2009). ANMs limit the form of the structural equations to be additive with respect to latent noise variables—but allow nonlinear interactions between causes. Janzing et al. (2009) used ANMs to devise a method for inferring a latent confounder between two observed variables. This is otherwise not possible without additional assumptions on the underlying causal model. ANMs have also been employed by Kilbertus et al. (2020) to investigate the sensitivity of counterfactual notions of fairness to the presence of unobserved confounding.

Our work proves that in certain classes of ANMs, the effect of a single-intervention can be identified from observational data and sets of interventions—even in the presence of unobserved confounders. Moreover, we show how to incorporate observed covariates in these ANMs to learn the heterogeneous effects of single-interventions conditioned on such covariates.

## 3 Identifiability of Single-Variable Interventional Effects

In this section, we provide identification proofs for single-variable interventional effects from observational data and joint interventions, for several model classes. Our theoretical analysis provides insights into the fundamental limitations of causal inference—and the assumptions that are required for identification.

**Problem Definition** We adopt the Structural Causal Model (SCM) framework as introduced by Pearl (2009). An SCM $\mathcal{M}$ is defined by $\langle \{C, X, Y\}, U, f, \mathsf{P}_U \rangle$, where $\{C, X, Y\}$ are endogenous variables separated into covariates $C$, treatments $X$, and the outcome $Y$, $U$ are exogenous variables (possibly confounders), $f$ are structural equations, and $\mathsf{P}_U$ defines a joint probability distribution over the exogenous variables.

The SCM $\mathcal{M}$ also induces a causal graph—where vertices represent endogenous variables, and edges represent structural equations. Vertices with outgoing edges to an endogenous variable $X_i$ are denoted as the parent set of this variable, or $\mathrm{PA}(X_i)$. Typically, the observed covariates $C$ causally influence the treatments as well as the outcome, and are a part of this set. Every endogenous variable $X_i$ (including $Y$) is then a function of its parents in the graph $\mathrm{PA}(X_i)$ and a latent noise term $U_i$, denoting the influence of factors external to the model:

$$X_i := f_i(\mathrm{PA}(X_i), U_i). \tag{1}$$

In *Markovian* SCMs, these latent noise terms are all mutually independent. However, in general, distinct noise terms can be correlated according to some global distribution $\mathsf{P}_U$. In this case, such correlation is due to the presence of *unobserved confounders*.

An intervention on variable $X_i$ is denoted by $\mathrm{do}(X_i = x_i)$, and it corresponds to replacing its structural equation with a constant, or removing all incoming edges in the causal graph. The core question we wish to answer in this work, is under which conditions the treatment effect of a single intervention can be disentangled from joint interventions and observational data. That is, given samples from the data regimes that induce

$$\mathbb{E}[Y|X_i = x_i, X_j = x_j, C = c], \text{ and } \mathbb{E}[Y|\mathrm{do}(X_i = x_i, X_j = x_j), C = c],$$

when can we learn conditional average causal effects

$$\mathbb{E}[Y|\mathrm{do}(X_i = x_i), X_j = x_j, C = c], \text{ or } \mathbb{E}[Y|X_i = x_i, \mathrm{do}(X_j = x_j), C = c]?$$

In what follows, we first show that this is not possible without restrictions on the causal model—a proof by counterexample. We then go on to prove, again by counterexample, that simply assuming ANMs without restrictions on the structure of the causal graph—the core assumption made by Saengkyongam and Silva (2020)—is not enough for this complementary research question. Finally, we prove identifiability of the treatment effect for ANMs without *causal* interactions between treatments. Note that this latter case does not mean treatments are independent: treatments can be influenced by observed covariates and unobserved confounders, and their interactions on the outcome are defined by the unrestricted function $f_Y$.

### 3.1 Unidentifiability for general SCMs

3

Table 2: Interventional distributions on $X_1$ under SCMs $\mathcal{M}$ and $\mathcal{M}'$ for 3.1.

| $\mathsf{P}_{\mathcal{M}}(Y, X_2 \mid \mathrm{do}(X_1))$ | | $Y = 0$ | $Y = 1$ |
|---|---|---|---|
| $\mathrm{do}(X_1 = 0)$ | $X_2 = 0$ | $1$ | $0$ |
| | $X_2 = 1$ | $0$ | $0$ |
| $\mathrm{do}(X_1 = 1)$ | $X_2 = 0$ | $1 - p$ | $0$ |
| | $X_2 = 1$ | $0$ | $p$ |

| $\mathsf{P}_{\mathcal{M}'}(Y, X_2 \mid \mathrm{do}(X_1))$ | | $Y = 0$ | $Y = 1$ |
|---|---|---|---|
| $\mathrm{do}(X_1 = 0)$ | $X_2 = 0$ | $1 - p$ | $0$ |
| | $X_2 = 1$ | $p$ | $0$ |
| $\mathrm{do}(X_1 = 1)$ | $X_2 = 0$ | $1 - p$ | $0$ |
| | $X_2 = 1$ | $0$ | $p$ |

