
Non-Parametric Path Analysis in Structural Causal Models

Junzhe Zhang and Elias Bareinboim
Purdue University, USA
{zhang745, eb}@purdue.edu

Abstract

One of the fundamental tasks in causal inference is to decompose the observed association between a decision X and an outcome Y into its most basic structural mechanisms. In this paper, we introduce counterfactual measures for effects along with a specific mechanism, represented as a path from X to Y in an arbitrary structural causal model. We derive a novel non-parametric decomposition formula that expresses the covariance of X and Y as a sum over unblocked paths from X to Y contained in an arbitrary causal model. This formula allows a fine-grained path analysis without requiring a commitment to any particular parametric form, and can be seen as a generalization of Wright’s decomposition method in linear systems (1923,1932) and Pearl’s non-parametric mediation formula (2001).

1 INTRODUCTION

Analyzing the relative strength of different pathways between a decision X and an outcome Y is a topic that has interested both scientists and practitioners across disciplines for many decades. Specifically, path analysis allows scientists to explain how Nature’s “black-box” works, and practically, it enables decision analysts to predict how an environment will change under a variety of policies and interventional conditions [Wright, 1923; Baron and Kenny, 1986; Bollen, 1989; Pearl, 2001].

More recently, understanding using causal inference tools how a black-box decision-making system operates has been a target of growing interest in the Artificial Intelligence community, most prominently in the context of Explainability, Transparency, and Fairness [Lu Zhang, 2017; Kusner *et al.*, 2017; Zafar *et al.*, 2017; Kilbertus *et al.*, 2017; Zhang and Bareinboim, 2018a]. For exam-

ple, consider the *standard fairness model* described in Fig. 1(a) that is concerned with the relation between a hiring decision (Y) and an applicant’s religious beliefs (X), which are *mediated* by the location (W), and *confounded* by the education background (Z) of the applicant.¹ Directed edges represent functional relations between variables. The relationship between X and Y is materialized through four different pathways in the system – the *direct* path $l_1 : X \rightarrow Y$, the *indirect* path $l_2 : X \rightarrow W \rightarrow Y$, and the *spurious* paths $l_3 : X \leftarrow Z \rightarrow Y$ and $l_4 : X \leftarrow Z \rightarrow W \rightarrow Y$.

Assuming, for simplicity’s sake, that the functional relationships are linear and U_{V_i} is an independent “error” associated with each variable V_i (called the linear-standard model), Fig. 1(a) shows the structural coefficients corresponding to each edge – i.e., the value of the variable Y is decided by the structural function $Y \leftarrow \alpha_{YX}X + \alpha_{YZ}Z + \alpha_{YW}W + U_Y$. The celebrated result known as Wright’s method of path coefficients [Wright, 1923, 1934], also known as Wright’s rule, allows one to express the covariance of X and Y , denoted by $\text{Cov}(X, Y)$, as the sum of the products of the structural coefficients along the paths from X to Y in the underlying causal model.² In particular, $\text{Cov}(X, Y)$ is equal to:

$$\underbrace{\alpha_{YX}}_{X \rightarrow Y} + \underbrace{\alpha_{WX}\alpha_{YW}}_{X \rightarrow W \rightarrow Y} + \underbrace{\alpha_{XZ}\alpha_{YZ}}_{X \leftarrow Z \rightarrow Y} + \underbrace{\alpha_{XZ}\alpha_{WZ}\alpha_{YW}}_{X \leftarrow Z \rightarrow W \rightarrow Y}. \quad (1)$$

Using the observational covariance matrix, the decomposition above allows one to answer some compelling questions about the relationship between X and Y in the underlying model. For instance, the product $\alpha_{WX}\alpha_{YW}$ explains how much the indirect discrimination through the location (the path l_2) accounts for the observed disparities in the religion composition among hired employees.

The path analysis method gained momentum in the so-

¹This specific setting has been called *standard fairness model* given its generality to representing a variety of decision-making scenarios [Zhang and Bareinboim, 2018a].

²For a survey on linear methods, see [Pearl, 2000, Ch. 5].

cial sciences during 1960's, becoming extremely popular in the form of the *mediation formula* in which the total effect of X on Y is decomposed into direct and indirect components [Baron and Kenny, 1986; Bollen, 1989; Duncan, 1975; Fox, 1980].³ The bulk of this literature, however, required a commitment to a particular parametric form, thus falling short of providing a general method for analyzing natural and social phenomena with nonlinearities and interactions [MacKinnon, 2008].

It took a few decades until this problem could be tackled in higher generality. In particular, the advent of non-parametric structural causal models (SCMs) allowed this leap, and a more fine-grained path-analysis with a much broader scope, including models with nonlinearities and arbitrarily complex interactions [Pearl, 2000, Ch. 7]. In particular, Pearl introduced the *causal mediation formula* for arbitrary non-parametric models, which decomposes the total effect $TE_{x_0, x_1}(Y) = E[Y_{x_1}] - E[Y_{x_0}]$, the difference between the causal effect of the intervention $do(x_1)$ and $do(x_0)$ ⁴, into what is now known as the natural direct (*NDE*) and indirect (*NIE*) effects [Pearl, 2001] (see also [Imai *et al.*, 2010, 2011; VanderWeele, 2015]). In the case of the specific linear-standard causal model,

$$TE_{0,1}(Y) = \underbrace{\alpha_{YX}}_{NDE} + \underbrace{\alpha_{WX}\alpha_{YW}}_{NIE}$$

for $x_0 = 0$ and $x_1 = 1$ levels. Remarkably, when compared with Eq. 1, *NDE* and *NIE* capture the effects along with the direct and indirect paths, but omits the spurious (non-causal) paths between X and Y (in this case, l_3, l_4). The mediation formula was recently generalized to account for these spurious paths (more akin to Wright's rules), which appears under the rubric of the *causal explanation formula* [Zhang and Bareinboim, 2018a]. This formula decomposes the total variation $TV_{x_0, x_1}(Y) = E[Y|x_1] - E[Y|x_0]$ (difference in conditional distributions) into counterfactual measures of the direct (*Ctf-DE*), indirect (*Ctf-IE*), and spurious (*Ctf-SE*) effects. In the linear-standard model, for $x_0 = 0, x_1 = 1$,

$$TV_{0,1}(Y) = \underbrace{\alpha_{YX}}_{Ctf-DE} + \underbrace{\alpha_{WX}\alpha_{YW}}_{Ctf-IE} + \underbrace{\alpha_{XZ}\alpha_{YZ} + \alpha_{XZ}\alpha_{WZ}\alpha_{YW}}_{Ctf-SE}$$

Despite the generality of such results, there are still outstanding challenges when performing path analysis in non-parametric models, i.e.: (1) Estimands are defined relative to specific values assigned to the treatment x_1 and its baseline x_0 , which may be difficult to select in some non-linear settings; (2) Mediators and confounders

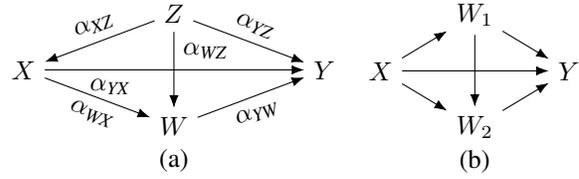


Figure 1: Causal diagrams for (a) the standard fairness model where X stands for the protected attribute, Y for the outcome, W the mediators, and Z the confounders; (b) the two-mediators setting where causal paths from X to Y are mediated by W_1, W_2 .

are collapsed and considered *en bloc*, leading to a coarse decomposition of the relationship between X and Y [Pearl, 2001; Vansteelandt and VanderWeele, 2012; Tchetgen and Shpitser, 2012; VanderWeele *et al.*, 2014; Daniel *et al.*, 2015; Zhang and Bareinboim, 2018a]; (3) Path-specific estimands are well-defined [Pearl, 2001; Avin *et al.*, 2005], but not in a way that they sum up to either the total effect (TE) or variation (TV), precluding the comparison of their relative strengths.

