

---

# On Measuring Causal Contributions via *do*-interventions

---

Yonghan Jung<sup>1</sup> Shiva Prasad Kasiviswanathan<sup>2</sup> Jin Tian<sup>3</sup>  
Dominik Janzing<sup>2</sup> Patrick Blöbaum<sup>2</sup> Elias Bareinboim<sup>4</sup>

## Abstract

Causal contributions measure the strengths of different causes to a target quantity. Understanding causal contributions is important in empirical sciences and data-driven disciplines since it allows to answer practical questions such as “what are the contributions of each cause to the effect?” In this paper, we develop a principled method for quantifying causal contributions. First, we provide desiderata in terms of axioms that causal contribution measures should satisfy and propose the *do*-Shapley values (inspired by *do*-interventions (Pearl, 2000)) as a method satisfying these properties. Next, we develop a criterion under which the *do*-Shapley values can be efficiently inferred from non-experimental data. Finally, we introduce *do*-Shapley estimators exhibiting consistency and statistical robustness. Simulation results corroborate with the proposed theory.

## 1. Introduction

Inferring causal effects is a fundamental problem throughout the data sciences since it can answer queries like “what would be the expected outcome if inputs had been fixed to certain values?”. There is a growing literature that investigates the conditions under which causal conclusions can be drawn from observational and experimental data (causal effect identification) (Pearl, 1995; Tian & Pearl, 2003; Huang & Valtorta, 2006; Shpitser & Pearl, 2006; Bareinboim & Pearl, 2012; 2016; Jaber et al., 2018; Lee et al., 2019; 2020; Lee & Bareinboim, 2020), and in estimating the identified causal functions from data (causal effect estimation) (Jung et al., 2020; Bhattacharya et al., 2020; Jung et al., 2021a;b; Bhattacharyya et al., 2020; 2021; Xia et al., 2021). Beyond these tasks, interpreting the results of causal inference, including answering “what is the most important cause of the

effects?”, or more generally, “what are the contributions of each cause to the effect?” are also of practical importance. Answering these queries falls under the task of measuring *causal contributions*, which aims to quantify the degree of contribution of different causes to a target effect. As a motivational example, consider the following scenario described in (Lundberg, 2021):

**Example 1.** *A video streaming service company has collected data that contains various features including sales call ( $S$ ), product needs ( $P$ ), interaction with customers ( $I$ ), monthly usage ( $M$ ), discounts provided ( $D$ ), last upgrade ( $L$ ), economic factors ( $E$ ), ad spend ( $A$ ), and bugs reported ( $B$ ). These features are causally related and affect the outcome: customers retention ( $Y$ ) (also, for further details, see Fig. 1). The company aims to measure the causal contributions of these features to the target effect – the expected customer retention if each feature had been fixed to a certain value (e.g., set to lower sales calls, higher product needs, etc.).*

Example 1 captures practical cases where the target quantity is related to the query “what would be the output if inputs had been fixed to certain values?”. This includes cases where the target quantity is a machine learning (ML) model’s output, which is derived by fixing inputs to specific values (see Remark 1).

In the area of explainable AI (XAI), there is a series of works concerned with measuring the contributions of features to an ML model output (Lundberg & Lee, 2017; Schwab & Karlen, 2019; Janzing et al., 2020b; Heskes et al., 2020; Covert et al., 2021). Most of these methods have focused on queries where the target quantity is induced from an *accessible* model – a model for target  $Y$  is said to be *accessible* if the model can be evaluated to obtain  $Y$  value for arbitrary input features – with less attention has been paid to settings where the target is induced by nature (i.e., the data-generating process is *inaccessible*, such as the customer retention in Example 1). Also, many existing techniques are based on the correlation between the features and the ML model output, including, (Lundberg & Lee, 2017; Frye et al., 2020) to cite a few. Another thread of this approach focused on measuring contributions based on causation (Schwab & Karlen, 2019; Janzing et al., 2020b; Heskes et al., 2020),

---

<sup>1</sup>Purdue University <sup>2</sup>Amazon <sup>3</sup>Iowa State University  
<sup>4</sup>Columbia University. Correspondence to: Yonghan Jung  
<jung222@purdue.edu>.

but often assumes that the data generating process for the target is known and accessible, allowing that an outcome corresponding to any arbitrary features can be generated. This rules out scenarios where the target quantity is induced from an inaccessible model (such as Example 1). A more detailed comparison with the existing literature is presented in Sec. 3.1.

In this paper, we generalize previous approaches to measure the causal contributions of each feature to a target effect induced by an inaccessible model and described by a joint interventional distribution. Our proposed method is applicable to the task of quantifying causal contributions of input features of an ML model prediction as well. More specifically, our contributions are as follows:

1. [Sec. 3] We *axiomatize* causal contribution measures. Specifically, we propose desiderata for causal contribution measures (a set of axioms), and introduce the *do*-Shapley, a Shapley value-based method (Shapley, 1953) specialized for quantifying the causal contributions described by *do*-interventions (Pearl, 2000).<sup>1</sup> Our axiomatic characterization provides a theoretical justification for using the *do*-Shapley for quantifying causal contributions.

2. [Sec. 4] We provide conditions under which the *do*-Shapley values can be inferred from observational data (*identifiability*) in polynomial time. Even if verifying the identifiability can be done through existing causal-effect identification results, determining the identifiability of *do*-Shapley values is, in practice, not computationally feasible. In particular, we introduce sufficient conditions under which the identifiability of *do*-Shapley values is determined in polynomial time.

3. [Sec. 5] We develop estimators for the *do*-Shapley values, exhibiting consistency, identifiability, and statistical robustness. We developed three estimators based on the *inverse probability weighting* (IPW) (Rosenbaum & Rubin, 1983), *outcome regression* (REG) (Rubin, 1979), and *double/debiased machine learning* (DML) (Chernozhukov et al., 2018), respectively. We also show that the DML estimator displays statistical robustness to model misspecification and bias.

4. [Sec. 6]. Finally, we present simulation results on these estimators that corroborate with the theory.

Due to space constraints, the proofs and other omitted details are provided in the appendix.

<sup>1</sup>The *do*-Shapley is a generalization of the *causal Shapley* (Heskes et al., 2020), which also uses the *do*-interventions to the case where the target quantity is induced by an inaccessible model.

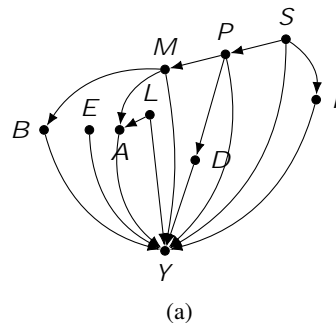


Figure 1: Causal graph for Example 1 (Lundberg, 2021).

## 2. Preliminaries

**Notation.** Each variable is represented with a capital letter ( $V$ ) and its realized value with a small letter ( $v$ ). We use bold letters  $\mathbf{V}$  and  $\mathbf{v}$  to denote a set of variables and their realized value, respectively. For any set  $S$ , we use  $|S|$  to denote its cardinality. Given a topological order over the vertices  $\mathbf{V} := \{V_1, \dots, V_n\}$  of a graph  $G$ , we will use  $\text{pre}(V_i)$  to denote the predecessors of  $V_i$  and use  $\text{pre}(v_i)$  to denote a realization of  $\text{pre}(V_i)$ ; i.e.,  $\text{pre}(v_i) = \mathbf{w}_i$  for  $\text{pre}(V_i) = \mathbf{W}_i$ . We use  $\text{Ch}(V_i)$  to represent the children of a variable  $V_i$  in  $G$ . For an index set  $[n] := \{1, \dots, n\}$  and a subset  $S \subseteq [n]$ , we use  $\mathbf{V}_S := \{V_k \mid k \in S\}$  and  $\mathbf{V}_{\bar{S}} := \{V_k \mid k \notin S\}$ . We use  $D$  to denote  $N$  samples from a distribution  $P$  over  $\mathbf{V}$ ; i.e.,  $D := \{V_{(i)}\}_{i=1}^N \sim P$ , where  $V_{(i)}$  denotes the  $i$ th sample. For a function  $f$ , we use  $\mathbb{E}[f(\mathbf{V})]$  as an expectation of  $f(\mathbf{V})$  over  $P$ , and  $\mathbb{E}_D[f(\mathbf{V})] := (1/N) \sum_{i=1}^N f(V_{(i)})$ . We use  $\|f(\mathbf{V})\|_2 := \sqrt{\mathbb{E}[(f(\mathbf{V}))^2]}$  to denote the  $L_2(P)$  norm of  $f(\mathbf{V})$ .  $O_P(\cdot)$  and  $o_P(\cdot)$  denotes the big O and little O in probability, respectively.

**Structural Causal Models.** We use the language of structural causal models (SCMs) as our basic semantical framework (Pearl, 2000). A structural causal model (SCM) is a tuple  $\mathcal{M} := \langle \mathbf{V}, \mathbf{U}, \mathbf{F}, P(\mathbf{u}) \rangle$ , where  $\mathbf{V}, \mathbf{U}$  are sets of endogenous (observables) and exogenous variables (latents) respectively,  $\mathbf{F}$  is a set of functions  $f_{V_i}$  one for each  $V_i \in \mathbf{V}$  where  $V_i = f_{V_i}(PA_{V_i}, U_{V_i})$  for some  $PA_{V_i} \subseteq \mathbf{V}$  and  $U_{V_i} \subseteq \mathbf{U}$ , and  $P(\mathbf{u})$  is a strictly positive probability measure for  $\mathbf{U}$ . Each SCM  $\mathcal{M}$  is associated to a causal diagram  $G$  over the node set  $\mathbf{V}$  where  $V_i \rightarrow V_j$  if  $V_i$  is an argument of  $f_{V_j}$ , and  $V_i \perp\!\!\!\perp V_j$  if the corresponding  $U_{V_i}$  and  $U_{V_j}$  are not independent. Performing an intervention  $\mathbf{X} = \mathbf{x}$  is represented through the *do*-operator,  $do(\mathbf{X} = \mathbf{x})$  (shortly,  $do(\mathbf{x})$ ), which encodes the operation of replacing the original equations of  $\mathbf{X}$  by the constant  $\mathbf{x}$  in the SCM  $\mathcal{M}$ , inducing a submodel  $\mathcal{M}_{\mathbf{x}}$  and an interventional distribution  $P(\mathbf{V} = \mathbf{v} \mid do(\mathbf{x}))$  (shortly,  $P(\mathbf{v} \mid do(\mathbf{x}))$ ) (Bareinboim et al., 2020).