For simplicity, but without loss of generality, we consider two treatment variables $\{X_1, X_2\}$ and no covariates. We will show that two different SCMs $\mathcal{M}$ and $\mathcal{M}'$ can yield identical observational distributions, as well as joint interventional distributions. They will also agree on the single-variable interventional distribution for

Table 1: SCMs for 3.1

| $\mathcal{M}$ | $\mathcal{M}'$ |
|---|---|
| $X_1 = U_1$ | $X_1 = U_1$ |
| $X_2 = X_1 U_2$ | $X_2 = U_2$ |
| $Y\ \ = X_1 X_2 U_Y$ | $Y\ \ = X_1 X_2 U_Y$ |

treatment $X_1$, but disagree on the single-variable effect of treatment $X_2$. As such, given data from sets of interventions, this example shows the treatment effect of single-variable interventions is not generally identifiable. Our SCMs are defined in Table 1, where $U_1 = U_2 = U_Y \sim \mathrm{Bernoulli}(p)$. The observational, joint, and single-variable interventional distributions on $X_2$ are identical under $\mathcal{M}$ and $\mathcal{M}'$; we defer them to the supplemental material. The interventional distribution on $X_1$ differs for $\mathcal{M}$ and $\mathcal{M}'$, as shown in Table 2. This counterexample shows that single-variable interventional effects are not identifiable for general SCMs, even in simple cases with two treatments.

## 3.2 Unidentifiability for general ANMs

An Additive Noise Model is an SCM where the influence of the latent noise variables is restricted to be additive in the structural equations. That is, Eq. 1 is restricted to the form:

$$X_i := f_i(\mathrm{PA}(X_i)) + U_i. \tag{2}$$

Following Saengkyongam and Silva (2020), we additionally assume the noise distribution to be a zero-centered Gaussian with an arbitrary covariance matrix: $\mathsf{P}_U \sim \mathcal{N}(0, \Sigma)$. Table 3 defines two SCMs that satisfy these assumptions. $\mathcal{M}$ and $\mathcal{M}'$ yield identical observational, joint interventional, and single-interventional distributions on $X_2$. However—they disagree on the causal effect of intervening on $X_1$. An intuitive underlying reason for this, is that we *can* identify the expression $\mathbb{E}[X_2 \mid X_1 = x_1] = f_2(x_1) + \mathbb{E}[U_2 \mid X_1 = x_1]$, but have no tools to disentangle the effects coming from the structural equation, $f_2(x_1)$, from those stemming from the additive noise, $\mathbb{E}[U_2 \mid X_1 = x_1]$. As a result, $\mathbb{E}[X_2 \mid \mathrm{do}(X_1 = x_1)] = f_2(x_1)$ is not identifiable, proving that general ANMs with unrestricted causal structures are insufficient for disentangling treatment effects. In this example, and in general SCMs of this structure, joint interventions *can* be disentangled for the *consequence* treatment $X_2$, but not for the *causing* treatment $X_1$:

**Theorem 1** (Identifiability of disentangled conditional average treatment effects in additive noise models with a causal dependency between treatments).
*Let $\mathcal{M} = \langle \{\boldsymbol{C}, \boldsymbol{X}, Y\}, \boldsymbol{U}, \boldsymbol{f}, \mathsf{P}_U \rangle$ be an SCM, where*

$$X_i = f_i(\boldsymbol{C}) + U_i, \qquad X_j = f_j(\boldsymbol{C}, X_i) + U_j, \qquad Y = f_Y(\boldsymbol{C}, \boldsymbol{X}) + U_Y,$$

*and $\mathsf{P}_U \sim \mathcal{N}(0, \Sigma)$. The estimand $\mathbb{E}[Y \mid do(X_j), C]$ is identifiable from the conjunction of two data regimes: (1) the observational distribution, and (2) the joint interventional distribution on $(X_i, X_j)$.*

We prove this by showing that the structural equations are identifiable from the joint interventional regime—whereas the necessary (co-)variances are identifiable from observational data. A formal proof is deferred to the supplementary material.

## 3.3 Identifiability for Symmetric ANMs

The crux of the problem of non-identifiability in Section 3.2 comes from the fact that treatment $X_1$ has a direct causal effect on treatment $X_2$. In many realistic applications, this might never occur.

4

Table 3: SCMs for 3.2, where $(U_1, U_2, U_Y)_{\mathcal{M}} \sim \mathcal{N}(0, \Sigma)$ and $(U_1, U_2, U_Y)_{\mathcal{M}'} \sim \mathcal{N}(0, \Sigma')$.

| $\mathcal{M}$ | $\mathcal{M}'$ |
|---|---|
| $X_1 = U_1$ | $X_1 = U_1$ |
| $X_2 = X_1 + U_2$ | $X_2 = 2X_1 + U_2$ |
| $Y = X_1 X_2 + U_Y$ | $Y = X_1 X_2 + U_Y$ |

where $\Sigma = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$, and $\Sigma' = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 1 \end{bmatrix}$.