This paper aims to circumvent these problems. In particular, we decompose the covariance of a treatment X and an outcome Y over effects along different mechanisms between X and Y . We define a set of novel counterfactual estimands for measuring the relative strength of a specific mechanism represented as a path from X to Y in an arbitrary causal model. These estimands lead to a non-parametric decomposition formula, which expresses the covariance $\text{Cov}(X, Y)$ as a sum of the unblocked paths from X to Y in the causal graph. This formula allows a more fine-grained analysis of the total observed variations of Y due to X (both through causal and spurious mechanisms) when compared to the state-of-art methods. More specifically, our contributions are: (1) counterfactual covariance measures for a specific pathway from X to Y (causal and spurious) in an arbitrary causal model (Defs. 8, 11-12); (2) non-parametric decomposition formulae of the covariance $\text{Cov}(X, Y)$ over paths from X to Y in the causal model (Thm. 5); (3) the identification formulae estimating the proposed path-specific decomposition from the passively-collected data in the standard model (Thms. 6-7).

2 PRELIMINARIES

In this section, we introduce notations used throughout the paper. We will use capital letters to denote variables (e.g., X), and small letters for their values (x). The abbreviation $P(x)$ represents the probabilities $P(X = x)$. For arbitrary sets A and B , let $A - B$ denote their differ-

³Just to give an idea of this popularity, Baron and Kenny's original paper counts more than 70,000 citations.

⁴By convention [Pearl, 2000], the post-interventional distribution is represented interchangeably by $P(y_x)$ and $P(y|do(x))$. General notation is discussed in the next section.

ence, and let $|A|$ be the dimension of set A . $V_{[i,j]}$ stands for a set $\{V_i, \dots, V_j\}$ (\emptyset if $i > j$). We use graphical family abbreviations: $An(X)_G$, $De(X)_G$, $Non-De(X)_G$, $Pa(X)_G$, $Ch(X)_G$, which stand for the set of ancestors, descendants, non-descendants, parents and children of X in G . We omit the subscript G when obvious.

The basic semantical framework of our analysis rests on *structural causal models* (SCM) [Pearl, 2000, Ch. 7; Bareinboim and Pearl, 2016]. A SCM M consists of a set of endogenous variables V (often observed) and exogenous variables U (often unobserved). The values of each $V_i \in V$ are determined by a structural function f_i taking as argument a combination of the other endogenous and exogenous variables (i.e., $V_i \leftarrow f_i(PA_i, U_i)$, $PA_i \subseteq V, U_i \subseteq U$). Values of U are drawn from a distribution $P(u)$. A SCM M is called *Markovian* when the exogenous are mutually independent and each $U_i \in U$ is associated with only one endogenous $V_i \in V$. If U_i is associated with two or more endogenous variables, M is called *semi-Markovian*.

Each recursive SCM M has an associated causal diagram in the form of a directed acyclic graph (DAG) G , where nodes represent endogenous variables and directed edges represent functional relations (e.g., Figs. 1-2). By convention, the exogenous U are not explicitly shown in the graph; a dashed-bidirected arrow between V_i and V_j indicates the presence of an unobserved confounder (UC) U_k affecting both V_i and V_j (e.g., the path $V_i \leftarrow U_k \rightarrow V_j$).

A path from X to Y is a sequence of edges which does not include a particular node more than once. It may go either along or against the direction of the edges. Paths of the form $X \rightarrow \dots \rightarrow Y$ are *causal* (from X to Y). We use d-separation and blocking interchangeably, following the convention in [Pearl, 2000]. Any unblocked path that is not causal is called *spurious*. The direct link $X \rightarrow Y$ is the *direct* path and all the other causal paths from X to Y are called *indirect*. The set of unblocked paths from X to Y given a set Z in a causal diagram G is denoted by $\Pi(X, Y|Z)_G$; causal, indirect, and spurious paths are denoted by $\Pi^c(X, Y|Z)_G$, $\Pi^i(X, Y|Z)_G$, and $\Pi^s(X, Y|Z)_G$ (G will be omitted when obvious). For a causal path g including nodes V_1, V_2 , we denote $g(V_1, V_2)$ a subpath of g from V_1 to V_2 .⁵

An intervention on a set of endogenous variables X and exogenous variables U_i , denoted by $do(x^*, u_i^*)$, is an operation where values of X, U_i are set to x^*, u_i^* , respectively, without regard to how they were ordinarily determined (X through f_X and U_i through $P(U_i)$). Formally, we can rewrite the definition of potential response [Pearl, 2000, Ch. 7.1] to account for operation on U_i , namely:

⁵Mediators (relative to X and Y) are a set of variables $W \subseteq De(X) \cap Non-De(Y)$ such that $|\Pi^s(X, Y|W)| = 0$.

Definition 1 (Potential Response). Let M be a SCM, X, Y sets of arbitrary variables in V , and U_i a set of arbitrary variables in U . Let $U_{-i} = U - U_i$. The potential response of Y to the intervention $do(x^*, u_i^*)$ in the situation $U = u$, denoted by $Y_{x^*, u_i^*}(u)$, is the solution for Y with $U_{-i} = u_{-i}, U_i = u_i^*$ in the modified submodel M_{x^*} where functions f_X are replaced by constant functions $X = x^*$, i.e., $Y_{x^*, u_i^*}(u) \triangleq Y_{M_{x^*}}(u_i^*, u_{-i})$.⁶

$Y_{x^*, u_i^*}(u)$ can be read as the counterfactual sentence “the value that Y would have obtained in situation $U_{-i} = u_{-i}$, had the treatment X been x^* and the situation U_i been u_i^* .” Averaging u over the distribution $P(u)$, we obtain a counterfactual random variable Y_{x^*, u_i^*} . If the values of x^*, u_i^* follow random variables X^*, U_i^* , we denote the resulting counterfactual Y_{X^*, U_i^*} .

3 A COARSE COVARIANCE DECOMPOSITION

In this section, we introduce counterfactual measures that will allow us to non-parametrically decompose the covariance $Cov(X, Y)$ in terms of direct, indirect and spurious pathways from X to Y . Given space constraints, all proofs are included in [Zhang and Bareinboim, 2018b].

If there exists no spurious path from X to Y , then treatment X is independent of the counterfactual Y_{x^*} , i.e., $(X \perp\!\!\!\perp Y_{x^*})$ [Pearl, 2000, Ch. 11.3.2]. The *spurious covariance* can then be defined as the correlation between the factual variable X and counterfactual Y_{x^*} .

Definition 2 (Spurious Covariance). The spurious covariance between treatment $X = x^*$ and outcome Y is:

$$Cov_{x^*}^s(X, Y) = Cov(X, Y_{x^*}). \quad (2)$$

Property 1. $|\Pi^s(X, Y)| = 0 \Rightarrow Cov_{x^*}^s(X, Y) = 0$.

The *causal covariance* can naturally be defined as the difference between the total and spurious covariance.

Definition 3 (Causal Covariance). The causal covariance of the treatment $X = x^*$ and the outcome Y is:

$$Cov_{x^*}^c(X, Y) = Cov(X, Y - Y_{x^*}). \quad (3)$$

Prop. 2 establishes the correspondence between the causal paths and the causal covariance – if there is no causal path from X to Y in the underlying model, the causal covariance equates to zero.

Property 2. $|\Pi^c(X, Y)| = 0 \Rightarrow Cov_{x^*}^c(X, Y) = 0$.

We consider more detailed measures corresponding to the different causal pathways, and first, the direct path:

⁶An alternative way to see that the replacement operation relative to U_i is to envision a system where U_i is observed.