**Causal Effect Identification.** Given a causal graph  $G$

over  $\mathbf{V}$ , an effect  $P(yjdo(\mathbf{x}))$  where  $\mathbf{X}, \mathbf{Y} \subseteq \mathbf{V}$  is *identifiable* if  $P(yjdo(\mathbf{x}))$  is computable from the distribution  $P(\mathbf{v})$  in any SCM  $\mathcal{M}$  that induces  $G$  (Pearl, 2000, p. 77). One key notion for identification is *confounded components* (in short, *C-components*): a set of nodes connected with a path composed solely of bi-directed edges  $V_i \leftrightarrow V_j$  (Tian & Pearl, 2003). For any  $\mathbf{C} \subseteq \mathbf{V}$ , the quantity  $Q[\mathbf{C}] := P(cjdo(\mathbf{v}nc))$ , called a *C-factor* (Tian & Pearl, 2003), is defined as an interventional distribution of  $\mathbf{C}$  under an intervention on  $\mathbf{V} \setminus \mathbf{C}$ . We use  $C(V_i)_G$  (shortly,  $C(V_i)$ ) to denote the *C-component* of  $V_i$  in  $G$ , the set of variables belonging to the same *C-component* as  $V_i$ . We use  $C(\mathbf{W}) := \bigcup_{V_i \in \mathbf{W}} C(V_i)$  to denote the *C-component* of a set  $\mathbf{W} \subseteq \mathbf{V}$ .

**Shapley Value.** The Shapley value (Shapley, 1953) seeks to allocate the contribution of each player  $i \in [n]$  on some function value  $f([n])$  given a value function  $\nu(S)$  that measures the value of coalition of players  $i \in S \subseteq [n]$ . The Shapley value is given by

$$\phi_i(\nu) := \sum_{S \subseteq [n], i \in S} \omega(S) \frac{f(S \cup \{i\}) - f(S)}{|S|+1}, \quad (1)$$

where  $\omega(S) := \frac{1}{n} \binom{n-1}{|S|}$ . The Shapley value is uniquely satisfying a set of some desiderata for fair allocation (Young, 1985) (See Appendix A for more details).

## 2.1. Problem Definition

We are given samples  $D$  drawn from a distribution  $P := P_{\mathcal{M}}$  and a compatible causal diagram  $G := G_{\mathcal{M}}$ , induced by the SCM  $\mathcal{M}$ , on topologically ordered variables  $(\mathbf{V}, Y)$ , with  $Y$  the final variable in the order. We assume that  $Y$  is bounded, and  $\mathbf{V}$  is a set of discrete variables.

Given  $(G, D, \mathbf{v})$  where  $\mathbf{v}$  is a realization of  $\mathbf{V}$ , the task is to measure the contribution of each  $v_i \in \mathbf{v}$  to the target causal effect<sup>2</sup>  $E[Yjdo(\mathbf{v})]$  based on their impact on  $Y$  if the SCM  $\mathcal{M}$  has fixed the value of the variable as  $V_i = v_i$ . We set  $E[Y] = 0$  without loss of generality. We make no assumptions regarding the data generating process of  $Y$  for generality. With the following additional assumption on  $f$ , our problem is reduced to the problem of attributing the importance of features in an ML prediction model:

**Remark 1 (Reduction to ML models).** *If  $Y$  is generated by a deterministic function, e.g.,  $Y$  is an output of an ML prediction model  $f$  s.t.  $Y := f(\mathbf{V})$ , then our task reduces to measuring the causal contribution of each features  $v_i \in \mathbf{v}$*

<sup>2</sup>We focus on the average causal effect  $E[Yjdo(\mathbf{v})]$ , one of the most commonly used quantities in practice. Our method is applicable for any function of the causal distribution  $P(yjdo(\mathbf{v}))$ . A condition whether the target quantity  $E[Yjdo(\mathbf{v})]$  (or  $P(yjdo(\mathbf{v}))$ ) can be determined using non-experimental data is discussed in Sec. 4.

on the ML prediction  $f(\mathbf{v})$ , since  $E[Yjdo(\mathbf{v})] = f(\mathbf{v})$ , the deterministic function.

## 3. Axioms for Causal Contribution

We start by asking the question: “What would be a good measure for causal contribution?” To answer this question, we propose the following desiderata inspired by previous works (Young, 1985; Friedman & Moulin, 1999; Sundararajan et al., 2017; Sundararajan & Najmi, 2020):

**Axiom 1 (Desiderata for Causal Contribution).** Causal contributions  $f\phi_{v_i}g_{i=1}^n$  is considered desirable if the following properties are satisfied:

1. **Perfect assignment:** The contributions are perfectly assigned; formally,  $E[Yjdo(\mathbf{v})] = \sum_{v_i \in \mathbf{v}} \phi_{v_i}$ .

2. **Causal irrelevance:** If  $V_i$  is causally irrelevant to  $Y$  for all witness  $\mathbf{w} \subseteq \mathbf{V} \setminus v_i$ ; formally,  $\delta y, P(yjdo(v_i, \mathbf{w})) = P(yjdo(\mathbf{w}))$ <sup>3</sup>, then  $\phi_{v_i} = 0$ .

3. **Causal symmetry:** If  $(v_i, v_j) \subseteq \mathbf{v}$  have the same causal explanatory power<sup>4</sup> to  $Y$ ,  $\delta \mathbf{w} \subseteq \mathbf{V} \setminus v_i, v_j, g(\delta y, P(yjdo(v_i, \mathbf{w})) = P(yjdo(v_j, \mathbf{w})))$ , then  $\phi_{v_i} = \phi_{v_j}$ .

4. **Causal approximation:** For any  $S \subseteq [n]$  and  $\mathbf{v}_S := f_{v_i \in S}g_{i \in S}$ ,  $\sum_{i \in S} \phi_{v_i}$  well approximates  $E[Yjdo(\mathbf{v}_S)]$ . Formally,  $f\phi_{v_i}g_{i=1}^n$  is a solution to the following weighted least square; i.e.,  $f\phi_{v_i}g_{i=1}^n = \arg \min_{\phi_{v_i}} \sum_{S \subseteq [n]} (E[Yjdo(\mathbf{v}_S)] - \sum_{i \in S} \phi_{v_i})^2 \omega(S)$  for some positive and bounded function  $\omega(S)$ .

The rationale behind Axioms 1 is the following: **(1) Perfect assignment** is a natural requirement since we aim to attribute the degree of contributions of each feature  $v_i \in \mathbf{v}$  to the target causal effect. **(2) Causal irrelevance** reflects a desire to understand the cause of the outcome by forcing zero contributions for variables not causing the outcome. **(3) Causal symmetry** enforces the equal contribution for a pair of features if they have the same causal explanatory power. **(4) Causal approximation** allows  $\phi_{v_i}$  to be interpreted as a proxy for the causal effect s.t.  $\sum_{i \in S} \phi_{v_i} \approx E[Yjdo(\mathbf{v}_S)]$  for any  $S \subseteq [n]$ .

Perhaps surprisingly, there is a unique causal contribution measure  $f\phi_{v_i}g$  satisfying the above four properties.

**Definition 1 (do-Shapley).** The *do-Shapley*<sup>5</sup> is a causal contribution measure  $f\phi_{v_i}g_{i=1}^n$  of  $\mathbf{v}$  on  $E[Yjdo(\mathbf{v})]$  w.r.t.  $G$

<sup>3</sup> $V_i$  is causally irrelevant to  $Y$  given  $\mathbf{V}_S$  if  $P(yjdo(v_i; \mathbf{v}_S)) = P(yjdo(\mathbf{v}_S))$  (Galles & Pearl, 1997, Def. 7).

<sup>4</sup>A causal explanatory power of  $\mathbf{X} = \mathbf{x}$  to  $Y = y$  is a measure of making  $Y = y$  ‘more likely’ if  $\mathbf{X}$  had been fixed to  $\mathbf{x}$ ; i.e.,  $P(yjdo(\mathbf{x})) > P(y)$  (Eva & Stern, 2019).

<sup>5</sup>Heskes et al. (2020) proposed the same equation for measuring contributions in the accessible model setting and referred to as the *causal Shapley*. In this paper, we use the term *do-Shapley* to make it clearer that the definition is based on the *do*-intervention.

defined as:

$$\phi_{v_i} := \sum_{S \subseteq [n] \setminus \{i\}} \omega(S) f \mathbb{E}[Y \text{do}(\mathbf{v}_{S,i})] - \mathbb{E}[Y \text{do}(\mathbf{v}_S)] g, \quad (2)$$

where  $\omega(S) := (1/n) \binom{n-1}{|S|}^{-1}$ .

**Theorem 1 (Uniqueness of the *do*-Shapley).** *The *do*-Shapley is a unique causal contribution measure satisfying all the properties in Axiom 1.*

**Remark 2.** *Thm. 1 is significant because Axiom 1 do not restrict the value function to any fixed form. Thm. 1 instead characterizes the *do*-Shapley as the unique causal contribution measure satisfying Axiom 1 among any arbitrary value functions and corresponding contribution measures, as in (Sundararajan & Najmi, 2020).*

The *do*-Shapley, as the name suggests, is a specialization of the Shapley value in Eq. (1) for  $\nu(S) = \mathbb{E}[Y \text{do}(\mathbf{v}_S)]$ . The *do*-Shapley can be alternatively viewed as a marginal causal effect of  $v_i \in \mathbf{v}$  (i.e.,  $\mathbb{E}[Y \text{do}(\mathbf{v}_{S,i})] - \mathbb{E}[Y \text{do}(\mathbf{v}_S)]$ ) weighted-averaging over a set  $S$ . The significance of Thm. 1 stems from that it codifies the guarantees of the *do*-Shapley, and provides a tool to compare and contrast with alternative contribution metrics.

**Remark 3 (Attribution of contributions for a subset of variables).** *It is worth noting that the *do*-Shapley allocates contributions to all  $v_i \in \mathbf{v}$  from a joint interventional distribution. In practice, assigning contributions exclusively to a subset  $\mathbf{x} \subseteq \mathbf{v}$  may lead to more interpretable results. For example, when  $\mathbf{X} := \text{Pa}(Y) \subseteq \mathbf{V}$ , assigning contributions only to the features in  $\mathbf{x} \subseteq \mathbf{v}$  might be more interpretable if it is needed that features indirectly affecting the outcome should be assigned zero contributions. Enforcing the *do*-Shapley to assign contributions only for the subset  $\mathbf{x}$  can be simply done (without loss of generality) by the following procedure: (1) Derive a causal graph  $G[\mathbf{X}]$  compatible with  $P(\mathbf{x})$  by applying the projection of a graph<sup>6</sup>; and (2) Compute the *do*-Shapley w.r.t.  $G[\mathbf{X}]$ . See Appendix B for more details.*

### 3.1. Relation with Other Work

In this section, we compare the *do*-Shapley in Def. 1 with other known methods aiming to measure contributions of features on the outcome. Table 1 summarizes the comparison.

**Conditional Shapley.** The *conditional Shapley* ( $\phi_{V_i}^{\text{cond}}$ ) is a specialization of the Shapley value with  $\nu(S) = \mathbb{E}[Y \text{do}(\mathbf{v}_S)]$

<sup>6</sup> $G[\mathbf{X}]$  is constructed as follow: For any  $V_i, V_j \in \mathbf{X}$ , (1) add a directed edge  $V_i \rightarrow V_j$  in  $G[\mathbf{X}]$  if there exists a directed path from  $V_i$  to  $V_j$  in  $G$  such that every vertex on the path is not in  $\mathbf{X}$ ; (2) add a bidirected edge  $V_i \leftrightarrow V_j$  in  $G[\mathbf{X}]$  if there exists a divergent path between  $V_i$  and  $V_j$  in  $G$  such that every vertex on the path is not in  $\mathbf{X}$  (Tian & Pearl, 2003).