Treatments will decidedly be correlated, but this can be encoded in the SCM via either the observed covariates or the unobserved confounders. Then still, as we have no restrictions on the functional form of the structural equations with respect to the treatments (i.e. $f_Y$), treatments can still yield highly nonlinear interaction effects on the outcome. As such, we now focus on the case where the SCM is devoid of causal links between treatments, but includes observed covariates as well as unobserved confounders. In this case, the structural equations take the following form:

$$X_i = f_i(\boldsymbol{C}) + U_i, \qquad Y = f_Y(\boldsymbol{C}, \boldsymbol{X}) + U_Y. \tag{3}$$

Figure 1 visualises the structure of the causal graph in this setting, for $K$ treatment variables.

**Theorem 2** (Identifiability of disentangled conditional average treatment effects in additive noise models with symmetric structure)**.**
*Let* $\mathcal{M} = \langle \{\boldsymbol{C}, \boldsymbol{X}, Y\}, \boldsymbol{U}, \boldsymbol{f}, \mathsf{P}_U \rangle$ *be an SCM, where*

$$X_i = f_i(\boldsymbol{C}) + U_i, \quad \forall i = 1, \dots, K,$$
$$Y = f_Y(\boldsymbol{C}, \boldsymbol{X}) + U_Y,$$

$C \perp\!\!\!\perp U$, *and* $\mathsf{P}_U \sim \mathcal{N}(0, \Sigma)$ *(following the DAG in Fig. 1). The estimand* $\mathbb{E}[Y|do(X_i), C]$ *is identifiable from the conjunction of two data regimes:*

1. *the observational distribution,*

2. *any interventional distribution on a set of treatments* $\boldsymbol{X}_{int} \subseteq \boldsymbol{X}$ *that contains* $X_i$: $X_i \in \boldsymbol{X}_{int}$.
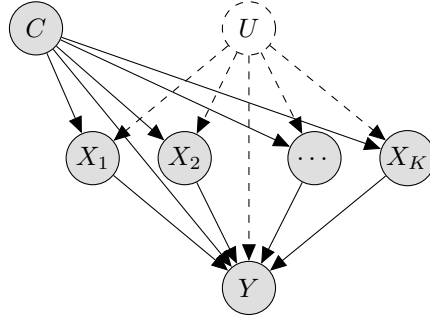


Figure 1: Causal graph for a symmetric SCM with observed covariates and unobserved confounders.

*Proof.* Our causal estimand is the effect of intervening on $X_i$. For notational convenience, we assume $k = i$, i.e. the intervention is on the last treatment. From the observational data regime, we can trivially obtain the joint $\mathsf{P}(\boldsymbol{C}, X_1, \dots, X_k, Y)$. As such, we can condition on the covariates and remaining treatment variables and then marginalise to obtain $\mathbb{E}[Y|do(X_k = x_k), C = c]$. We can rewrite our causal estimand as follows:

$$\begin{aligned} &\mathbb{E}[Y|C = c, X_1 = x_1, \dots, X_{k-1} = x_{k-1}, do(X_k = x_k)] \\ =&f_Y(c, x_1, \dots, x_k) + \mathbb{E}[U_Y|C = c, X_1 = x_1, \dots, X_{k-1} = x_{k-1}]. \end{aligned} \tag{4}$$

From the joint interventional data regime, we have access to the following expectation:

$$\mathbb{E}[Y|C = c, do(X_1 = x_1, \dots, X_k = x_k)] = f_Y(c, x_1, \dots, x_k). \tag{5}$$

Subtracting Eq. 5 from Eq. 4 shows that we only need to provide identifiability for the conditional expectation on the outcome noise, given the remaining treatment variables and the observed covariates:

$$\begin{aligned} \mathbb{E}[U_Y|C =&c, X_1 = x_1, \dots, X_{k-1} = x_{k-1}] = \Sigma_{u_y} \Sigma_{u_x}^{-1} \boldsymbol{u}_x, \\ &\text{where } \Sigma_{u_y} = \begin{bmatrix} \sigma_{Y1} & \dots & \sigma_{Y(k-1)} \end{bmatrix}, \\ &\Sigma_{u_x} = \begin{bmatrix} \sigma_1^2 & \dots & \sigma_{1(k-1)} \\ \vdots & \ddots & \vdots \\ \sigma_{(k-1)1} & \dots & \sigma_{k-1}^2 \end{bmatrix}, \\ &\text{and } \boldsymbol{u}_x = \begin{bmatrix} x_1 - f_1(c) & \dots & x_{k-1} - f_{k-1}(c) \end{bmatrix}^\mathsf{T}. \end{aligned} \tag{6}$$

5

(a) Only $\mathbb{E}[Y|C = c, \mathrm{do}(X_2 = x_2)]$ is generally identifiable.

(b) Only $\mathbb{E}[Y|C = c, \mathrm{do}(X_3 = x_3)]$ is generally identifiable.