Definition 4 (Direct Covariance). Given a semi-Markovian model M , let the set W be the mediators between X and Y . The pure ($\text{Cov}_{x^*}^{dp}(X, Y)$) and total ($\text{Cov}_{x^*}^{dt}(X, Y)$) direct covariance of the treatment $X = x^*$ on the outcome Y are defined respectively as

$$\text{Cov}_{x^*}^{dp}(X, Y) = \text{Cov}(X, Y - Y_{x^*, W}), \quad (4)$$

$$\text{Cov}_{x^*}^{dt}(X, Y) = \text{Cov}(X, Y_{W_{x^*}} - Y_{x^*}). \quad (5)$$

By the composition axiom [Pearl, 2000, Ch. 7.3], Eqs. 4 and 5 can be explicitly written as follows ⁷:

$$\begin{aligned} \text{Cov}(X, Y - Y_{x^*, W}) &= \text{Cov}(X, Y_{X, W} - Y_{x^*, W}), \\ \text{Cov}(X, Y_{W_{x^*}} - Y_{x^*}) &= \text{Cov}(X, Y_{X, W_{x^*}} - Y_{x^*, W_{x^*}}). \end{aligned}$$

The counterfactual pure direct covariance (Eq. 4) is shown graphically in Fig. 2, where (a) corresponds to the Y -side, and (b) to the $Y_{x^*, W}$ -side. Note that from the mediator W perspective, X remains at the level that it would naturally have attained, while the “direct” input from X to Y varies from its natural level (Fig. 2a) to $do(x^*)$ (b). The change of the outcome Y thus measures the effect of the direct path. A similar analysis also applies to the total direct covariance (Eq. 5).

Property 3. $\text{Cov}_{x^*}^{dp}(X, Y) = \text{Cov}_{x^*}^{dt}(X, Y) = 0$ if X is not a parent of Y (i.e., $X \notin Pa(Y)$).

We can turn around the definitions of direct covariance and provide operational estimands for indirect paths.

Definition 5 (Indirect Covariance). Given a semi-Markovian model M , let the set W be the mediators between X and Y . The pure ($\text{Cov}_{x^*}^{ip}(X, Y)$) and total ($\text{Cov}_{x^*}^{it}(X, Y)$) indirect covariance of the treatment $X = x^*$ on the outcome Y are defined respectively as:

$$\text{Cov}_{x^*}^{ip}(X, Y) = \text{Cov}(X, Y - Y_{W_{x^*}}), \quad (6)$$

$$\text{Cov}_{x^*}^{it}(X, Y) = \text{Cov}(X, Y_{x^*, W} - Y_{x^*}). \quad (7)$$

Eqs. 6 and 7 correspond to the indirect paths, since they capture the covariance of X and Y , but only via paths mediated by W . The first argument of Y is the same in both halves of the contrast, but this value can either be x^* (Eq. 7) or at the level that X would naturally attain without intervention (Eq. 6).

Property 4. $|\Pi^i(X, Y)| = 0 \Rightarrow \text{Cov}_{x^*}^{ip}(X, Y) = \text{Cov}_{x^*}^{it}(X, Y) = 0$.

Putting these definitions together, we can prove a general non-parametric decomposition of $\text{Cov}(X, Y)$:

⁷Consider Eq. 4 as an example. For any $U = u$, $Y_{X(u), W(u)}(u) = Y_{x^*, w}(u)$ if $X(u) = x^*, W(u) = w$. By the composition axiom, $X(u) = x^*, W(u) = w$ implies $Y(u) = Y_{x^*, w}(u)$, which in turn gives $Y_{X(u), W(u)}(u) = Y(u)$. Averaging u over $P(u)$, we obtain $Y_{X, W} = Y$.

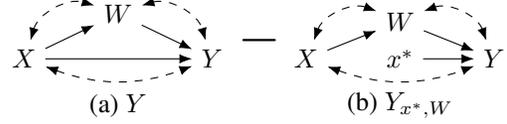


Figure 2: The graphical representation of measuring the pure direct covariance $\text{Cov}_{x^*}^{dp}(X, Y)$.

Theorem 1. $\text{Cov}(X, Y)$, $\text{Cov}_{x^*}^s(X, Y)$ and $\text{Cov}_{x^*}^c(X, Y)$ obey the following non-parametric relationship:

$$\text{Cov}(X, Y) = \text{Cov}_{x^*}^c(X, Y) + \text{Cov}_{x^*}^s(X, Y), \quad (8)$$

where $\text{Cov}_{x^*}^c(X, Y) = \text{Cov}_{x^*}^{dp}(X, Y) + \text{Cov}_{x^*}^{it}(X, Y) = \text{Cov}_{x^*}^{dt}(X, Y) + \text{Cov}_{x^*}^{ip}(X, Y)$.

In other words, the covariance between X and Y can be partitioned into its corresponding direct, indirect, and spurious components. In particular, Thm. 1 coincides with Eq. 1 in the linear-standard model.

Corollary 1. In the linear-standard model, for any x^* , $\text{Cov}_{x^*}^s(X, Y)$, $\text{Cov}_{x^*}^{dp}(X, Y)$, $\text{Cov}_{x^*}^{dt}(X, Y)$, $\text{Cov}_{x^*}^{ip}(X, Y)$ and $\text{Cov}_{x^*}^{it}(X, Y)$ are equal to:

$$\text{Cov}_{x^*}^s(X, Y) = \alpha_{XZ}\alpha_{YZ} + \alpha_{XZ}\alpha_{WZ}\alpha_{YW},$$

$$\text{Cov}_{x^*}^{dp}(X, Y) = \text{Cov}_{x^*}^{dt}(X, Y) = \alpha_{YX},$$

$$\text{Cov}_{x^*}^{ip}(X, Y) = \text{Cov}_{x^*}^{it}(X, Y) = \alpha_{WX}\alpha_{YW}.$$

Corol. 1 says that the proposed decomposition (Thm. 1) does not depend on the value of $do(x^*)$ in the linear model of Fig. 1(a), which is not achievable in previous value-specific decompositions [Pearl, 2001; Zhang and Bareinboim, 2018a].⁸

4 DECOMPOSING CAUSAL RELATIONS

We now focus on the challenge of decomposing the causal covariance into more elementary components. We use the two-mediators setting (Fig. 1(b)) as example, where X and Y are connected through four causal paths: through both W_1, W_2 ($g_1 : X \rightarrow W_1 \rightarrow W_2 \rightarrow Y$), only through W_1 ($g_2 : X \rightarrow W_1 \rightarrow Y$), only through W_2 ($g_3 : X \rightarrow W_2 \rightarrow Y$), and directly ($g_4 : X \rightarrow Y$). Our goal is to decompose the $\text{Cov}_{x^*}^c(X, Y)$ over the paths $g_{[1,4]}$. Our analysis applies to semi-Markovian models, without loss of generality, and the Markovian example (Fig. 1(b)) is used for simplicity of the exposition.

⁸For the nonlinear models, the decomposing terms (e.g., $\text{Cov}_{x^*}^s(X, Y)$) are still sensitive to the target level $do(x^*)$. To circumvent the challenges of picking a specific decision value, one could assign a randomized treatment $do(x^* \sim P(X))$, where $P(X)$ is the distribution over the treatment X induced by the underlying causal model.

For a node $S_i \in Pa(Y)$ and a set of causal paths π , the edge $S_i \rightarrow Y$ defines a funnel operator $\triangleleft_{S_i \rightarrow Y}$, which maps from π to the set of paths $\triangleleft_{S_i \rightarrow Y}(\pi)$ obtained from π by replacing all paths of the form $X \rightarrow \dots \rightarrow S_i \rightarrow Y$ with $X \rightarrow \dots \rightarrow S_i$, and removing all the other paths. As an example, for $\pi = \{g_1, g_2, g_3\}$, $\triangleleft_{W_2 \rightarrow Y}(\pi) = \{g_1(X, W_2), g_3(X, W_2)\}$, where $g_1(X, W_2)$ is the subpath $X \rightarrow W_1 \rightarrow W_2$ and $g_3(X, W_2)$ is the subpath $X \rightarrow W_2$. We next formalize the notion of path-specific interventions, which isolates the influence of the intervention $do(x^*)$ passing through a subset π of causal paths from X , denoted by $do(\pi[x^*])$ (a similar notion has been introduced by [Pearl, 2001], and then [Avin *et al.*, 2005; Shpitser and Tchetgen, 2016]).