	<i>Causality</i>	<i>Inaccessibility</i>	<i>Axioms</i>
<b>Conditional</b>	<b>X</b>	×	<b>X</b>
<b>Marginal</b>	×	<b>X</b>	×
<b>Causal</b>	×	<b>X</b>	×
<b>ICC</b>	×	×	<b>X</b>
<b><i>do</i>-Shapley</b>	×	×	×

Table 1: Summary of comparisons of the **conditional**, **marginal**, **causal** Shapley values, and the **ICC** with our method (***do*-Shapley**) w.r.t. consideration of *causality*, capability in handling outcomes induced by an *inaccessible model* (e.g., Example 1), and having justification from *axioms*.

(Lundberg & Lee, 2017; Frye et al., 2020). The conditional Shapley measures contributions based on associations rather than causation. In general, the conditional Shapley doesn’t match with the *do*-Shapley; The causal irrelevance property doesn’t hold in the conditional Shapley (see Example C.1).

**Marginal Shapley.** The *marginal Shapley* is another widely used contribution measure in the XAI in which the target variable is a model prediction  $Y = f(\mathbf{V})$ , where  $f$  is a deterministic (refer Remark 1) and *accessible* prediction model. The marginal Shapley is a specialization of the Shapley value with  $\nu(S) = \mathbb{E}[f(\mathbf{v}_S, \mathbf{V}_{\bar{S}})]$  (Janzing et al., 2020b). The marginal Shapley is known to satisfy certain desiderata in attributing the feature importance (Sundararajan & Najmi, 2020). With access to the model  $f$ , and a particular graphical assumption that features are not causally affecting each other, the marginal Shapley matches with the *do*-Shapley (Janzing et al., 2019, Eq. (14)). In general settings where features are causally related as in Example 1, the marginal Shapley doesn’t match with the *do*-Shapley.

**Causal Shapley.** The *causal Shapley* (Heskes et al., 2020) is most closely related to the *do*-Shapley. Specifically, (Heskes et al., 2020) proposed the same equation as the *do*-Shapley proposed here for measuring the contributions when the outputs are generated by the *accessible* models, and the graph is unknown (only a partial topological ordering of the graph is known). While the *do*-Shapley doesn’t have a restriction that the output is induced by the accessible models and is defined specifically on causal with bidirected edges induced by SCM  $\mathcal{M}$  (See Sec. 2) for which rich theories on causal effect identification and estimation are available.

**Intrinsic Causal Contribution (ICC).** Janzing et al. (2020b) proposed a new method called *Intrinsic Causal Contribution* (ICC) ( $\phi_{V_i}^{\text{icc}}$ ) to measure the causal contribution under the setting where the causal graph is Markovian, and the structural functions are invertible in the sense that the noise values can be reconstructed from the observations. The ICC relies on so-called a structure-based intervention,

which intervenes to features while keeping a causal structure and a joint distribution unaffected, to measure the contribution of  $V_i$  on  $Y$ . By doing so, the ICC can measure the contribution of  $V_i$  on  $Y$  that is not via upstream variables. However, there is no axiomatic characterization of the ICC to the best of our knowledge. It is easy to show that ICC does not satisfy the causal symmetry property (see Example C.2).

**Other Contribution Measures.** Wang et al. (2021a) focused on measuring the relevance of paths in a causal graph to a target node, whereas Singal et al. (2021) provided a recursive approach to capture the flow of importance through the graph. The causal influence defined in Janzing et al. (2013) is based on an operation called ‘deletion of edges’ and measures the relevance of edges with respect to the joint distribution, but not the relevance of edges for a certain target node. Schamberg et al. (2020) describes a generalization of the information-theoretic approach of Janzing et al. (2013) which quantifies relevance of *paths or edges* for a target node, based on operations on edges. Under some particular graphical assumptions, e.g., *flat graphs*, (Singal et al., 2021, Def. 8), the path/edge-based Shapley values (Wang et al., 2021a; Singal et al., 2021) match with the *do*-Shapley. In general, however, the link between these lines of work is yet to be fully established.

#### 4. Identification of the *do*-Shapley

In this section, we investigate the question of evaluating the *do*-Shapley values. To evaluate the *do*-Shapley, expressing  $E[Yjdo(\mathbf{v}_S)]$  as a functional of an observational distribution  $P$  using  $G$  is essential because we are only given non-experimental dataset  $D$  drawn from the observational distribution  $P$ . For each  $S \subseteq [n]$ , complete causal effect identification algorithms  $E[Yjdo(\mathbf{v}_S)]$  are already available (Tian & Pearl, 2003; Huang & Valtorta, 2006; Shpitser & Pearl, 2006). A major practical challenge still remains, however, in using them because determining the identifiability for all subsets  $S \subseteq [n]$  takes *exponential* computation time. In this section, we address this challenge in determining the identifiability by presenting a graphical criterion where the identifiability can be determined in polynomial time, which makes this procedure feasible in practice. Formally,

**Definition 2 (Identifiability & Feasibility).** The *do*-Shapley values  $f\phi_{V_i} \mathcal{G}_{i=1}^n$  w.r.t.  $G$  are said to be *identifiable* if all elements in  $fE[Yjdo(\mathbf{v}_S)] \mathcal{G}_S \subseteq [n]$  are identifiable in the causal graph  $G$ . The identification of the *do*-Shapley values are said to be (computationally-) *feasible* if the identification can be done in  $O(\text{poly}(n))$ .

Since naïvely applying the existing causal effect identification algorithms to determine the identifiability of the *do*-Shapley values is not computationally feasible (requires  $O(2^n)$  computations), we provide a simple sufficient graph-

ical criterion under which determining the *do*-Shapley identifiability is feasible. We start with a definition (refer Sec. 2 for  $C$ -component,  $C$ -factor):

**Definition 3 ( $C$ -partition).** For a set of variables  $\mathbf{X} \subseteq \mathbf{V}$ ,  $f\mathbf{X}_k \mathcal{G}_{k=1}^c$  is said to be the  $C$ -partition if  $\mathbf{X} = \bigcup_{k=1}^c \mathbf{X}_k$  (where  $\mathbf{X}_a \cap \mathbf{X}_b = \emptyset$  for  $a \neq b$ ) where  $\forall k \in [c], \mathbf{X}_k$  is a set s.t. any two pairs  $X_i, X_j \in \mathbf{X}_k$  are in the same  $C$ -component. in  $G$ .

**Theorem 2 (Identifiability & Feasibility of *do*-Shapley).** The *do*-Shapley is identifiable if no variable in  $V_i \in f\mathbf{V} \mathcal{G}$  is connected to its child  $\text{Ch}(V_i)$  by bidirected paths in  $G$ . Suppose  $Y$  is not connected by bidirected paths. In this case, for any  $S \subseteq [n]$ ,

$$E[Yjdo(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] Q[\mathbf{V} \cap \mathbf{V}_S],$$

where  $Q[\mathbf{V} \cap \mathbf{V}_S] := Q[\mathbf{V} \cap \mathbf{V}_S](\mathbf{v})$  is given as

$$Q[\mathbf{V} \cap \mathbf{V}_S] = \frac{P(\mathbf{v})}{Q[C(\mathbf{V}_S)]} \prod_{k=1}^c \sum_{\mathbf{s}_k} Q[C(\mathbf{S}_k)],$$

where  $Q[C(\mathbf{V}_S)] = \prod_{V_a \in C(\mathbf{V}_S)} P(v_a | \text{pre}(v_a))$  is a  $C$ -factor of a  $C$ -component  $\mathbf{V}_S$  ( $C(\mathbf{V}_S)$ );  $f\mathbf{S}_k \mathcal{G}_{k=1}^c$  is a  $C$ -partition of  $\mathbf{V}_S$ ; and  $Q[C(\mathbf{S}_k)] := \prod_{V_a \in C(\mathbf{S}_k)} P(v_a | \text{pre}(v_a))$  is a  $C$ -factor of a  $C$ -component  $C(\mathbf{S}_k)$  for  $\mathbf{S}_k$ .

Fig. 2a is an example graph where a graph satisfies the conditions in Thm. 2. Specifically, for all computations  $\mathbf{V}_S \subseteq \mathbf{V} := fV_1, V_2, V_3 \mathcal{G}$ , the causal effects are identified through Thm. 2 as

$$\begin{aligned} E[Yjdo(\mathbf{v}_S)] &= \begin{cases} \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(v_2 | v_1, v_3) P(\mathbf{v}_S), & \text{if } S \subseteq \{1, 3\}, \\ \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(\mathbf{v}_{\bar{S}}), & \text{if } S \subseteq \{1, 2, 3\}, \\ \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(\mathbf{v}_{\bar{S}} | \mathbf{v}_S), & \text{if } S \subseteq \{1, 2, 3\}, \\ E[Yj\mathbf{v}] & \text{if } S = \{1, 2, 3\}. \end{cases} \end{aligned} \quad (3)$$

Thm. 2 suggests the feasibility of the *do*-Shapley values since the proposed graphical criteria (checking whether  $V_i$  and  $\text{Ch}(V_i)$  are connected by bidirected paths) can be done in  $O(n^3)$  by applying the breadth-first-search for each variable  $V_i \in \mathbf{V}$ .

To demonstrate the applicability of Theorem 2, we provide two special cases which are commonly considered in the literature:

1. **Markovian case:** All latent variables in the SCM is independent; i.e.,  $G$  is given as a DAG (Janzing et al., 2013; 2019; Heskes et al., 2020; Basu, 2020; Wang et al., 2021b; Singal et al., 2021).

2. **Direct-cause case:** No pair of variables  $(V_i, V_j) \subseteq \mathbf{V}$  ( $i \neq j$ ) is connected by a directed path, no  $V_i$  are connected to  $Y$  via bidirected edges, and no directed edge from  $Y$  to  $V_i$  exists (i.e, only  $V_i \rightarrow Y$  is allowed) (Janzing et al., 2020a;b).

For each of these cases, the identification result in Theorem 2 can be simplified as follows.

**Corollary 1 (Identification – Markovian).** *In the Markovian case,  $E[Yjdo(\mathbf{v}_S)]$  is given as*

$$E[Yjdo(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}_S, \mathbf{v}_{\bar{S}}] \prod_{i \notin S} P(v_i | pre(v_i)).$$

**Corollary 2 (Identification – Direct-cause).** *In the Direct-cause case,  $E[Yjdo(\mathbf{v}_S)]$  is given as*

$$E[Yjdo(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}_S, \mathbf{v}_{\bar{S}}] P(\mathbf{v}_{\bar{S}}).$$

Figs. (2b,2c) provide example graphs for Markovian and Direct-cause cases. For Fig. 2b, Coro. 1 gives

$$E[Yjdo(\mathbf{v}_S)] = \begin{cases} \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(\mathbf{v}_{\bar{S}} | \mathbf{v}_S), & \text{if } S \subseteq \{1, 3, \bar{1}, 3g\}, \\ \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(\mathbf{v}_{\bar{S}}), & \text{if } S \subseteq \{1, 2, \bar{2}, 3g, \bar{1}, 2gg\}, \\ E[Yj\mathbf{v}] & \text{if } S = \{1, 2, 3g\}. \end{cases} \quad (4)$$

For Fig. 2c, Coro. 2 gives

$$E[Yjdo(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Yj\mathbf{v}] P(\mathbf{v}_{\bar{S}}) \quad (5)$$

for all  $\mathbf{v}_S = \{v_1, v_2, v_3g\}$ .