(c) All $\mathbb{E}[Y|C = c, \mathrm{do}(X_i = x_i)]$ are generally identifiable.
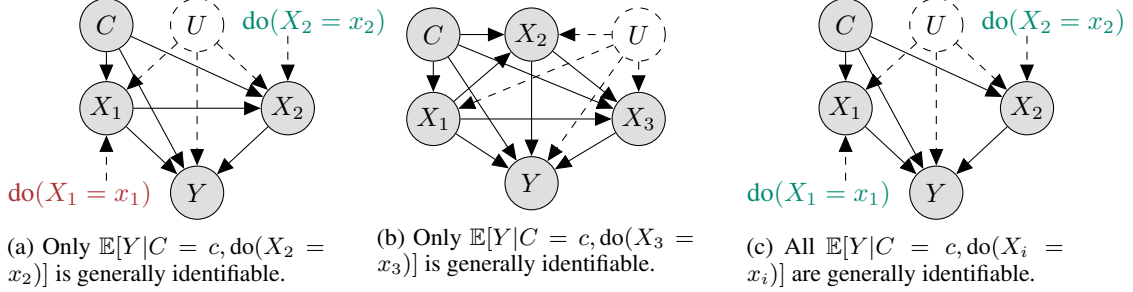
Figure 2: Causal Graphs illustrating under which conditions the single-variable causal effect on the outcome is identifiable from the observational and joint interventional data regimes.

Here, $\sigma_{ij}$ denotes the covariance between noise variables $U_i$ and $U_j$, and $\sigma_{Yi}$ denotes the covariance between the outcome $U_Y$ and $U_i$. There are two types of unidentified factors in these expressions: the **structural equations**, that is, the $f_i$'s encapsulated in $\boldsymbol{u}_x$ from Eq. 6, and the parameters of the **noise distribution**, which are encapsulated in $\Sigma_{u_x}$ and $\Sigma_{u_y}$. We tackle these in what follows:

**Identifying the structural equations.** As a direct result of our model definition, we can obtain $\mathbb{E}[X_i|C = c] = f_i(c)$ from the observational data regime. This follows because in the DAG of Figure 1 $C$ and $U_i$ are independent for all $i$. As such, the structural equations (and thus $\boldsymbol{u}_x$ from Eq. 6) are identifiable.

**Identifying the noise distribution.** For any pair of treatment variables, we can obtain $\mathbb{E}[X_i|C = c, X_j = x_j] = f_i(c) + \mathbb{E}[U_i|X_j = x_j]$. The latter term in this expression can then be rewritten as $\mathbb{E}[U_i|U_j = x_j - f_j(c)]$, for fixed values of $x_j$ and $c$. Because we have shown the structural equations $(f_i, f_j)$ to be identifiable, we have that $\mathbb{E}[U_i|U_j]$ is identifiable, which gives us the covariance $\sigma_{ij}$. Hence, the entirety of $\Sigma_{u_x}$ is identifiable. The same procedure can be used for every entry of $\Sigma_{u_y}$, and the covariances $\sigma_{Yi}$ are identifiable as a result.

It follows naturally that $\mathbb{E}[Y|C = c, X_1 = x_1, \ldots, X_{k-1} = x_{k-1}, \mathrm{do}(X_k = x_k)]$ is identifiable. As we have data from the observational regime, any marginalisation of this query is identifiable as well, which concludes the proof. $\qquad\square$

In this section, we have studied for several classes of SCMs whether causal effects can be disentangled. Our results provide insights into the fundamental limits of learning and inference, and help crystallise which assumptions are necessary and sufficient to make disentanglement of effects from sets of interventions feasible. Figure 2 visualises and summarises our key results denoting *which* single-variable causal effects can be disentangled in ANMs with zero-mean Gaussian noise. In what follows, we present a learning methodology and validate our theoretical identifiability results with empirical observations.

## 4 Estimating SCMs from Observational and Joint Interventional Data

In this section, we introduce our methodology for estimating SCMs and providing estimates for treatment effects under *any* set of interventions. Estimating an SCM from a combination of observational and interventional data boils down to (1) estimating the structural equations, and (2) estimating the noise distribution. We extend the Expectation-Maximisation-style iterative algorithm proposed by Saengkyongam and Silva (2020) to handle observed covariates, and to *disentangle* causal effects instead of *combining* them. The resulting method is not limited to learning single-variable causal effects, as the results from Saengkyongam and Silva (2020) allow us to extend these newly learned single-variable effects to sets of interventions. As a result, our learning method is able to generalise to sets of interventions that were never observed in the training data, with the only restriction that every variable that makes up the set was part of *some* intervention set in the training data.