Definition 6 (Path-Specific Potential Response). For a SCM M and an arbitrary variable $Y \in V$, let π be a set of causal paths. Let X be the source variables of paths in π . Further, let $X_{\pi \rightarrow Y} = \{X_i : \forall X_i \in X, X_i \rightarrow Y \in \pi\}$ and $S = (Pa(Y)_G \cap V) - X_{\pi \rightarrow Y}$. The π -specific potential response of Y to the intervention $do(\pi[x^*])$ in the situation $U = u$, denoted by $Y_{\pi[x^*]}(u)$, is defined as:

$$Y_{\pi[x^*]}(u) = \begin{cases} Y_{x^*_{\pi \rightarrow Y}, S_{\triangleleft_{S \rightarrow Y}(\pi)[x^*]}(u)} & \text{if } \pi \neq \emptyset \\ Y(u) & \text{otherwise} \end{cases}$$

where $S_{\triangleleft_{S \rightarrow Y}(\pi)[x^*]}(u)$ is a set of π -specific potential response $\{S_i \in S : S_i \rightarrow Y \in \pi[x^*]\}$.⁹

Despite the non-trivial notation, the π -specific counterfactual $Y_{\pi[x^*]}$ is simply assigning the treatment $do(x^*)$ exclusively to the causal paths in π , while allowing all the other causal paths to behave naturally. This contrasts with the counterfactual Y_{x^*} , which can be seen as assigning the treatment $do(x^*)$ to all causal paths from X to Y . For instance, repeatedly applying Def. 6 to $g_1 : X \rightarrow W_1 \rightarrow W_2 \rightarrow Y$ (see [Zhang and Bareinboim, 2018b, Sec. 2.1]), we obtain the g_1 -specific potential response $Y_{g_1[x^*]}$ as

$$Y_{g_1[x^*]} = Y_{X, W_1, W_2, W_2, W_1, x^*} = Y_{W_2, W_1, x^*}.$$

The intervention $do(g_1[x^*])$ can be visualized more immediately through its graphical representation (Fig. 3(b)) – the treatment $do(x^*)$ is assigned throughout g_1 while all the other paths are kept at the level that it would have attained “naturally” following X . The difference of the outcome Y (induced by $do(g_1[x^*])$) and the unintervened Y (Fig. 3(a)) measures the relative strength of g_1 itself, which leads to the following definition.

Definition 7 (Pure Path-Specific Causal Covariance). For a semi-Markovian model M and an arbitrary causal path g from X , the pure g -specific causal covariance of the treatment $X = x^*$ on the outcome Y is defined as:

$$\text{Cov}_g^c(x^*)(X, Y) = \text{Cov}(X, Y - Y_{g[x^*]}). \quad (9)$$

⁹For a single causal path g , let $Y_{g[x^*]}(u) = Y_{\{g\}[x^*]}(u)$. Averaging u over $P(u)$, we obtain a random variable $Y_{\pi[x^*]}$.

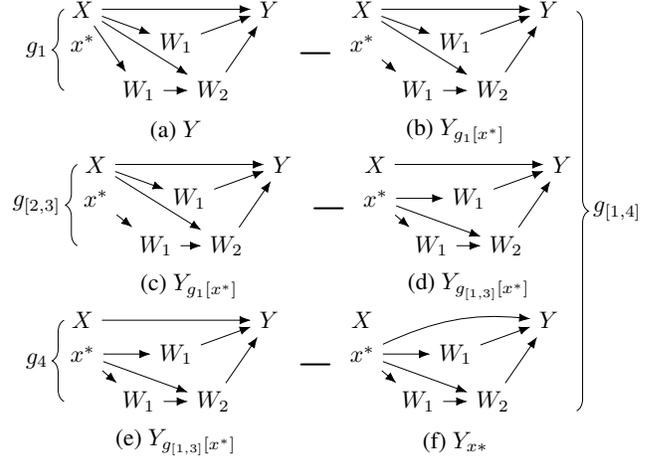


Figure 3: Graphical representations of the causal covariance specific to g_1 (a-b), $g_{[2,3]}$ (c-d) and g_4 (e-f).

In the previous example, more explicitly, the pure g_1 -specific causal covariance is equal to (Fig. 3(a-b)):

$$\text{Cov}_{g_1[x^*]}^c(X, Y) = \text{Cov}(X, Y - Y_{W_2, W_1, x^*}). \quad (10)$$

For $U = u$, the counterfactual $Y_{\emptyset[x^*]}(u)$ stands for the values of Y when all causal paths are under the natural regime. Eq. 9 can then be rewritten as:

$$\text{Cov}_g^c(x^*)(X, Y) = \text{Cov}(X, Y_{\emptyset[x^*]} - Y_{g[x^*]}).$$

The pure path-specific causal covariance for g can be seen as a function of the difference between two path-specific potential response $Y_{\pi_0[x^*]}$ and $Y_{\pi_1[x^*]}$ such that $g \notin \pi_0$ and $\pi_1 = \pi_0 \cup \{g\}$ (i.e., the different between π_1 and π_0 is g). The difference $Y_{\pi_1[x^*]} - Y_{\pi_0[x^*]}$, therefore, measures precisely the effects of $do(x^*)$ along the target causal path g . Def. 7 can be generalized to account for the path-specific covariance in terms of path-differences.

Definition 8 (Path-Specific Causal Covariance). For a semi-Markovian model M and an arbitrary causal path g from X , let π be a function mapping g to a set of causal paths $\pi(g)$ from X such that $g \notin \pi(g)$. The g -specific causal covariance of the treatment $X = x^*$ on the outcome Y is defined as:

$$\text{Cov}_{g[x^*]}^c(X, Y)_\pi = \text{Cov}(X, Y_{\pi(g)[x^*]} - Y_{\pi(g) \cup \{g\}[x^*]}).$$

The following property establishes the correspondence between a causal path and its path-specific estimand.

Property 5. $g \notin \Pi^c(X, Y) \Rightarrow \text{Cov}_{g[x^*]}^c(X, Y)_\pi = 0$.

Prop. 5 follows immediately as a corollary of Lem. 1, which implies that the counterfactuals $Y_{\pi(g)[x^*]}$ and $Y_{\pi(g) \cup \{g\}[x^*]}$ define the same variable over U if g is not a causal path from X to Y .

Lemma 1. $g \notin \Pi^c(X, Y) \Rightarrow Y_{\pi(g)[x^*]}(u) = Y_{\pi(g) \cup \{g\}[x^*]}(u)$.

Considering again the model in Fig. 1(b), let $g_{[i,j]} = \{g_k\}_{i \leq k \leq j}$ (\emptyset if $i > j$). Recall that $g_4 = \{X \rightarrow Y\}$, and note that the g_4 -specific causal covariance can be computed using $\pi(g_4) = g_{[1,3]}$, which yields:

$$\begin{aligned} \text{Cov}_{g_4[x^*]}^c(X, Y)_\pi &= \text{Cov}(X, Y_{g_{[1,3]}[x^*]} - Y_{g_{[1,4]}[x^*]}) \\ &= \text{Cov}(X, Y_{W_{1,x^*}, W_{2,x^*}} - Y_{x^*}), \end{aligned} \quad (11)$$

which coincides with the direct effect (Eq. 5 with $W = \{W_1, W_2\}$). Fig. 3(e-f) shows a graphical representation of this procedure.