## 5. Estimation of the *do*-Shapley

Estimating the *do*-Shapley values in Eq. (2) is computationally and statistically challenging because (1) Iterating over all  $S \subseteq [n]$  takes time exponential in  $n$ , and (2) Estimating  $E[Yjdo(\mathbf{v}_S)]$  might be vulnerable to bias due to finiteness of the sample dataset. In this section, we design computationally efficient and statistically robust estimators for the *do*-Shapley values to overcome these challenges, using three different techniques. For ease of presentation, we focus only on the Markovian & Direct-cause cases discussed in Sec. 4.

We first introduce estimators leveraging the idea of the inverse probability weighting (IPW) (Rosenbaum & Rubin, 1983). Our construction of the IPW estimator is based on the following result.

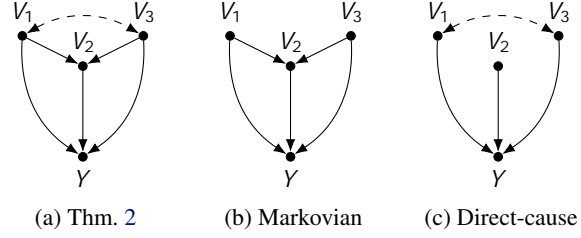


Figure 2: Example graphs for Thm. 2 and two special cases: Markovian and Direct-cause.

**Lemma 1 (Representation using IPW).** *Let  $S = \{m_1, \dots, m_s\} \subseteq [n]$  denote an index set for  $\mathbf{V}_S$ . Let*

$$\omega_k^S := \prod_{r=1}^k \mathbb{1}_{V_{m_r}}(V_{m_r}) / h_r^S, \text{ for } k = s, \dots, 1;$$

$$\omega^S := \mathbb{1}_{\mathbf{v}_S}(\mathbf{V}_S) / h^S,$$

where  $h_r^S := P(V_{m_r} | pre(V_{m_r}))$  and  $h^S := P(\mathbf{V}_S | \mathbf{V}_{\bar{S}})$ . Then,  $E[Yjdo(\mathbf{v}_S)] = E[Y\omega]$  where  $\omega = \omega_k^S$  for the Markovian case, and  $\omega = \omega^S$  for the Direct-cause case.

Using Lemma 1, we construct the IPW estimators.

**Definition 4 (IPW for  $E[Yjdo(\mathbf{v}_S)]$ ).** The IPW estimator  $T^{\text{IPW}}(S)$  for  $E[Yjdo(\mathbf{v}_S)]$  is constructed as:

1. Split  $D$  randomly into two halves:  $D_0$  and  $D_1$ ;
2. Let  $\hat{\omega}_{s,p}^S, \hat{\omega}_p^S$  denote estimators for  $\omega_s^S, \omega^S$  from  $D_p \subseteq D_0, D_1$ , respectively.
3. For each  $p \in \{0, 1\}$ , set

$$T_p^{\text{IPW}}(S) := \begin{cases} E_{D_1-p} [Y \hat{\omega}_{s,p}^S] & \text{(Markovian)} \\ E_{D_1-p} [Y \hat{\omega}_p^S] & \text{(Direct-cause)} \end{cases}$$

4.  $T^{\text{IPW}}(S) := \frac{1}{2} (T_0^{\text{IPW}}(S) + T_1^{\text{IPW}}(S))$ .

The data-splitting (also known as sample-splitting) technique (Klaassen, 1987; Robins & Ritov, 1997; Robins et al., 2008; Zheng & van der Laan, 2011; Chernozhukov et al., 2018) will be employed in constructing all *do*-Shapley estimators discussed in this section. Without data-splitting, some restriction on the complexity of the estimator function class must be imposed to guarantee statistical consistency.

We introduce estimators leveraging the idea of outcome regression (REG) (Rubin, 1979). Our REG estimator is based on the following result.

**Lemma 2 (Representation using REG).** *Let  $S = \{m_1, \dots, m_s\} \subseteq [n]$  denote an index set for  $\mathbf{V}_S$ . Let*

$\theta_{S,1}^S := Y$ . For  $k = s, s-1, \dots, 1$ ,

$$\begin{aligned}\theta_{k,2}^S &:= E[\theta_{k,1}^S | \mathcal{V}_{m_k}, \text{pre}(V_{m_k})] \\ \theta_{k,1,1}^S &:= E[\theta_{k,1}^S | \mathcal{V}_{m_k}, \text{pre}(V_{m_k})], \\ \theta_a^S &:= E[Y | \mathcal{V}_S, \mathbf{V}_{\bar{S}}], \\ \theta_b^S &:= E[Y | \mathcal{V}_S, \mathbf{V}_{\bar{S}}].\end{aligned}$$

Then,  $E[Y | \mathcal{J}do(\mathbf{v}_S)] = E[\theta]$  where  $\theta = \theta_{0,1}^S$  for the Markovian case, and  $\theta = \theta_a^S$  for the Direct-cause case.

We construct the REG estimator based on Lemma 2.

**Definition 5 (REG for  $E[Y | \mathcal{J}do(\mathbf{v}_S)]$ ).** The REG estimator  $T^{\text{reg}}(S)$  for  $E[Y | \mathcal{J}do(\mathbf{v}_S)]$  is constructed as:

1. Split  $D$  randomly into two halves:  $D_0$  and  $D_1$ .
2. Let  $\hat{\theta}_{k,2,p}^S, \hat{\theta}_{k,1,1,p}^S, \hat{\theta}_a^S$  denote an estimator for  $\theta_{k,2}^S, \theta_{k,1,1}^S, \theta_a^S$  from  $D_p \stackrel{i.i.d.}{\sim} D_0, D_1$ , respectively.
3. For each  $p \stackrel{i.i.d.}{\sim} \{0, 1\}$ ,

$$T_p^{\text{reg}}(S) := \begin{cases} E_{D_1, p} \left[ \hat{\theta}_{0,1,p}^S \right] & \text{(Markovian)} \\ E_{D_1, p} \left[ \hat{\theta}_a^S \right] & \text{(Direct-cause).} \end{cases}$$

4.  $T^{\text{reg}}(S) := \bar{f} T_0^{\text{reg}}(S) + T_1^{\text{reg}}(S)g/2$ .

For IPW and REG estimators to be consistent, one needs to estimate each individual functional (called *nuisances*) including  $E[Y | \mathcal{V}_S, \mathbf{v}_{\bar{S}}]$  or  $P(v_i | \text{pre}(v_i))$  consistently. A desirable robust estimator is one that converges to the ground-truth at a fast rate even when estimates for nuisances are misspecified (i.e., wrongly specified) or converging relatively slowly. *Double/Debiased Machine Learning (DML)* (Chernozhukov et al., 2017) is a recently introduced technique to construct such estimators.

**Lemma 3 (Representation using DML).** Let

$$\eta^S := \begin{cases} \bar{f} \theta_{0,1}^S g + \bar{f} \theta_{k,1}^S, \theta_{k,2}^S g_{k=1} + \bar{f} h_r^S g_{r=1} & \text{(Markovian)} \\ \bar{f} \theta_a^S, \theta_b^S, h^S g & \text{(Direct-cause),} \end{cases}$$

defined in Defs. (4, 5) above, and

$$V_S(\mathbf{V}^\theta; \eta^S) := \begin{cases} \theta_{0,1}^S + \sum_{k=1}^s \omega_k^S (\theta_{k,1}^S - \theta_{k,2}^S) & \text{(Markovian)} \\ \theta_a^S + \omega^S (Y - \theta_b^S) & \text{(Direct-cause),} \end{cases}$$

where  $\omega_k^S := \prod_{r=1}^k \mathbb{1}_{\mathcal{V}_{m_r}}(V_{m_r})/h_r^S$  and  $\omega^S := \mathbb{1}_{\mathcal{V}_S}(\mathbf{V}_S)/h^S$ . Then,  $E[Y | \mathcal{J}do(\mathbf{v}_S)] = E[V_S(\mathbf{V}^\theta; \eta^S)]$ .

We construct the DML estimators based on Lemma 3:

**Definition 6 (DML for  $E[Y | \mathcal{J}do(\mathbf{v}_S)]$ ).** The DML estimator  $T^{\text{dml}}(S)$  is constructed as:

1. Split  $D$  randomly into two halves:  $D_0$  and  $D_1$ ;

**Algorithm 1** *do*-Shapley( $M, T^{\text{est}}(\cdot)$ )

- 1: **Input:**  $M$ , Estimators  $T^{\text{est}}(\cdot)$  in Defs. (4,5,6).
- 2: **Output:** Estimates  $\hat{f}_{\phi_{V_i}}^S, g_{I=1}^S$ .
- 3: Initialize  $\hat{\phi}_{V_i} = 0$  for all  $V_i \in \mathbf{V}$ .
- 4: **for**  $j = 1$  **to**  $M$  **do**
- 5:     Generate the random permutation  $\pi$  over  $[n]$ .
- 6:     **for**  $i = 1$  **to**  $n$  **do**
- 7:          $\hat{\phi}_{V_i} = \hat{\phi}_{V_i} + T^{\text{est}}(\hat{f}_i, \text{pre}(\pi_j(i))) - T^{\text{est}}(\text{pre}(\pi_j(i)))$
- 8:     **end for**
- 9: **end for**
- 10: **return**  $\hat{f}_{\phi_{V_i}}^S / M g_{I=1}^S$

2. Construct  $\hat{\eta}_p^S$ , estimates of  $\eta^S$  from  $D_p, p \stackrel{i.i.d.}{\sim} \{0, 1\}$ .

3.  $T_p^{\text{dml}}(S) := E_{D_1, p} [V_S(\mathbf{V}; \hat{\eta}_p^S)]$  for  $p \stackrel{i.i.d.}{\sim} \{0, 1\}$ .

4.  $T^{\text{dml}}(S) := \bar{f} T_0^{\text{dml}}(S) + T_1^{\text{dml}}(S)g/2$ .