Say we intervene on a subset of treatments $\boldsymbol{X}_{\mathrm{int}} \subseteq \boldsymbol{X}$, and $\boldsymbol{X}_{\mathrm{obs}} \equiv \boldsymbol{X} \setminus \boldsymbol{X}_{\mathrm{int}}$. In general, we can decompose a causal query with interventions on $\boldsymbol{X}_{\mathrm{int}}$ as follows:

$$\mathbb{E}[Y|\boldsymbol{C}; \mathrm{do}(\boldsymbol{X}_{\mathrm{int}}); \boldsymbol{X}_{\mathrm{obs}}] = f_Y(\boldsymbol{C}; \boldsymbol{X}) + \mathbb{E}[U_Y|\boldsymbol{X}_{\mathrm{obs}}]. \qquad (7)$$

---

**Algorithm 1** SCM Estimation for Symmetric ANMs

---

**Input:** Dataset $\mathcal{D}$
**Output:** Parameter estimates $\widehat{\theta}, \widehat{\Sigma}$
 1: Initialise $\widehat{\theta}$ and $\widehat{\Sigma}$
 2: **while** not converged **do**
 3:     // Solve for $\theta$ with fixed $\widehat{\Sigma}$
 4:     Optimise log-likelihood in Eq. 9
 5:     // Solve for $\Sigma$ with fixed $\widehat{\theta}$
 6:     Estimate $\widehat{\Sigma}$ from $\widehat{U} = \boldsymbol{x} - \boldsymbol{f}(\boldsymbol{x}; \widehat{\theta})$
 7: **return** $\widehat{\theta}, \widehat{\Sigma}$

---

Samples consist of observed values for all endogenous variables. For convenience, we denote a sample by $\boldsymbol{x}$. Suppose we have data from $d$ different data regimes, corresponding to $d$ different sets of interventions (the empty set $\emptyset$ corresponds to the observational regime). The full dataset consists of samples and their corresponding interventions $\mathcal{D} = \{(\boldsymbol{X}_{\mathrm{int}}; \boldsymbol{x})\}$.

**Estimating the structural equations.** We parameterise the structural equations with $\theta$, denoted as $f(\cdot; \theta)$. The Gaussian likelihood with covariance matrix $\Sigma$ is denoted as $\mathsf{P}_U(\cdot; \Sigma)$.

The likelihood for a single endogenous variable $x_i$ is defined as $L(x_i; \theta, \Sigma) = \mathsf{P}_U(x_i - f_i(\mathrm{PA}(x_i); \theta); \Sigma)$. The likelihood for a sample $\boldsymbol{x}$ is defined as the sum of the likelihoods for every endogenous variable that was *not intervened on* in that sample:

$$L(\boldsymbol{x}; \boldsymbol{X}_{\mathrm{int}}, \theta, \Sigma) = \sum_{x_i \in \boldsymbol{X} \setminus \boldsymbol{X}_{\mathrm{int}}} L(x_i; \theta, \Sigma). \tag{8}$$

Naturally, the log-likelihood of the dataset is then:

$$\ell(\mathcal{D}; \theta, \Sigma) = \sum_{(\boldsymbol{X}_{\mathrm{int}}, \boldsymbol{x}) \in \mathcal{D}} \log L(\boldsymbol{x}; \boldsymbol{X}_{\mathrm{int}}, \theta, \Sigma). \tag{9}$$

In principle, any iterative optimisation procedure can be used to maximise Eq. 9 when we fix the covariance $\Sigma$. Typically, an appropriate method is chosen with respect to the model parameterisation.

**Estimating the noise distribution.** Following Saengkyongam and Silva (2020), we note that the maximum likelihood estimate for the covariance matrix of a multivariate normal can be directly computed from the sample. This allows for efficient closed-form computation of this step. We additionally note that the assumption of zero-mean Gaussian noise simplifies the proof and yields an analytical solution for this step, but this is not a general limitation of the method. Indeed—what matters is that we can estimate the conditional noise $\mathbb{E}[U_Y | \boldsymbol{X}_{\mathrm{obs}}]$, which can in principle just as well be estimated via a different model.

Algorithm 1 shows the full iterative procedure that we adopt to estimate the parameters of a symmetric ANM. At inference time, Eq. 7 allows us to estimate the outcome under any—even potentially unseen—set of interventions.