The path-specific quantity given in Def. 8 has another desirable property, namely, the causal covariance $\text{Cov}_x^c(X, Y)$ can be decomposed as a summation over causal paths from X to Y . To witness, first let an order over $\Pi^c(X, Y)$ be $\mathcal{L}^c : g_1 < \dots < g_n$. For a path $g_i \in \Pi^c(X, Y)$, the order \mathcal{L}^c defines a function \mathcal{L}_π^c which maps from g_i to a set of paths $\mathcal{L}_\pi^c(g_i)$ that precede g_i in \mathcal{L}^c , i.e., $\mathcal{L}_\pi^c(g_i) = g_{[1, i-1]}$. We derive in the sequel a path-specific decomposition formula for the causal covariance relative to an order \mathcal{L}^c .

Theorem 2. For a semi-Markovian model M , let \mathcal{L}^c be an order over $\Pi^c(X, Y)$. For any x^* , the following non-parametric relationship hold:

$$\text{Cov}_x^c(X, Y) = \sum_{g \in \Pi^c(X, Y)} \text{Cov}_{g[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c}.$$

Thm. 2 can be demonstrated in the model of Fig. 1(a). Let an order \mathcal{L}^c over $g_{[1,4]}$ be $g_i < g_j$ if $i < j$. First note that the path-specific causal covariance of g_2 ($\text{Cov}_{g_2[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c}$) and g_3 ($\text{Cov}_{g_3[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c}$) are equal to, respectively,

$$\text{Cov}\left(X, Y_{W_{2W_{1,x^*}}} - Y_{W_{2W_{1,x^*}}, W_{1,x^*}}\right) \quad (12)$$

$$\text{Cov}\left(X, Y_{W_{2W_{1,x^*}}, W_{1,x^*}} - Y_{W_{1,x^*}, W_{2,x^*}}\right) \quad (13)$$

The causal covariance $\text{Cov}_x^c(X, Y)$ can then be decomposed as the sum of Eqs. 10-13, respectively, g_1, g_4, g_2, g_3 . Fig. 3 describes this decomposition procedure: we measure the difference of the outcome Y as the intervention $do(x^*)$ propagates through paths g_1, g_2, g_3, g_4 . The sum of these differences thus equate to the total influence of the intervention $do(x^*)$ to the outcome Y , i.e., the causal covariance $\text{Cov}_{x^*}^c(X, Y)$.

5 DECOMPOSING SPURIOUS RELATIONS

We introduce in the sequel a new strategy to decompose the spurious covariance (Def. 2), which will play a cen-

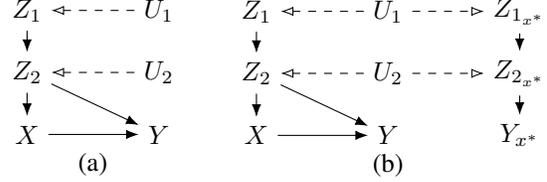


Figure 4: Causal diagrams for (a) the one-confounder setting where X and Y are confounded by the variable Z_2 , of which Z_1 is a parent node; (b) the twin network for the model of (a) under $do(x^*)$.

tral role in the analysis of the spurious relations relative to the pair X, Y . The spurious covariance measures the correlation between the observational X and the counterfactual Y_{x^*} (Def. 2). We will employ in our analysis the twin network [Balke and Pearl, 1994; Pearl, 2000, Sec. 7.1.4], which is a graphical method to analyzing the relation between observational and counterfactual variables.

Consider the causal model M in Fig. 4(a), for example, where the exogenous variables $\{U_1, U_2\}$ are shown explicitly. Its twin network is the union of the model M (factual) and the submodel M_{x^*} (counterfactual) under intervention $do(x^*)$, which is shown in Fig. 4(b). The factual (M) and counterfactual (M_{x^*}) worlds share only the exogenous variables (in this case, U_1, U_2), which constitute the invariances shared across worlds. In this twin network, the observational X and the counterfactual Y_{x^*} are connected through two paths: one through U_1 and the other through U_2 . These paths correspond to two pathways from X to Y in the original causal diagram: $\tau_1 : X \leftarrow Z_2 \leftarrow Z_1 \leftarrow U_1 \rightarrow Z_1 \rightarrow Z_2 \rightarrow Y$, and $\tau_2 : X \leftarrow Z_2 \leftarrow U_2 \rightarrow Z_2 \rightarrow Y$.

Note that when considering the corresponding paths in the original graph (Fig. 4(a)), these paths (τ_1, τ_2) are not necessarily simple, i.e., they may contain a particular node more than once. Furthermore, each path can be partitioned into a pair of causal paths (say, g_l, g_r) from a common source $U_i \in U$ (e.g., τ_1 consists of a pair (g_{l_1}, g_{r_1}) , where $g_{l_1} : U_1 \rightarrow Z_1 \rightarrow Z_2 \rightarrow X$, and $g_{r_1} : U_1 \rightarrow Z_1 \rightarrow Z_2 \rightarrow Y$). Indeed, these non-simple paths are referred to as *treks* in the causal inference literature, which usually has been studied in the context of linear models [Spirtes *et al.*, 2001; Sullivant *et al.*, 2010].

Definition 9 (Trek). A trek τ in G (from X to Y) is an ordered pair of causal paths (g_l, g_r) with a common exogenous source $U_i \in U$ such that $g_l \in \Pi^c(U_i, X)$ and $g_r \in \Pi^c(U_i, Y)$. The common source U_i is called the top of the trek, denoted $top(g_l, g_r)$. A trek is spurious if $g_r \in \Pi^c(U_i, Y|X)$, i.e., g_r is a causal path from U_i to Y that is not intercepted by X .

We denote the set of treks from X to Y in G by $\mathcal{T}(X, Y)_G$ and spurious treks by $\mathcal{T}^s(X, Y)_G$ (G will be omitted when obvious). We introduce next an estimand for a specific spurious trek. For a spurious trek $\tau = (g_l, g_r)$ with $U_i = \text{top}(\tau)$, first let X_{g_l} denote the path-specific potential response $X_{g_l[U_i^l]}$, where U_i^l is an i.i.d. draw from the distribution $P(U_i)$. Similarly, let $Y_{x^*, g_r} = Y_{x^*, g_r[U_i^r]}$ ¹⁰, where $U_i^r \sim P(U_i)$. Pure trek-specific covariance can then finally be defined.

Definition 10 (Pure Trek-Specific Spurious Covariance). For a semi-Markovian model M and a spurious trek $\tau = (g_l, g_r)$ with $U_i = \text{top}(g_l, g_r)$, the pure τ -specific covariance of the treatment $X = x^*$ on the outcome Y is defined as:

$$\text{Cov}_{\tau[x^*]}^{ts}(X, Y) = \text{Cov}(X - X_{g_l}, Y_{x^*} - Y_{x^*, g_r}).$$

In words, the differences $X - X_{g_l}$ and $Y_{x^*} - Y_{x^*, g_r}$ are simply measuring the effects of the causal paths g_l and g_r (Lem. 1), while the $\text{Cov}(\cdot)$ operator is in charge of compounding them. (In the extreme case when g_l or g_r are disconnected, the pure τ -specific spurious covariance will equate to zero.) For example, the pure τ_1 -specific spurious covariance $\text{Cov}_{\tau_1[x^*]}^{ts}(X, Y)$ in Fig. 4(a) is

$$\text{Cov}(X - X_{g_{l_1}}, Y_{x^*} - Y_{x^*, g_{r_1}}). \quad (14)$$

Note that the counterfactuals $X_{g_{l_1}}$ and $Y_{x^*, g_{r_1}}$ assign the randomized interventions $do(U_1^l), do(U_1^r)$ to the paths g_{l_1}, g_{r_1} , respectively. By Def. 6, Eq. 14 is equal to:

$$\text{Cov}(X - X_{U_1^l}, Y_{x^*} - Y_{x^*, U_1^r}).$$

This quantity can be more easily seen through its graphical representation, see Fig. 5 (top). The main idea is to decompose U_1 into two independent components U_1^l, U_1^r (Fig. 5b), which is then contrasted with the world in which U_1 is kept intact (a).^{11 12} We note that by Def. 6, $X = X_\emptyset$ and $Y_{x^*} = Y_{x^*, \emptyset}$. The pure τ_1 -specific spurious covariance can be written as:

$$\text{Cov}_{\tau_1[x^*]}^{ts}(X, Y) = \text{Cov}(X_\emptyset - X_{g_{l_1}}, Y_{x^*, \emptyset} - Y_{x^*, g_{r_1}}).$$

More generally, the pure trek-specific spurious covariance for $\tau = (g_l, g_r)$ measures the covariance of variables $X_{\pi_l} - X_{\pi_l \cup \{g_l\}}$ and $Y_{x^*, \pi_r} - Y_{x^*, \pi_r \cup \{g_r\}}$, where $\pi_l (\pi_r)$ is an arbitrary set of causal paths from U that does not contain $g_l (g_r)$. This observational will be useful later on, which leads to the trek-specific spurious covariance.