Based on estimators in Defs. (4,5,6), we now propose a computationally efficient estimator for the *do*-Shapley values based on random permutations:

**Definition 7 (*do*-Shapley estimators – Two cases).** Let  $T^{\text{est}}(S) \stackrel{i.i.d.}{\sim} \bar{f} T^{\text{ipw}}(S), T^{\text{reg}}(S), T^{\text{dml}}(S)g$  denote an estimator for  $E[Y | \mathcal{J}do(\mathbf{v}_S)]$  defined in Defs. (4,5,6), respectively. The *do*-Shapley estimator is given as

$$\phi_{V_i}^{\text{est}} := \frac{1}{M} \sum_{j=1}^M (T^{\text{est}}(\hat{f}_i, \text{pre}_{\pi_j}(i)) - T^{\text{est}}(\text{pre}_{\pi_j}(i))),$$

where  $M$  is the number of randomly generated permutation of  $[n]$ ,  $\pi_j$  denotes the  $j$ th permutation, and  $\text{pre}_{\pi_j}(i)$  is the set of elements that precedes  $i$  in  $\pi_j$ .

A systematic procedure for constructing *do*-Shapley estimators is provided in Algorithm 1. The following theorem summarizes the error analyses of all the three *do*-Shapley estimators.

**Theorem 3 (Bias Analysis).** Let  $\bar{f}_{\pi_j} g_{j=1}^M$  denote  $M$  randomly generated permutations of  $[n]$ . For the fixed index  $i$ , let  $S_{j,0} := \text{pre}_{\pi_j}(i)$  and  $S_{j,1} := \hat{f}_i g[S_{j,0}]$ . Let  $\bar{f}_{\eta^{S_{j,0}}}, \bar{\eta}^{S_{j,1}}, g_{j=1}^M$  denote  $L_2$ -consistent estimates for all nuisances  $\bar{f}_{\eta^{S_{j,0}}}, \eta^{S_{j,1}}, g_{j=1}^M$  defined in Def. 6. Let  $R_{M:N} := O_P(M^{-1/2} + N^{-1/2})$ . Let  $e(\hat{g}) := k\hat{g} - g$  denote an error for a nuisance estimates for any  $\hat{g} \stackrel{i.i.d.}{\sim} \hat{\eta}$  and  $g \in \eta$ . For the *do*-Shapley estimators defined in Def. 7, suppose the estimators  $T^{\text{est}}(S)$  are bounded. Let  $\epsilon_{V_i}^{\text{est}} := \phi_{V_i}^{\text{est}} - \phi_{V_i}$  (where  $\text{est} \in \{\text{ipw}, \text{reg}, \text{dml}\}$ ).

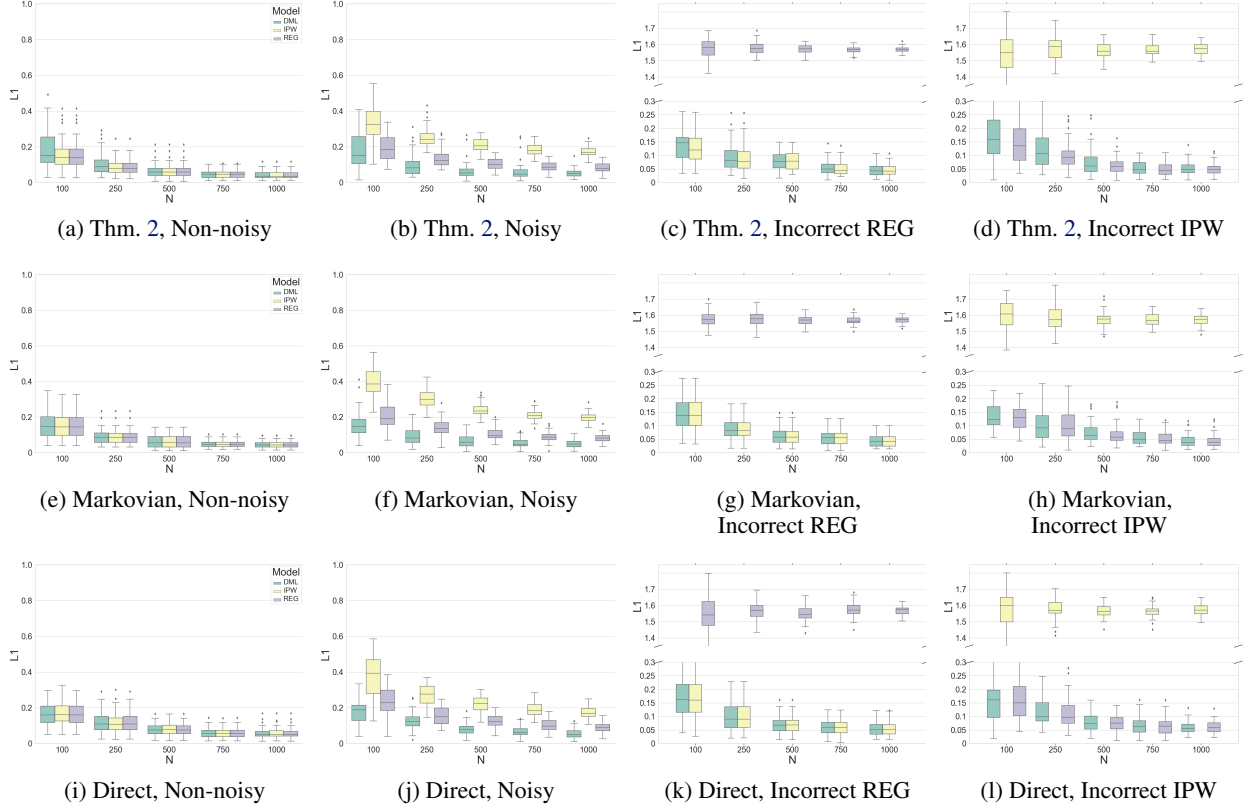


Figure 3: The L1-error plots. Plots are rendered in high resolution and can be zoomed in.

Under the Markovian case,

$$\begin{aligned}\epsilon_{V_i}^{\text{ipw}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} e(\hat{\omega}_{S_j^{j,p}}) g, \\ \epsilon_{V_i}^{\text{reg}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} e(\hat{\theta}_{0;1}^{S_j^{j,p}}) g, \\ \epsilon_{V_i}^{\text{dml}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} \sum_{k=1}^{S_j} e(\hat{h}_k^{S_j^{j,p}}) e(\hat{\theta}_{k;2}^{S_j^{j,p}}) g.\end{aligned}$$

Under the Direct-cause case,

$$\begin{aligned}\epsilon_{V_i}^{\text{ipw}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} e(\hat{\omega}^{S_j^{j,p}}) g, \\ \epsilon_{V_i}^{\text{reg}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} e(\hat{\theta}_2^{S_j^{j,p}}) g, \\ \epsilon_{V_i}^{\text{dml}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{\tau \geq 0; 1 \leq j \leq M} e(\hat{h}^{S_j^{j,p}}) e(\hat{\theta}_b^{S_j^{j,p}}) g.\end{aligned}$$

**Remark 4 (Properties of the Proposed Estimators).** Error analyses in Thm. 3 exhibit consistency of IPW, REG,

DML estimators. Specifically, if nuisances are consistently estimated,  $\epsilon_{V_i}^{\text{ipw}} = \epsilon_{V_i}^{\text{reg}} = \epsilon_{V_i}^{\text{dml}} = O_P(1)$ , indicating that the estimators converge to the true quantity. Furthermore, the result presents the statistical robustness property of the DML. In particular, the DML estimates  $\hat{\phi}_{V_i}^{\text{dml}}$  converges to the true value if either  $e(\hat{h}_k^{S_j^{j,p}})$  or  $e(\hat{\theta}_{k;2}^{S_j^{j,p}})$  under Markovian, and either  $e(\hat{h}^{S_j^{j,p}})$  or  $e(\hat{\theta}_b^{S_j^{j,p}})$  under Direct-cause are accurate (doubly robustness). Also,  $\hat{\phi}_{V_i}^{\text{dml}}$  converges at the root- $N$  rate if all nuisances  $\hat{h}_k^{S_j^{j,p}}, \hat{\theta}_{k;2}^{S_j^{j,p}}$  under Markovian, and all nuisances  $\hat{h}^{S_j^{j,p}}, \hat{\theta}_b^{S_j^{j,p}}$  under Direct-cause converge at least at  $N^{-1/4}$  rate (debiasedness).

## 6. Experiments

In this section, we empirically compare the performance of the proposed *do*-Shapley estimators from the previous section. Details of the experiments and a different simulation example are provided in Appendices E and F.

**Experimental Setup.** We use synthetic datasets based on Figs. (2a, 2b, 2c) where each figures matches with Thm. 2, Markovian, and Direct-cause cases. We note that causal effects are identified as in Eqs. (3, 4, 5), respectively. Even if no known estimators for Thm. 2 exist generally, we note



that Eq. (3) is in an amenable form for which results in Sec. 5 are applicable. Throughout the simulation, we denote  $\bar{f}\phi_{V_i}g_{i=1}^n$  as the ground-truth *do*-Shapley values.

**Comparison Between Estimators.** We compare the three estimators (IPW, REG, DML), denoted by  $\bar{f}\phi_{V_i}^{\text{ipw}}, \bar{f}\phi_{V_i}^{\text{reg}}, \bar{f}\phi_{V_i}^{\text{dml}}g$  respectively, for scenarios depicted in graphs in Figs. (2a, 2b, 2c). Nuisances are estimated using gradient boosting model (Friedman, 2001).

Let  $\phi_{V_i,k}^{\text{est}} \geq \bar{f}\phi_{V_i,k}^{\text{dml}}, \phi_{V_i,k}^{\text{ipw}}, \phi_{V_i,k}^{\text{reg}}g$  denote an estimated importance of the  $i$ th feature of  $j$ th samples (i.e.,  $V_{i:k} \geq \mathbf{V}^{(k)} \geq D$ ). We assess the quality of the estimator by computing the  $L_1$  error as  $L_1(\text{est}, k) := (1/n) \sum_{i=1}^n |\phi_{V_i;k}^{\text{est}} - \phi_{V_i;k}|$  (where  $n$  is the number of features). We ran the simulation for 50 randomly generated sets of samples; i.e.,  $k \geq \{1, 2, \dots, 50\}$ , and with sample size  $N := \{100, 250, 500, 750, 1000\}$  to observe convergence behaviors of estimators. We fix  $M = 20$ . We refer the box-plot for  $L_1(\text{est}, k)$  as the ‘L1-error plot’.

For all  $\{ \text{Thm. 2, Markovian, Direct-cause} \}$  cases, we compare the performances of the three *do*-Shapley estimators for (1) ‘Non-noisy’ where no noises are introduced in the model; (2) ‘Noisy’ where a ‘converging noise’  $\epsilon$ , decaying at a  $N^{-\alpha}$  rate (i.e.,  $\epsilon \sim \text{Normal}(N^{-\alpha}, N^{-2\alpha})$ ) for  $\alpha = 1/4$ , is added to the estimated nuisance to control the convergence rate, following the technique in (Kennedy, 2020); (3) ‘Incorrect REG’ where the model for the REG estimator in Def. 5 is wrongly specified; and (4) ‘Incorrect IPW’ the model for the IPW estimator in Def. 4 is wrongly specified.