## 5   Empirical Validation and Discussion

In this section, we empirically validate the effectiveness of our method in estimating SCMs from observational and joint-interventional data, and assess the accuracy of the inferred outcomes under varying sets of interventions. We adopt a simulation setup, which gives us the freedom to vary the *true* underlying SCM and observe performance differences among competing methods. The structural equation functions $f_i, f_Y$ in Eq. 3 are polynomials with second-order interactions to illustrate the effectiveness of the learning method even when treatment effects are highly non-linear, and the optimisation surface is highly non-convex. We use the well-known Adam optimiser in our
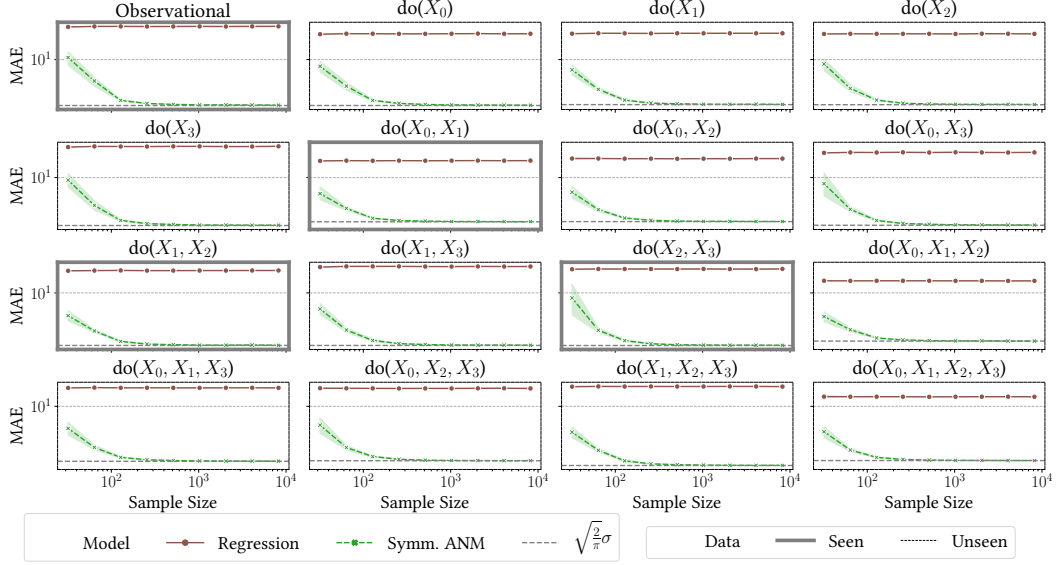
Figure 3: Mean Absolute Error for estimating the outcome under varying sets of interventions, for varying sample sizes, comparing the baseline regression approach with our proposed method. We empirically validate that our estimator is unbiased and consistent, exhibiting optimal predictions even with modest training sample sizes.

experiments (Kingma and Ba, 2014). As a baseline learning method we adopt a regression model that estimates the outcome from its direct treatments using pooled data from different regimes. Here, the presence of unobserved confounders severely hinders the method to provide accurate estimates under varying sets of interventions, further motivating the need for causal models. The research questions we wish to answer from empirical results, are the following:

**RQ1** Is our method able to accurately estimate the outcome under previously unseen sets of interventions?

**RQ2** Is our method able to accurately estimate the parameters of the structural equations and the noise distribution?

We let $|\boldsymbol{C}| = |\boldsymbol{X}| = 4$, yielding $\theta_i \in \mathbb{R}^{11} \quad \forall f_i$, and $\theta_y \in \mathbb{R}^{37}$ for the parameterisation of $f_y$. The parameters for the structural equations are randomly sampled from a uniform distribution over $[-2, +2]$. The covariance matrix $\Sigma \in \mathbb{R}^{5 \times 5}$ is uniformly sampled from a uniform distribution over $[-1, +1]$, and then ensured to be positive semi-definite through Ledoit-Wolf shrinkage (Ledoit and Wolf, 2004). We repeat this process 10 times with varying random seeds and report confidence intervals over obtained measurements.

We vary the size of the training sample $|\mathcal{D}| \in \{2^i \forall i = 3, \ldots, 11\}$, where the set of intervened treatments $\boldsymbol{X}_{\text{int}}$ is one of $\{\emptyset; \{X_0, X_1\}; \{X_1, X_2\}; \{X_2, X_3\}; \}$. As such, no single-variable interventions and the majority of the possible sets of interventions are never observed. From this data, we estimate an SCM using the procedure laid out in Algorithm 1.

Then, for every possible $(2^{|\boldsymbol{X}|})$ set of interventions, we simulate $10\,000$ samples and estimate the outcome based on our estimated SCM using Eq. 7. We report the Mean Absolute Error (MAE) between our estimated outcome and the true outcome. Note that, as the true outcome is Gaussian, the optimal estimate is the location of that Gaussian, which will have an expected deviation of $\sqrt{\frac{2}{\pi}}\sigma$ where $\sigma$ is the standard deviation on the noise parameter $U_Y$.

**Estimating outcomes (RQ1).** Figure 3 visualises the results from the procedure laid out above, increasing the size of the training sample over the x-axis and reporting the MAE on the y-axis. Note that both axes are logarithmically scaled, and the shaded regions indicate 95% confidence intervals. The plots clearly indicate that our method is able to provide accurate and close to optimal estimates
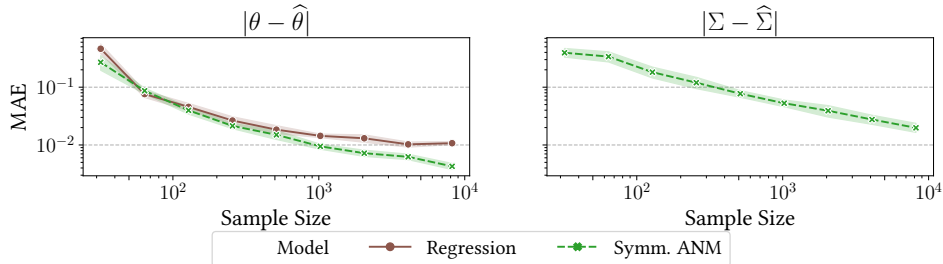
Figure 4: Mean Absolute Error for estimating the underlying SCM under varying sets of interventions, for varying sample sizes, comparing the baseline regression approach with our proposed method. We empirically validate that our estimator is unbiased and consistent, yielding accurate estimates even at modest training sample sizes.