¹⁰ $Y_{x^*, g_r[U_i^r]}$ is the g_r -specific potential response of Y to $do(g_r[U_i^r])$ in the submodel M_{x^*} .

¹¹This operation can be seen as the parallel to the pure path-specific covariance (Def. 7), with the distinct requirement that the replacement operator, used to generate the differences, is not relative to the observed X , but the corresponding U_i .

¹²To avoid clutter, Fig. 5 is a projected version of the original twin network focused on the relevant quantities (w.l.g.).

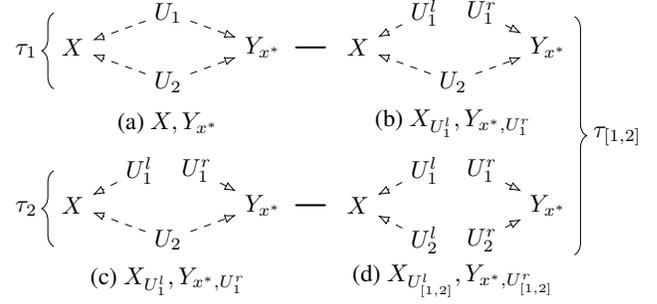


Figure 5: The decomposition procedure of the spurious covariance over the spurious treks τ_1, τ_2 (Thm. 3).

Definition 11 (Trek-Specific Spurious Covariance). For a semi-Markovian model M , let τ be a spurious trek (g_l, g_r) and π is a function mapping τ to a pair $\pi(\tau) = (\pi_l, \pi_r)$ where π_l and π_r are sets of causal paths from U such that $g_l \notin \pi_l$ and $g_r \notin \pi_r$. The τ -specific spurious covariance of the treatment $X = x^*$ on the outcome Y , denoted by $\text{Cov}_{\tau[x^*]}^{ts}(X, Y)_\pi$, is defined as

$$\text{Cov}(X_{\pi_l} - X_{\pi_l \cup \{g_l\}}, Y_{x^*, \pi_r} - Y_{x^*, \pi_r \cup \{g_r\}}).$$

The next proposition establishes the relationship between Def. 11 and the corresponding spurious treks. This property can be seen as a necessary condition for any measure of strength for spurious relations.

Property 6. $\tau \notin \mathcal{T}^s(X, Y) \Rightarrow \text{Cov}_{\tau[x^*]}^{ts}(X, Y)_\pi = 0$.

As an example of Def. 11, the trek τ_2 in Fig. 4(a) consists of paths $g_{l_2} : U_2 \rightarrow Z_2 \rightarrow X$ and $g_{r_2} : U_2 \rightarrow Z_2 \rightarrow Y$. If we set $\pi(\tau_2) = (\{g_{l_1}\}, \{g_{r_1}\})$, the τ_2 -specific spurious covariance can be measured by $\text{Cov}_{\tau_2[x^*]}^{ts}(X, Y)_\pi$, i.e.,

$$\text{Cov}(X_{g_{l_1}} - X_{g_{l_1, [2]}}, Y_{x^*, g_{r_1}} - Y_{x^*, g_{r_1, [2]}}) \quad (15)$$

$$= \text{Cov}(X_{U_1^l} - X_{U_{[1,2]}^l}, Y_{x^*, U_1^r} - Y_{x^*, U_{[1,2]}^r}). \quad (16)$$

Eq. 16 is graphically represented in Fig. 5(c-d), where the effect of the trek τ_2 is measured. In words, the difference between Fig. 5(c) and (d) is the effect of the causal paths g_{l_2} and g_{r_2} when U_2 is kept intact versus when divided into two independent components (U_2^l, U_2^r) .

Armed with the definition of trek-specific spurious covariance, we can finally study the decomposability of the spurious covariance $\text{Cov}_{x^*}^s(X, Y)$ (Def. 2). First, let $U^s \subseteq U$ denote the maximal set of exogenous variables that simultaneously affect variables X and Y_{x^*} (common exogenous ancestors), and let an order over U^s be $\mathcal{L}_u^s : U_1 < \dots < U_n$. For each $U_i \in U^s$, let $\mathcal{L}_{l_i}^s$ be an order $g_{l_1}^i < \dots < g_{l_n}^i$ over the set $\Pi^c(U_i, X)$. Similarly, we define $\mathcal{L}_{r_i}^s$ for $\Pi^c(U_i, Y|X)$. The tuple $\mathcal{L}^s = \langle \mathcal{L}_u^s, \{(\mathcal{L}_{l_i}^s, \mathcal{L}_{r_i}^s)\}_{1 \leq i \leq |U^s|} \rangle$ thus defines an order

for the spurious treks $\mathcal{T}^s(X, Y)$. We denote \mathcal{L}_π^s a function which maps from a trek τ to sets of paths $\mathcal{L}_\pi^s(\tau)$ covered by the spurious treks preceding τ in \mathcal{L}^s . Formally, given a spurious trek $\tau = (g_{l_j}^i, g_{r_k}^i)$, $\mathcal{L}_\pi^s(\tau)$ is equal to

$$(\Pi^c(U_{[1, i-1]}, X) \cup g_{l_{[1, j-1]}}^i, \Pi^c(U_{[1, i-1]}, Y|X) \cup g_{r_{[1, k-1]}}^i).$$

We are now ready to derive the decomposition formula for the spurious covariance $\text{Cov}_{x^*}^s(X, Y)$.

Theorem 3. *For a semi-Markovian model M , let $\mathcal{L}^s = \langle \mathcal{L}_u^s, \{(\mathcal{L}_{l_i}^s, \mathcal{L}_{r_i}^s)\}_{1 \leq i \leq |U^s|} \rangle$ be an order over spurious treks $\mathcal{T}^s(X, Y)$. For any x^* , the following non-parametric relationship hold:*

$$\text{Cov}_{x^*}^s(X, Y) = \sum_{\tau \in \mathcal{T}^s(X, Y)} \text{Cov}_{\tau[x^*]}^{ts}(X, Y) \mathcal{L}_\pi^s$$

For example, in the model of Fig. 4(a), $U^s = \{U_1, U_2\}$. τ_1 (τ_2) is the only spurious trek associated with U_1 (U_2). If we consider the order \mathcal{L}^s such that $\mathcal{L}_u^s : U_1 < U_2$, Thm. 3 dictates that $\text{Cov}_{x^*}^s(X, Y)$ should be decomposed as the sum of Eqs. 14 and 15. Fig. 5 shows the graphical representation of this decomposition procedure: we measure the change of the covariance between X and Y_{x^*} as we disconnect the relations going through τ_1 (associated with U_1) and τ_2 (U_2), sequentially. The sum of these changes thus equates to the correlations of X and Y along the spurious pathways, i.e., the spurious covariance $\text{Cov}_{[x^*]}^s(X, Y)$. (See [Zhang and Bareinboim, 2018b, Sec. 2] for more examples.)

6 NON-PARAMETRIC PATH ANALYSIS

In this section, we put the results of the previous sections together and derive a general path-specific decomposition for the covariance of the treatment X and the outcome Y without assuming any specific parametric form.