**Experimental Results.** The L1-error plots for all cases are presented in Fig. 3. For the non-noisy setting, performances of all three models {DML, REG, IPW} are similar. In the noisy setting where the estimated nuisances are controlled to converge at  $N^{-1/4}$  rate, the DML estimators outperform the other two estimators by achieving a fast convergence with the smallest variance. This result corroborates the robustness property of the DML (Remark 4). Also, the DML estimator exhibits the doubly robustness property; the estimator converges in both the ‘Incorrect IPW’ and ‘Incorrect REG’ settings where each corresponding nuisance is wrongly specified.

**Contrasting with Conditional Shapley.** We contrast the *do*-Shapley and conditional Shapley in the non-noisy setting. We compare the importance ranking measured by the true *do*-Shapley with the ranks from the *do*-DML and conditional Shapley through the Spearman’s rank correlation. The correlation is close to 1 if two ranks are similar and -1 if the ranks are opposite. The true data generating function is  $Y = 3V_1 + 0.4V_2 + V_3 + U_Y$  and the true-*do*-Shapley identifies  $V_1$  having the largest coefficient as the most important. As shown in Table 2, the *do*-DML-Shapley ranks the feature importance closer to the true rank. As noted,

*do*-DML-Shapley identifies  $V_1$  as the most important.

	Thm. 2	Markovian	Direct
DML	1.0	0.8	0.93
Conditional	-0.28	-0.74	0.52

Table 2: Comparison of the rank correlation.

## 7. Conclusion

We proposed the *do*-Shapley for measuring causal contribution and provided a theoretical justification through the axiomatic characterization (Thm. 1). Next, we developed conditions under which *do*-Shapley values can be inferred from non-experimental data in polynomial time (Thm. 2). We then proposed three *do*-Shapley estimators (IPW, REG, DML) that are consistent. We showed that the DML estimator has additional robustness property called doubly robustness and debiasedness (Thm. 3). We expect the proposed contribution measure will help empirical scientists to answer “*what are the contributions of each cause to the effect?*”

## Acknowledgements

We thank the reviewers for their feedback and help in improving this manuscript. This work was done in part while Jin Tian was visiting the Simons Institute for the Theory of Computing. Elias Bareinboim and Yonghan Jung were partially supported by funding from the NSF, ONR, DoE, Amazon, JP Morgan, and The Alfred P. Sloan Foundation.

## References

- Bareinboim, E. and Pearl, J. Causal inference by surrogate experiments: z-identifiability. In *In Proceedings of the 28th Conference on Uncertainty in Artificial Intelligence*, pp. 113–120. AUAI Press, 2012.
- Bareinboim, E. and Pearl, J. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.
- Bareinboim, E., Correa, J. D., Ibeling, D., and Icard, T. On pearl’s hierarchy and the foundations of causal inference. In *Probabilistic and Causal Inference: The Works of Judea Pearl*, pp. 507–556. 2020.
- Basu, D. On shapley credit allocation for interpretability. *arXiv preprint arXiv:2012.05506*, 2020.
- Bhattacharya, R., Nabi, R., and Shpitser, I. Semiparametric inference for causal effects in graphical models with hidden variables. *arXiv preprint arXiv:2003.12659*, 2020.

- Bhattacharyya, A., Gayen, S., Kandasamy, S., Maran, A., and Variyam, V. N. Learning and sampling of atomic interventions from observations. In *International Conference on Machine Learning*, pp. 842–853. PMLR, 2020.
- Bhattacharyya, A., Gayen, S., Kandasamy, S., Raval, V., and Vinodchandran, N. Efficient inference of interventional distributions. *arXiv preprint arXiv:2107.11712*, 2021.
- Charnes, A., Golany, B., Keane, M., and Rousseau, J. Extremal principle solutions of games in characteristic function form: Core, chebychev and shapley value generalizations. 1988.
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., and Newey, W. Double/debiased/neyman machine learning of treatment effects. *American Economic Review*, 107(5):261–65, 2017.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal*, 21(1), 2018.
- Covert, I., Lundberg, S., and Lee, S.-I. Explaining by removing: A unified framework for model explanation. *Journal of Machine Learning Research*, 22(209):1–90, 2021.
- Eva, B. and Stern, R. Causal explanatory power. *The British Journal for the Philosophy of Science*, 70(4):1029–1050, 2019.
- Friedman, E. and Moulin, H. Three methods to share joint costs or surplus. *Journal of economic Theory*, 87(2): 275–312, 1999.
- Friedman, J. H. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pp. 1189–1232, 2001.
- Frye, C., Rowat, C., and Feige, I. Asymmetric shapley values: incorporating causal knowledge into model-agnostic explainability. *Advances in Neural Information Processing Systems*, 33, 2020.
- Galles, D. and Pearl, J. Axioms of causal relevance. *Artificial Intelligence*, 97(1-2):9–43, 1997.
- Heskes, T., Sijben, E., Bucur, I. G., and Claassen, T. Causal shapley values: Exploiting causal knowledge to explain individual predictions of complex models. *Advances in Neural Information Processing Systems*, 33, 2020.
- Huang, Y. and Valtorta, M. Pearl’s calculus of intervention is complete. In *Proceedings of the 22nd Conference on Uncertainty in Artificial Intelligence*, pp. 217–224. AUAI Press, 2006.
- Jaber, A., Zhang, J., and Bareinboim, E. Causal identification under markov equivalence. In *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence*, 2018.
- Janzing, D., Balduzzi, D., Grosse-Wentrup, M., and Schölkopf, B. Quantifying causal influences. *The Annals of Statistics*, 41(5):2324–2358, 2013.
- Janzing, D., Budhathoki, K., Minorics, L., and Blöbaum, P. Causal structure based root cause analysis of outliers. *arXiv preprint arXiv:1912.02724*, 2019.
- Janzing, D., Blöbaum, P., Minorics, L., and Faller, P. Quantifying causal contributions via structure preserving interventions. *arXiv preprint arXiv:2007.00714*, 2020a.
- Janzing, D., Minorics, L., and Blöbaum, P. Feature relevance quantification in explainable ai: A causal problem. In *International Conference on Artificial Intelligence and Statistics*, pp. 2907–2916. PMLR, 2020b.
- Jung, Y., Tian, J., and Bareinboim, E. Learning causal effects via weighted empirical risk minimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Jung, Y., Tian, J., and Bareinboim, E. Estimating identifiable causal effects on markov equivalence class through double machine learning. In *Proceedings of the 38th International Conference on Machine Learning*, 2021a.
- Jung, Y., Tian, J., and Bareinboim, E. Estimating identifiable causal effects through double machine learning. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021b.
- Kennedy, E. H. Optimal doubly robust estimation of heterogeneous causal effects. *arXiv preprint arXiv:2004.14497*, 2020.
- Kennedy, E. H., Balakrishnan, S., G’Sell, M., et al. Sharp instruments for classifying compliers and generalizing causal effects. *Annals of Statistics*, 48(4):2008–2030, 2020.
- Klaassen, C. A. Consistent estimation of the influence function of locally asymptotically linear estimators. *The Annals of Statistics*, pp. 1548–1562, 1987.
- Koster, J. T. et al. Marginalizing and conditioning in graphical models. *Bernoulli*, 8(6):817–840, 2002.
- Lattimore, T. and Szepesvári, C. *Bandit algorithms*. Cambridge University Press, 2020.

- Lee, S. and Bareinboim, E. Causal effect identifiability under partial-observability. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Lee, S., Correa, J. D., and Bareinboim, E. General identifiability with arbitrary surrogate experiments. In *Proceedings of the 35th Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2019.
- Lee, S., Correa, J., and Bareinboim, E. Generalized transportability: Synthesis of experiments from heterogeneous domains. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, 2020.
- Lundberg, S. Be careful when interpreting predictive models in search of causal insights. 2021. URL <https://bit.ly/3gcZmgl>.
- Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*, pp. 4768–4777, 2017.
- Molnar, C. *Interpretable machine learning*. Lulu. com, 2020.
- Pearl, J. Causal diagrams for empirical research. *Biometrika*, 82(4):669–710, 1995.
- Pearl, J. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, NY, 2000. 2nd edition, 2009.
- Robins, J. A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Mathematical modelling*, 7(9-12):1393–1512, 1986.
- Robins, J., Li, L., Tchetgen, E., van der Vaart, A., et al. Higher order influence functions and minimax estimation of nonlinear functionals. In *Probability and statistics: essays in honor of David A. Freedman*, pp. 335–421. Institute of Mathematical Statistics, 2008.
- Robins, J. M. and Ritov, Y. Toward a curse of dimensionality appropriate (coda) asymptotic theory for semi-parametric models. *Statistics in medicine*, 16(3):285–319, 1997.
- Rosenbaum, P. R. and Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- Rubin, D. B. Using multivariate matched sampling and regression adjustment to control bias in observational studies. *Journal of the American Statistical Association*, 74(366a):318–328, 1979.
- Schamberg, G., Chapman, W., Xie, S.-P., and Coleman, T. P. Direct and indirect effects—an information theoretic perspective. *Entropy*, 22(8):854, 2020.
- Schwab, P. and Karlen, W. Cxplain: Causal explanations for model interpretation under uncertainty. *Advances in Neural Information Processing Systems* 32, pp. 10220–10230, 2019.
- Shapley, L. Annals of mathematics study no. 28. 1953.
- Shapley, L. S. and Shubik, M. A method for evaluating the distribution of power in a committee system. *American political science review*, 48(3):787–792, 1954.
- Shpitser, I. and Pearl, J. Identification of joint interventional distributions in recursive semi-markovian causal models. In *Proceedings of the 21st AAAI Conference on Artificial Intelligence*, pp. 1219, 2006.
- Singal, R., Michailidis, G., and Ng, H. Flow-based attribution in graphical models: A recursive shapley approach. Available at SSRN 3845526, 2021.
- Štrumbelj, E. and Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems*, 41(3):647–665, 2014.
- Sundararajan, M. and Najmi, A. The many shapley values for model explanation. In *International Conference on Machine Learning*, pp. 9269–9278. PMLR, 2020.
- Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.
- Tian, J. and Pearl, J. On the identification of causal effects. Technical Report R-290-L, 2003.
- Wang, J., Wiens, J., and Lundberg, S. Shapley flow: A graph-based approach to interpreting model predictions. In Banerjee, A. and Fukumizu, K. (eds.), *The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 721–729. PMLR, 2021a.
- Wang, L., Zhang, Y., Richardson, T. S., and Robins, J. M. Estimation of local treatment effects under the binary instrumental variable model. *Biometrika*, 2021b.
- Winter, E. The shapley value. *Handbook of game theory with economic applications*, 3:2025–2054, 2002.
- Xia, K., Lee, K.-Z., Bengio, Y., and Bareinboim, E. The causal-neural connection: Expressiveness, learnability, and inference. *Advances in Neural Information Processing Systems*, 34:10823–10836, 2021.
- Young, H. P. Monotonic solutions of cooperative games. *International Journal of Game Theory*, 14(2):65–72, 1985.
- Zheng, W. and van der Laan, M. J. Cross-validated targeted minimum-loss-based estimation. In *Targeted Learning*, pp. 459–474. Springer, 2011.