for the outcome, under *any* possible set of interventions. Indeed, the plots labelled *observational*, $do(X_0, X_1)$, $do(X_1, X_2)$ and $do(X_2, X_3)$ show the accuracy of our method to estimate the outcome under interventions that occurred in the training data—but the remaining twelve plots relate to previously *unseen* data regimes. While our results allow us to *disentangle* joint interventions, the results of Saengkyongam and Silva (2020) allow us to *combine* interventions. As such, our method incorporates (and subsumes) theirs, in order to generalise to arbitrary sets of interventions. In contrast, the regression method that is oblivious to confounders fails to accurately estimate the outcome under any data regime—even those that are seen in the training data, or those where confounders yield no influence on the outcome (i.e. all treatments are intervened on). As such, the reported results validate that the proposed method provides an unbiased and consistent estimator for the disentangled causal effect, learned from sets of interventions in the presence of unobserved confounders.

**Estimating SCMs (RQ2).**   We visualise the results with respect to the accuracy of the parameter estimates in Figure 4. These measurements are obtained using the same procedure and runs as for RQ1. Increasing the size of the training sample over the x-axis, the leftmost plot shows the accuracy of the estimated parameters for $f_Y$. The rightmost plot shows the accuracy of the estimated covariance matrix for the multivariate normal that defines the noise distribution. As the baseline regression method is oblivious to the noise distribution, it is not included in the latter plot. Both axes are logarithmically scaled, and the shaded regions indicate 95% confidence intervals. We observe that our method is able to efficiently and effectively estimate the underlying causal model—empirically validating the identifiability results presented in this work.

## 6   Conclusions and Future Work

In this work, we motivated the need for methods that can disentangle the effects of single interventions from jointly applied interventions. As multiple interventions are bound to interact in possibly complex ways, this is a challenging task; even more so in the presence of unobserved confounders. First, we proved that such disentanglement is not possible in the general setting, even when we restrict the influence of the unobserved confounders to be additive in nature. By restricting the structure of the causal graph to be symmetric—void of edges between treatments—we showed that an identifiability result *can* be acquired. Additionally, we showed how to incorporate observed covariates into an existing learning method for *joint* interventional effects, and have empirically demonstrated how it can estimate the outcome under *arbitrary* sets of interventions, even those unseen in the training sample.

In future work, we wish to tackle the case where the noise distribution is not restricted to a zero-mean multivariate Gaussian. As we have hinted at in Section 4, this assumption provides a closed-form expression for the conditional expectation of the outcome noise given the observed treatments—but universal function approximators could be leveraged to obtain similar guarantees for more general model classes. Additionally, we wish to further validate our method on real-world data containing jointly applied interventions, as well as a ground truth for the outcome on single-variable interventions.

9

# References

Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352.

Correa, J. D., Lee, S., and Bareinboim, E. (2021). Nested counterfactual identification from arbitrary surrogate experiments. *arXiv preprint arXiv:2107.03190*.

Egami, N. and Imai, K. (2018). Causal interaction in factorial experiments: Application to conjoint analysis. *Journal of the American Statistical Association*.

Ghassami, A., Sani, N., Xu, Y., and Shpitser, I. (2021). Multiply robust causal mediation analysis with continuous treatments. *arXiv preprint arXiv:2105.09254*.

Gilligan-Lee, C. (2020). Causing trouble. *New Scientist*, 246(3279):32–35.

Hoyer, P., Janzing, D., Mooij, J., Peters, J., and Schölkopf, B. (2009). Nonlinear causal discovery with additive noise models. In *Twenty-Second Annual Conference on Neural Information Processing Systems (NIPS 2008)*, pages 689–696. Curran.

Ilse, M., Forré, P., Welling, M., and Mooij, J. M. (2021). Efficient causal inference from combined observational and interventional data through causal reductions. *arXiv preprint arXiv:2103.04786*.

Janzing, D., Peters, J., Mooij, J., and Schölkopf, B. (2009). Identifying confounders using additive noise models. In *25th Conference on Uncertainty in Artificial Intelligence (UAI 2009)*, pages 249–257. AUAI Press.

Kilbertus, N., Ball, P. J., Kusner, M. J., Weller, A., and Silva, R. (2020). The sensitivity of counterfactual fairness to unmeasured confounding. In *Uncertainty in Artificial Intelligence*, pages 616–626. PMLR.

Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Ledoit, O. and Wolf, M. (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis*, 88(2):365–411.

Lee, S., Correa, J. D., and Bareinboim, E. (2020). General identifiability with arbitrary surrogate experiments. In *Uncertainty in Artificial Intelligence*, pages 389–398. PMLR.

Lopez, M. J. and Gutman, R. (2017). Estimation of causal effects with multiple treatments: a review and new ideas. *Statistical Science*, pages 432–454.