We start by noting that each spurious path from X to Y corresponds to a unique set of spurious treks that start on X and end in Y . Recall that a spurious path l can be seen as a pair of causal paths (g_l, g_r) , where the only node shared among g_l and g_r is the common source. For example, the spurious path $l : X \leftarrow Z_2 \rightarrow Y$ is a pair (g_l, g_r) such that $g_l : Z_2 \rightarrow X$ and $g_r : Z_2 \rightarrow Y$. We can thus define a rule f which maps a trek $\tau \in \mathcal{T}^s(X, Y)$ to a spurious path $l \in \Pi^s(X, Y)$. For $\tau = (g_l, g_r)$, let V_t be the most distant recurring node from $\text{top}(g_l, g_r)$ such that V_t is the only node shared among subpaths $g_l(V_t, X)$ and $g_r(V_t, Y)$; the pair $(g_l(V_t, X), g_r(V_t, Y))$ corresponds to a path l in $\Pi^s(X, Y)$. As an example, the trek τ_1 in Fig. 4(a) has $V_t = Z_2$, which corresponds to the spurious path $l : X \leftarrow Z_2 \rightarrow Y$, and similarly, $f(\tau_1) = l$ as well as $f(\tau_2) = l$. Lem. 2 shows that the rule f forms a valid surjective function.

Lemma 2. *For a semi-Markovian model M , for each spurious trek $\tau \in \mathcal{T}^s(X, Y)$, there always exists a unique most distant recurring node V_t .*

For a spurious path l , let $\mathcal{T}^s(l) = f^{-1}(l)$ denote its corresponding treks. Specifically, if $l \notin \Pi^s(X, Y)$, then for each $\tau \in \mathcal{T}^s(l)$, we must have $\tau \notin \mathcal{T}^s(X, Y)$. For instance, if the spurious l in Fig. 4(a) is disconnected, e.g., $Z_2 \not\rightarrow X$, treks τ_1, τ_2 are both disconnected as well. From this observation, we could naturally define the spurious covariance of a path l as a sum over treks in $\mathcal{T}^s(l)$.

Definition 12 (Path-Specific Spurious Covariance). For a semi-Markovian model M with an associated causal diagram G , let l be an arbitrary spurious path in G . Let π be a function that maps a trek $\tau = (g_l, g_r) \in \mathcal{T}^s(l)$ to a pair $\pi(\tau) = (\pi_l, \pi_r)$, where π_l and π_r are arbitrary sets of causal paths from U such that $g_l \notin \pi_l$ and $g_r \notin \pi_r$. The l -specific spurious covariance of the treatment $X = x^*$ on the outcome Y is defined as

$$\text{Cov}_{l[x^*]}^s(X, Y)_\pi = \sum_{\tau \in \mathcal{T}^s(l)} \text{Cov}_{\tau[x^*]}^{ts}(X, Y)_\pi$$

Property 7. $l \notin \Pi^s(X, Y) \Rightarrow \text{Cov}_{l[x^*]}^s(X, Y)_\pi = 0$.

The surjectivity of the function f assures that the set $\{\mathcal{T}^s(l)\}_{l \in \Pi^s(X, Y)}$ forms a partition over the spurious treks $\mathcal{T}^s(X, Y)$. From Thm. 3, it follows immediately that the path-specific spurious covariance (Def. 12) has the property that expresses the spurious covariance $\text{Cov}_{x^*}^s(X, Y)$ as a sum over $\Pi^s(X, Y)$.

Theorem 4. *For a semi-Markovian model M , let $\mathcal{L}^s = \langle \mathcal{L}_u^s, \{(\mathcal{L}_{l_i}^s, \mathcal{L}_{r_i}^s)\}_{1 \leq i \leq |U^s|} \rangle$ be an order over spurious treks $\mathcal{T}^s(X, Y)$. For any x^* , the following non-parametric relationship hold:*

$$\text{Cov}_{x^*}^s(X, Y) = \sum_{l \in \Pi^s(X, Y)} \text{Cov}_{l[x^*]}^s(X, Y) \mathcal{L}_\pi^s$$

As an example, the path $l : X \leftarrow Z_2 \rightarrow Y$ in Fig. 4(a) corresponds to $\mathcal{T}^s(l) = \{\tau_1, \tau_2\}$. For an arbitrary order \mathcal{L}^s , Thm. 4 is applicable and immediately yields $\text{Cov}_{x^*}^s(X, Y) = \text{Cov}_{l[x^*]}^s(X, Y) \mathcal{L}_\pi^s$, which means that the path l accounts for all the spurious relations between X and Y . In other words, the spurious joint variability of X and Y is fully explained by the variance of Z_2 , which is a function of the exogenous variables U_1 (through τ_1) and U_2 (through τ_2).

Thms. 1-2 and 4 together lead to a general path-specific decomposition formula, which allows one to non-parametrically decompose the covariance $\text{Cov}(X, Y)$ over all open paths from X to Y in the underlying model.

Theorem 5 (Path-Specific Decomposition). *For a semi-Markovian model M , let \mathcal{L}^c be an order over $\Pi^c(X, Y)$*

and $\mathcal{L}^s = \langle \mathcal{L}^u, \{(\mathcal{L}_{l_i}^s, \mathcal{L}_{r_i}^s)\}_{1 \leq i \leq |U^s|} \rangle$ be an order over $\mathcal{T}^s(X, Y)$. For any x^* , the following non-parametric relationship hold:

$$\begin{aligned} \text{Cov}(X, Y) &= \sum_{l \in \Pi^c(X, Y)} \text{Cov}_{l[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c} \\ &+ \sum_{l \in \Pi^s(X, Y)} \text{Cov}_{l[x^*]}^s(X, Y)_{\mathcal{L}_\pi^s}. \end{aligned} \quad (17)$$

We illustrate the use of Thm. 5 using the model discussed in Sec. 1 (Fig. 1(a)). Recall that X and Y are connected through the causal paths l_1, l_2 and spurious paths l_3, l_4 . Note that $U^s = \{U_Z\}$ spuriously affects the treatment X through the path $g_l = U_Z \rightarrow Z \rightarrow X$, and the outcome Y through the paths $g_{r_1} = U_Z \rightarrow Z \rightarrow Y$ and $g_{r_2} = U_Z \rightarrow Z \rightarrow W \rightarrow Y$. Let order \mathcal{L}^c be $l_1 < l_2$ and \mathcal{L}^s be $g_{r_1} < g_{r_2}$. For any level x^* , Thm. 5 equates the covariance $\text{Cov}(X, Y)$ to the sum of $\{\text{Cov}_{l_i[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c}\}_{i=1,2}$ and $\{\text{Cov}_{l_i[x^*]}^s(X, Y)_{\mathcal{L}_\pi^s}\}_{i=3,4}$, which can be written as

$$\begin{aligned} &\underbrace{\text{Cov}(X, Y - Y_{x^*, W})}_{l_1: X \rightarrow Y} + \underbrace{\text{Cov}(X, Y_{x^*, W} - Y_{x^*})}_{l_2: X \rightarrow W \rightarrow Y} \\ &+ \underbrace{\text{Cov}(X - X_{U_Z^l}, Y_{x^*} - Y_{x^*, W_{x^*} Z_{U_Z^r}})}_{l_3: X \leftarrow Z \rightarrow Y} \\ &+ \underbrace{\text{Cov}(X - X_{U_Z^l}, Y_{x^*, W_{x^*} Z_{U_Z^r}} - Y_{x^*, U_Z^r})}_{l_4: X \leftarrow Z \rightarrow W \rightarrow Y}, \end{aligned} \quad (18)$$

which are all well-defined, computable from the structural causal model [Def. 1; Pearl, 2000, Sec. 7.1].

7 IDENTIFYING PATH-SPECIFIC DECOMPOSITION

By and large, identifiability is one of the most studied topics in causal inference. It is acknowledged in the literature that obtaining identifiability may be non-trivial even in the context of less granular measures of causal effects, including quantities without nested counterfactual and following the analysis of Pearl's do-calculus.