---

## Appendix – On Measuring Causal Contributions via *do*-interventions

---

### A. Fundamentals of the Shapley Value

The Shapley value (Shapley, 1953) in Eq. (1) seeks to allocate the contribution of each  $i \in [n]$  on some function value  $f([n])$  given a coalition function  $\nu(S)$  that measures the value of coalition of values of players  $i \in S$  (where  $\nu([n]) = f([n])$ ). The Shapley value uniquely satisfying the following desiderata:

**Theorem A.1 (Axiomatization of the Shapley Value (Shapley, 1953; Shapley & Shubik, 1954; Young, 1985)).** *For any subset  $S$  of the players indexed  $[n] = \{1, 2, \dots, n\}$  and the value function of  $S$ , denoted  $\nu(S)$ , the Shapley value of the player  $i$ , denoted  $\phi_i = \phi_i(\nu)$ , equals*

$$\phi_i(\nu) := \frac{1}{n} \sum_{S \subseteq [n], i \in S} \binom{n-1}{|S|-1} \nu(S \cup \{i\}) - \nu(S), \quad (\text{A.1})$$

is the unique attribution methods satisfying the following axioms (properties):

1. **Efficiency:**  $\sum_{i \in [n]} \phi_i = \nu([n])$ ;
2. **Dummy:** For some  $i \in [n]$ , if  $\nu(S \cup \{i\}) = \nu(S)$  for all  $S \subseteq [n] \setminus \{i\}$ , then  $\phi_i = 0$ ;
3. **Symmetry:** For some distinct  $(i, j) \in [n]$ , if  $\nu(S \cup \{i\}) = \nu(S \cup \{j\})$  for all  $S \subseteq [n] \setminus \{i, j\}$ , then  $\phi_i = \phi_j$ ;
4. **Linearity:** For all  $i \in [n]$ , for any two coalition functions  $\nu_1$  and  $\nu_2$ ,  $\phi_i(\nu_1 + \nu_2) = \phi_i(\nu_1) + \phi_i(\nu_2)$ .

### B. Details on Remark 3

Given a semi-Markovian causal graph  $G$ , a realized vector  $(\mathbf{v}, y)$  corresponding to a set of variables  $(\mathbf{V}, Y)$  and its subset  $\mathbf{x} \subseteq \mathbf{v}$  corresponding to a set of variables  $\mathbf{X} \subseteq \mathbf{V}$ , a procedure for assigning contributions only to  $x_i \in \mathbf{x}$  is the following:

1. Construct a graph  $G[\mathbf{X}]$  composed of nodes in  $\mathbf{X}$  and edges added as follows (Tian & Pearl, 2003).
  - (a) add a directed edge  $V_i \rightarrow V_j$  in  $G[\mathbf{C}]$  if there exists a directed path from  $V_i$  to  $V_j$  in  $G$  such that every vertex on the path is not in  $\mathbf{C}$ ;
  - (b) add a bidirected edge  $V_i \leftrightarrow V_j$  in  $G[\mathbf{C}]$  if there exists a divergent path between  $V_i$  and  $V_j$  in  $G$  such that every vertex on the path is not in  $\mathbf{C}$ .
2. Construct the *do*-Shapley w.r.t.  $f(y, \mathbf{x}g)$  on  $G[\mathbf{X}]$ . Specifically, for all  $x_i \in \mathbf{x}$

$$\phi_{x_i} := \sum_{\mathbf{x}_S \subseteq \mathbf{x} \setminus x_i} \omega^{\mathbf{x}}(S) \{E[Y | do(\mathbf{x}_S; i)] - E[Y | do(\mathbf{x}_S)]\} g, \quad (\text{B.1})$$

where  $\omega^{\mathbf{x}}(S) := (1/|\mathbf{x}|) \binom{|\mathbf{x}|-1}{|S|-1}$ .

Then,  $f\phi_{x_i} g_{x_i \in \mathbf{x}}$  is a unique causal contribution measure:

**Proposition S.1.**  $f\phi_{x_i} g_{x_i \in \mathbf{x}}$  is a unique causal contribution measure w.r.t.  $f(y, \mathbf{x}g)$  on  $G$ .

*Proof.* It suffices to show that  $G[\mathbf{X}]$  is a graph corresponding to  $P(\mathbf{X})$ , because of  $\phi_{x_i}$  is the *do*-Shapley value defined on a graph corresponding to  $P(\mathbf{X})$ . By (Koster et al., 2002),  $G[\mathbf{X}]$  is a graph corresponding to  $P(\mathbf{X})$ .  $\square$

## C. Relation with Other Work - Examples

In this section, we provide examples to demonstrate that other types of Shapley values doesn't satisfy the Axiom 1. We first note that the conditional Shapley doesn't satisfy the causal irrelevance property in Axiom 1.

**Example C.1 (Causal Irrelevance Property doesn't hold for the conditional Shapley (Janzing et al., 2020b)).** Consider  $G = \overleftarrow{V_1} \text{---} V_2 \text{---} Y$  where  $V_1, V_2 \in \{0, 1\}$ , and the bidirected edge means the existence of hidden confounders. Suppose  $P(v_1, v_2) = 1/2$  whenever  $v_1 = v_2$ . Note  $V_1$  and  $Y$  is causally irrelevant. Causal irrelevance property doesn't hold in the conditional Shapley. Specifically, for any  $v_1, v_2$ ,  $E[Y|v_1] - E[Y] = v_1 - 1/2 \neq 0$ , which leads that  $\phi_{V_1}^{\text{cond}} \neq 0$ . In contrast,  $E[Y|do(v_1)] - E[Y] = E[Y|do(v_1, v_2)] - E[Y|do(v_2)] = 0$ . Therefore,  $\phi_{V_1} = 0$ , implying that *do*-Shapley satisfies the causal irrelevance property, unlike to the conditional Shapley.

The ICC doesn't satisfy the causal symmetry property in Axiom 1.

**Example C.2 (Causal Symmetry Property doesn't hold for the ICC Approach).** Consider a following SCM  $\mathcal{M}$ : For all binary variables  $U_{V_1}, U_{V_2}, U_Y, V_1, V_2, Y \in \{0, 1\}$ ,  $P(U_1 = 1) = 0.5$ ,  $P(U_2 = 1) = 0.2$ , and  $P(U_Y = 1) = 0.8$ . Also,  $V_1 = f_{V_1}(U_{V_1}) = U_{V_1}$ ;  $V_2 = f_{V_2}(V_1, U_2) = V_1 \oplus U_2$ ; and  $Y = f_Y(V_2, U_Y) = V_2 \oplus U_Y$ . A corresponding causal diagram is  $G = \overleftarrow{V_1} \text{---} V_2 \text{---} Y, \overleftarrow{U_{V_1}} \text{---} V_1, \overleftarrow{U_{V_2}} \text{---} V_2, \overleftarrow{U_Y} \text{---} Y$  for all  $V \in \{V_1, V_2, Y\}$ . Let  $y = v_1 = v_2 = 1$ . Then,  $P(y|do(v_1)) = P(y|v_1) = 0.8$ ,  $P(y|do(v_2)) = P(y|v_2) = 0.8$ ,  $P(y|do(v_1, v_2)) = P(y|v_2) = 0.8$ , and  $P(y) = 0.65$ . We first note that  $v_1$  and  $v_2$  have the same causal explanatory power to  $Y$  since  $P(y|do(v_1)) = P(y|do(v_2)) = 0.8$ . Also, the *do*-Shapley values for  $v_1, v_2$  are the same as  $\phi_{v_1} = \phi_{v_2} = 0.075$ , which exhibits the causal symmetry. To compute the ICC of the features  $v_1 = v_2 = 1$ , we fix  $u_1 = 1$  and  $u_2 = 0$ , which makes  $v_1 = v_2 = 1$ . Let  $\phi_{V_i}^{\text{icc}}$  denote the ICC of  $v_i$ . Then,  $\phi_{V_1}^{\text{icc}} = 0.225$  and  $\phi_{V_2}^{\text{icc}} = 0.075$  even if  $v_1, v_2$  have the same causal explanatory power. This implies that the causal symmetry doesn't hold.

## D. Proofs

We provide complete proofs and additional missing details here.

### D.1. Proofs from Section 3

We use

$$\nu_{do}(S) := E[Y|do(\mathbf{v}_S)]$$

in the proof.

**Theorem D.1** (Restated Theorem 1). *The *do*-Shapley is a unique causal contribution measure satisfying all the properties in Axiom 1.*

*Proof.* We first prove that *do*-Shapley satisfies all the properties in Axiom 1.

**Lemma S.1 (Soundness of *do*-Shapley).** *The *do*-Shapley satisfies all properties in Axiom 1.*

*Proof.* First, consider the **perfect assignment** property. By the result of (Štrumbelj & Kononenko, 2014), we can represent the *do*-Shapley as

$$\phi_{v_i}(\nu_{do}) = \frac{1}{n!} \sum_{\pi \in \Pi([n])} \nu_{do}(\text{pre}(\pi(i))) - \nu_{do}(\text{pre}(\pi(i-1)))$$

where  $\Pi([n])$  is a set of all possible permutations of  $[n]$ ,  $\pi$  is an individual permutation in  $\Pi([n])$ , and  $\text{pre}(\pi(i)) := \{k \in [n] \mid k < i \text{ in } \pi([n])\}$ . Then,

$$\begin{aligned} \sum_{i=1}^n \phi_{v_i}(\nu_{do}) &= \frac{1}{n!} \sum_{\pi \in \Pi([n])} \sum_{i=1}^n \nu_{do}(\text{pre}(\pi(i))) - \nu_{do}(\text{pre}(\pi(i-1))) \\ &= \frac{1}{n!} \sum_{\pi \in \Pi([n])} \nu_{do}(\pi(n)) - \nu_{do}(\pi(1)) \\ &= \nu_{do}([n]) - \nu_{do}(\emptyset) = E[Y|do(\mathbf{v})] - E[Y] = E[Y|do(\mathbf{v})]. \end{aligned}$$

Now we consider the **causal irrelevance** property. Suppose  $V_i$  is causally irrelevant to  $Y$  in expectation for all witness  $w \models \nu_{do}^i$ . Then, the equality  $\nu_{do}(S \setminus i) = \nu_{do}(S) = 0$  holds immediately for all  $S \subseteq [n] \setminus i$ .