Nandy, P., Maathuis, M. H., and Richardson, T. S. (2017). Estimating the effect of joint interventions from observational data in sparse high-dimensional settings. *The Annals of Statistics*, 45(2):647–674.

Parbhoo, S., Bauer, S., and Schwab, P. (2021). Ncore: Neural counterfactual representation learning for combinations of treatments. *arXiv preprint arXiv:2103.11175*.

Parbhoo, S., Wieser, M., Roth, V., and Doshi-Velez, F. (2020). Transfer learning from well-curated to less-resourced populations with hiv. In *Machine Learning for Healthcare Conference*, pages 589–609. PMLR.

Pearl, J. (2009). *Causality*. Cambridge university press.

Saengkyongam, S. and Silva, R. (2020). Learning joint nonlinear effects from single-variable interventions in the presence of hidden confounders. In Peters, J. and Sontag, D., editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of Machine Learning Research*, pages 300–309. PMLR.

Schwab, P., Linhardt, L., Bauer, S., Buhmann, J. M., and Karlen, W. (2020). Learning counterfactual representations for estimating individual dose-response curves. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5612–5619.

## A Unidentifiability under Unconstrained SCMs

These tables hold the full distributions for the counterexample in Section 3.1, showing that single-variable interventional effects are unidentifiable from observational and joint interventional data for unconstrained SCMs.

Table 4: Distributions under both SCMs $\mathcal{M}$ and $\mathcal{M}'$.

(a) Observational joint distribution.

| $\mathsf{P}(X_1, X_2, Y)$ | $Y = 0$ | $Y = 1$ |
|---|---|---|
| $X_1, X_2 = 0, 0$ | $1 - p$ | $0$ |
| $X_1, X_2 = 0, 1$ | $0$ | $0$ |
| $X_1, X_2 = 1, 0$ | $0$ | $0$ |
| $X_1, X_2 = 1, 1$ | $0$ | $p$ |

(b) Joint interventional distribution.

| $\mathsf{P}(Y \mid \mathrm{do}(X_1, X_2))$ | $Y = 0$ | $Y = 1$ |
|---|---|---|
| $\mathrm{do}(X_1 = 0, X_2 = 0)$ | $1$ | $0$ |
| $\mathrm{do}(X_1 = 0, X_2 = 1)$ | $1$ | $0$ |
| $\mathrm{do}(X_1 = 1, X_2 = 0)$ | $1$ | $0$ |
| $\mathrm{do}(X_1 = 1, X_2 = 1)$ | $1 - p$ | $p$ |

(c) Interventional distribution on $X_2$.

| $\mathsf{P}(Y, X_1 \mid \mathrm{do}(X_2))$ | | $Y = 0$ | $Y = 1$ |
|---|---|---|---|
| $\mathrm{do}(X_2 = 0)$ | $X_1 = 0$ | $1 - p$ | $0$ |
| | $X_1 = 1$ | $p$ | $0$ |
| $\mathrm{do}(X_2 = 1)$ | $X_1 = 0$ | $1 - p$ | $0$ |
| | $X_1 = 1$ | $0$ | $p$ |

## B Identifiability with Causally Dependent Treatments

This section provides a proof for identifiabiltiy of single-variable causal effects when there is a causal dependency among treatments (**Theorem 1**). Observational and joint interventional data are not sufficient in this case to identify causal effects on *all* treatments – but we can identify the causal effect of intervening on the *consequence* treatment instead of the *causing* treatment.

*Proof.* Our causal estimand is the effect of intervening on $X_j$. We can rewrite our causal estimand as follows:

$$\mathbb{E}[Y | X_i = x_i, \mathrm{do}(X_j = x_j), C = c] = f_Y(c, x_i, x_j) + \mathbb{E}[U_Y | C = c, X_i = x_i]. \tag{10}$$

From the joint interventional data regime, we have access to the following expectation:

$$\mathbb{E}[Y | C = c, \mathrm{do}(X_i = x_i, X_j = x_j)] = f_Y(c, x_i, x_j). \tag{11}$$

Subtracting Eq. 11 from Eq. 10 shows that we only need to provide identifiability for the conditional expectation on the outcome noise, given the observed value for treatment $X_i$ and the observed covariates $C$:

$$\mathbb{E}[U_Y | C = c, X_i = x_i] = \mathbb{E}[U_Y | U_i = x_i - f_i(c)] = \frac{\sigma_{Yi}}{\sigma_{ii}}(x_i - f_i(c)). \tag{12}$$

Here, the first step comes from our SCM definition, and the second step comes from the fact that we assume the noise distribution to be a zero-centered multivariate Gaussian. As such, we need to identify the function $f_i$, the variance on the noise variable $U_i$, and the covariance between $U_i$ and $U_Y$. We obtain $\mathbb{E}[X_i | C = c] = f_i(c)$ directly from the observational data regime. This makes the noise variable $U_i = X_i - f_i(C)$ identifiable. Now, as a result, we can identify its variance $\sigma_{ii}$ and covariance $\sigma_{Yi}$ from the observational data regime, which concludes the proof. $\qquad\square$