In this section, we start the study of identifiability conditions for when the path-specific decomposition formula (Thm. 5) can be estimated from data, when the SCM is not fully known. We'll analyze the causal model discussed in Fig. 1(a) given its generality and potential to encode more complex models. The main assumption encoded in this model is Markovianity, i.e., that all exogenous variables are independent. We show next that identifiability can be obtained under these assumptions.

Theorem 6. *The path-specific decomposition of Eq. 18 is identifiable if the distributions $P(x, y_{x^*}), P(x, y_{x^*, W})$ and $P(x, y_{x^*, W_{x^*}, Z_{U_Z^r}})$ are identifiable. Specifically, in the model of Fig. 1(a), $P(x, y_{x^*}), P(x, y_{x^*, W})$, and*

$P(x, y_{x^, W_{x^*}, Z_{U_Z^r}})$ can be estimated, respectively, from the observational distribution $P(x, y, z, w)$ as follows:*

$$P(x, y_{x^*}) = \sum_{z, w} P(y|x^*, w, z)P(w|x^*, z)P(x, z)$$

$$P(x, y_{x^*, W}) = \sum_{z, w} P(y|x^*, z, w)P(x, z, w)$$

$$P(x, y_{x^*, W_{x^*}, Z_{U_Z^r}}) = \sum_{z, z', w} P(y|x^*, z, w)P(w|x^*, z')P(x, z')P(z)$$

Note that all the quantities listed in Thm. 6 are expressible in terms of conditional distributions and do not involve any counterfactual (simple nor nested), which are readily estimable from the observational distribution. As an example, the l_2 -specific causal covariance $\text{Cov}_{l_2[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c}$ in Eq. 18 can be written as $\text{Cov}(X, Y_{x^*, W}) - \text{Cov}(X, Y_{x^*})$, which are computed from the counterfactual distributions $P(x, y_{x^*})$ and $P(x, y_{x^*, W})$, respectively. These distributions can be estimated from the observational distribution $P(x, y, z, w)$ following Thm. 6. Indeed, the path-specific decomposition formula (Thm. 5) is identifiable in the model of Fig. 1(a) regardless of the order \mathcal{L}^c and \mathcal{L}^s . (For other decompositions, see [Zhang and Bareinboim, 2018b].)

We further considered the identifiability conditions for the path-specific decomposition formula when the more stringent assumption that the underlying structural functions are linear is imposed.

Theorem 7. *Under the assumption of linearity and the assumption of Fig. 1(a), for any arbitrary orders \mathcal{L}^c and \mathcal{L}^s , for any x , the path-specific covariance of l_1, l_2, l_3 and l_4 are equal to:*

$$\text{Cov}_{l_1[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c} = \alpha_{YX}, \quad \text{Cov}_{l_2[x^*]}^c(X, Y)_{\mathcal{L}_\pi^c} = \alpha_{WX}\alpha_{YW}$$

$$\text{Cov}_{l_3[x^*]}^s(X, Y)_{\mathcal{L}_\pi^s} = \alpha_{XZ}\alpha_{YZ}, \quad \text{Cov}_{l_4[x^*]}^s(X, Y)_{\mathcal{L}_\pi^s} = \alpha_{XZ}\alpha_{WZ}\alpha_{YW}$$

The parameters α can be estimated from the corresponding (partial) regression coefficients [Pearl, 2000, Ch. 5].

Clearly, after applying Thm. 7 to Eq. 18, the resulting decomposition coincides with Wright's method of path coefficients in the linear-standard model (Eq. 1).

8 CONCLUSIONS

We introduced novel covariance-based counterfactual measures to account for effects along with a specific path from a treatment X to an outcome Y (Defs. 8, 11-12). We developed machinery to allow, for the first time, the non-parametric decomposition of the covariance of X and Y as a summation over the different pathways in the underlying causal model (Thm. 5). We further provided identification conditions under which the decomposition formula can be estimated from data (Thm. 6-7).

Acknowledgments

Bareinboim and Zhang are supported in parts by grants from NSF IIS-1704352 and IIS-1750807 (CAREER).

References

- C. Avin, I. Shpitser, and J. Pearl. Identifiability of path-specific effects. In *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence IJCAI-05*, pages 357–363, Edinburgh, UK, 2005. Morgan-Kaufmann Publishers.
- A. Balke and J. Pearl. Counterfactual probabilities: Computational methods, bounds, and applications. In R. Lopez de Mantaras and D. Poole, editors, *Uncertainty in Artificial Intelligence 10*, pages 46–54. Morgan Kaufmann, San Mateo, CA, 1994.
- E. Bareinboim and J. Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113:7345–7352, 2016.
- Reuben M Baron and David A Kenny. The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of personality and social psychology*, 51(6):1173, 1986.
- K.A. Bollen. *Structural Equations with Latent Variables*. John Wiley, NY, 1989.
- RM Daniel, BL De Stavola, SN Cousens, and Stijn Vansteelandt. Causal mediation analysis with multiple mediators. *Biometrics*, 71(1):1–14, 2015.
- O.D. Duncan. *Introduction to Structural Equation Models*. Academic Press, New York, 1975.
- J. Fox. Effect analysis in structural equation models. *Sociological Methods and Research*, 9(1):3–28, 1980.
- Kosuke Imai, Luke Keele, and Teppei Yamamoto. Identification, inference and sensitivity analysis for causal mediation effects. *Statist. Sci.*, 25(1):51–71, 02 2010.
- Kosuke Imai, Luke Keele, Dustin Tingley, and Teppei Yamamoto. Unpacking the black box of causality: Learning about causal mechanisms from experimental and observational studies. *American Political Science Review*, 105(4):765–789, 2011.
- Niki Kilbertus, Mateo Rojas Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*, pages 656–666, 2017.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pages 4069–4079, 2017.
- Xintao Wu Lu Zhang, Yongkai Wu. A causal framework for discovering and removing direct and indirect discrimination. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pages 3929–3935, 2017.
- D.P. MacKinnon. *An Introduction to Statistical Mediation Analysis*. Lawrence Erlbaum Associates, New York, 2008.
- J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000.
- J. Pearl. Direct and indirect effects. In *Proc. of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, pages 411–420. Morgan Kaufmann, CA, 2001.
- Ilya Shpitser and Eric Tchetgen Tchetgen. Causal inference with a graphical hierarchy of interventions. *Annals of statistics*, 44(6):2433, 2016.
- Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press, 2001.
- Seth Sullivant, Kelli Talaska, and Jan Draisma. Trek separation for gaussian graphical models. *The Annals of Statistics*, pages 1665–1685, 2010.
- Eric J Tchetgen Tchetgen and Ilya Shpitser. Semiparametric theory for causal mediation analysis: efficiency bounds, multiple robustness, and sensitivity analysis. *Annals of statistics*, 40(3):1816, 2012.
- Tyler J VanderWeele, Stijn Vansteelandt, and James M Robins. Effect decomposition in the presence of an exposure-induced mediator-outcome confounder. *Epidemiology*, 25(2):300, 2014.
- Tyler VanderWeele. *Explanation in Causal Inference: Methods for Mediation and Interaction*. Oxford University Press, New York, 2015.
- Stijn Vansteelandt and Tyler J VanderWeele. Natural direct and indirect effects on the exposed: effect decomposition under weaker assumptions. *Biometrics*, 68(4):1019–1027, 2012.
- S. Wright. The theory of path coefficients: A reply to Niles’ criticism. *Genetics*, 8:239–255, 1923.
- Sewall Wright. The method of path coefficients. *The annals of mathematical statistics*, 5(3):161–215, 1934.
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th International Conference on WWW*, pages 1171–1180, 2017.
- Junzhe Zhang and Elias Bareinboim. Fairness in decision-making – the causal explanation formula. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, 2018a.
- Junzhe Zhang and Elias Bareinboim. Non-parametric path analysis in structural causal models. Technical Report R-34, AI Lab, Purdue University., 2018b.