Next we consider the **causal symmetry** property. Suppose  $v_i, v_j$  has the same causal explanatory power w.r.t. any witnesses  $w \models \nu_{do}^i, \nu_{do}^j$ . This leads  $\nu_{do}(S \setminus i) = \nu_{do}(S \setminus j)$  for any  $S \subseteq [n] \setminus \{i, j\}$ . Then,

$$\begin{aligned} v_i(\nu_{do}) &= \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i\}} \binom{n-1}{|S|}^{-1} \nu_{do}(S \setminus i) - \nu_{do}(S) \\ &= \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i, j\}} \binom{n-1}{|S|}^{-1} \nu_{do}(S \setminus i) - \nu_{do}(S) + \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i, j\}} \binom{n-1}{|S|+1}^{-1} \nu_{do}(S \setminus i; j) - \nu_{do}(S \setminus i; j) \\ &= \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i, j\}} \binom{n-1}{|S|}^{-1} \nu_{do}(S \setminus j) - \nu_{do}(S) + \frac{1}{n} \sum_{S \subseteq [n] \setminus \{i, j\}} \binom{n-1}{|S|+1}^{-1} \nu_{do}(S \setminus i; j) - \nu_{do}(S \setminus i; j) \\ &= \frac{1}{n} \sum_{S \subseteq [n] \setminus \{j\}} \binom{n-1}{|S|}^{-1} \nu_{do}(S \setminus j) - \nu_{do}(S) = v_j(\nu_{do}); \end{aligned}$$

where the third equality holds since  $(v_i, v_j)$  has the same causal explanatory power.

Now, we prove that the *do*-Shapley satisfies the **causal approximation** property by showing that there exists  $\omega(S)$  that makes the *do*-Shapley as the solution of the weighted least square problem defined in Axiom 1. For a coalition function  $\nu(S)$  (see the ‘‘Shapley value’’ paragraph in Sec. 2), it’s known that there exists a specific weight function  $\omega(S)$  that makes the Shapley value in Eq. (1) as the solution of the following WLS problem:  $\arg \min_{\omega} \sum_{S \subseteq [n]} (\nu(S) - \sum_{i \in S} \phi_{V_i}^\omega)^2 \omega(S)$  (by (Charnes et al., 1988, Thm. 4) and (Lundberg & Lee, 2017, Theorem 2)). This implies that such an  $\omega(S)$  is the weight function that makes the *do*-Shapley as the solution of the weighted least square problem defined in Axiom 1.  $\square$

We now show the other direction that a measure satisfying all properties in Axiom 1 is the *do*-Shapley.

**Lemma S.2 (Completeness of *do*-Shapley).** *A vector  $\phi_{V_i} \in \mathbb{R}^{\mathcal{S}}$  satisfying Axiom 1 is the *do*-Shapley value.*

*Proof.* Throughout the proof, we will define a *canonical SCM* as follow: Let  $T \subseteq [n]$  denote any fixed index set. A SCM is called *canonical* for  $T$  if  $\mathbb{E}[Y | do(\mathbf{V}_S = 1)] = 1$  iff  $T \subseteq S$ , and 0 otherwise. We use  $\nu_{do}^T(S)$  denote the causal coalition function induced by the canonical SCM. Note  $\nu_{do}^T(S) = 1$  iff  $T \subseteq S$ , and 0 otherwise, by the definition of the canonical SCM.

We first note that a vector  $\phi_{V_i}$  that satisfies the *causal approximation* property can be represented as a linear function of  $\nu_{do}(S)$ , because  $\phi_{V_i}$  is a solution of the weighted least square linear regression problem. Therefore,

$$\phi_{V_i} = \sum_{S \subseteq [n]} a_S^i \nu_{do}(S). \quad (\text{D.1})$$

for some constants  $a_S^i$ .

Now we focus on the **causal irrelevance** property. Suppose  $T \subseteq [n] \setminus i$ . For any  $S \subseteq [n]$ ,  $(T \subseteq S) \Rightarrow (T \subseteq S \setminus i)$ . With  $i \notin T$ ,  $(T \subseteq S) \Rightarrow (T \subseteq S \setminus i)$ . Therefore,  $\nu_{do}^T(S) = \nu_{do}^T(S \setminus i)$  for all  $S \subseteq [n]$ . Then, by the causal irrelevance property,  $\phi_{V_i}(\nu_{do}^T) = 0$  if  $T \subseteq [n] \setminus i$ . Then,  $\phi_{V_i}(\nu_{do}^{[n] \setminus i}) = a_{[n]}^i + a_{[n] \setminus i}^i = 0$ .

Suppose it has been shown that  $a_{T \setminus i}^i + a_T^i = 0$  for  $T \subseteq [n] \setminus i$  such that  $|T| = k$  for some  $k$ . Then, for any  $S \subseteq [n] \setminus i$  such that  $|S| = k-1$ ,

$$\begin{aligned} \phi_{V_i}(\nu_{do}^S) &= \sum_{T \subseteq [n]} a_T^i \nu_{do}^S(T) = \sum_{\substack{T \subseteq [n] \\ T \subseteq S}} a_T^i = \sum_{\substack{T \subseteq [n] \setminus i \\ T \subseteq S}} (a_{T \setminus i}^i + a_T^i) \\ &= \left\{ \sum_{\substack{T \subseteq [n] \setminus i \\ T \subseteq S \text{ but } T \not\subseteq S}} (a_{T \setminus i}^i + a_T^i) \right\} + (a_{S \setminus i}^i + a_S^i) = a_{S \setminus i}^i + a_S^i, \end{aligned}$$

where the first equality by Eq. (D.1), the second by the property of the canonical SCM, the third and fourth by the standard algebra, and the fifth by the inductive hypothesis. Since  $S \setminus [n]n\bar{f}i\bar{g}$ , by causal irrelevance property,  $\phi_{v_i}(\nu_{do}^S) = 0$ . This implies that  $a_{S \setminus [n]n\bar{f}i}^i + a_S^i = 0$ . Therefore, for any  $T \setminus [n]n\bar{f}i$ ,  $a_{T \setminus [n]n\bar{f}i}^i + a_T^i = 0$ .

Fix  $p_T^i := a_{T \setminus [n]n\bar{f}i}^i = -a_T^i$ . Then,

$$\phi_{v_i}(\nu_{do}) = \sum_{T \setminus [n]n\bar{f}i} a_T^i \nu_{do}(T) = \sum_{T \setminus [n]n\bar{f}i} (a_{T \setminus [n]n\bar{f}i}^i \nu_{do}(T \setminus [n]n\bar{f}i) + a_T^i \nu_{do}(T)) = \sum_{T \setminus [n]n\bar{f}i} p_T^i (\nu_{do}(T \setminus [n]n\bar{f}i) - \nu_{do}(T)).$$

Now we focus on the **causal symmetry** property. Suppose  $v_i$  and  $v_j$  have the same causal explanatory power with any given witness  $\mathbf{w} \setminus \{v_i, v_j\}$  in the canonical SCM for  $[n]$ ; i.e.,  $\nu_{do}^{[n]}(S \setminus [n]n\bar{f}i) = \nu_{do}^{[n]}(S \setminus [n]n\bar{f}j)$  for  $S \setminus [n]n\bar{f}i, j\bar{g}$ . We note that  $\phi_{v_i}(\nu_{do}^{[n]}) = p_{[n]n\bar{f}i}^i = \phi_{v_j}(\nu_{do}^{[n]}) = p_{[n]n\bar{f}j}^j$ . This implies that there exists  $p_{n-1} := p_{[n]n\bar{f}i}^i = p_{[n]n\bar{f}j}^j$ .

Again, suppose  $v_i, v_j$  have the same causal explanatory power with any given witness  $\mathbf{w} \setminus \{v_i, v_j\}$  in the canonical SCM for  $[n]n\bar{k}$  for any fixed  $k \notin \bar{f}i, j\bar{g}$ . Then,

$$\phi_{v_i}(\nu_{do}^{[n]n\bar{k}}) = p_{[n]n\bar{f}i;k\bar{g}}^i + p_{n-1} = \phi_{v_j}(\nu_{do}^{[n]n\bar{k}}) = p_{[n]n\bar{f}j;k\bar{g}}^j + p_{n-1}.$$

This implies that there exists a constant  $p_{n-2} := p_{[n]n\bar{f}i;k\bar{g}}^i = p_{[n]n\bar{f}j;k\bar{g}}^j$ . By repeating this, we can have a  $p_1, \dots, p_{n-1}$  where  $p_{jTj}$  is a constant applying to all  $p_T^j$  for any  $T \setminus [n]n\bar{f}i$ . Therefore, there are constants  $\bar{p}_{jTj} p_{jTj}$  such that

$$\phi_{v_i} := \sum_{T \setminus [n]n\bar{f}i} p_{jTj} (\nu_{do}(T \setminus [n]n\bar{f}i) - \nu_{do}(T)).$$

Finally, we focus on the **perfect assignment** property. An attribution  $\phi_{v_i}$  satisfies the perfect assignment property if and only if  $\sum_{i \geq 1} p_{n-1} = 1$ , and for any nonempty  $T \setminus [n]$ ,  $\sum_{i \geq 1} p_{jTj} = \sum_{j \notin T} p_{jTj}$  (Winter, 2002, Chap. 7, Theorem 11). This gives  $p_{n-1} = 1/n$ , and for any nonempty  $T \setminus [n]$ ,  $\sum_{j \notin T} p_{jTj} = (n - |T|) p_{jTj}$ . Then, a closed form for  $p_{jTj}$  is given as

$$p_{jTj} = \frac{(n - |T|)! |T|!}{n!} = \frac{1}{n} \binom{n-1}{|T|}.$$

□

Taking a conjunction of Lemmas (S.1,S.2) completes the proof of the Theorem D.1. □

## D.2. Proofs from Section 4

**Theorem D.2** (Restated Theorem 2). *The do-Shapley is identifiable if no variable in  $V_i \setminus \{v_i\}$  is connected to its child  $\text{Ch}(V_i)$  by bidirected paths in  $G$ . Suppose  $Y$  is not connected by bidirected paths. In this case, for any  $S \setminus [n]$ ,*

$$\mathbb{E}[Y | do(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} \mathbb{E}[Y | \mathbf{v}] Q[\mathbf{V} | \mathbf{v}_S],$$

where  $Q[\mathbf{V} | \mathbf{v}_S] := Q[\mathbf{V} | \mathbf{v}_S](\mathbf{v})$  is given as

$$Q[\mathbf{V} | \mathbf{v}_S] = \frac{P(\mathbf{v})}{Q[C(\mathbf{v}_S)]} \prod_{k=1}^c \sum_{\mathbf{s}_k} Q[C(\mathbf{s}_k)],$$

where  $Q[C(\mathbf{v}_S)] = \prod_{V_a \in C(\mathbf{v}_S)} P(v_a | \text{pre}(v_a))$  is a  $C$ -factor of a  $C$ -component  $\mathbf{V}_S$  ( $C(\mathbf{v}_S)$ );  $\{ \mathbf{s}_k \}_{k=1}^c$  is a  $C$ -partition of  $\mathbf{V}_S$ ; and  $Q[C(\mathbf{s}_k)] := \prod_{V_a \in C(\mathbf{s}_k)} P(v_a | \text{pre}(v_a))$  is a  $C$ -factor of a  $C$ -component  $C(\mathbf{s}_k)$  for  $\mathbf{s}_k$ .

*Proof.* We prove the following, which would imply the above theorem.

**Proposition S.1 (Generalized Tian’s Adjustment – Complete identification criteria for  $P(\mathbf{V}jdo(\mathbf{X}))$ ).**  $P(\mathbf{V}jdo(\mathbf{X}))$  is identifiable if  $\delta X_a \not\subseteq \mathbf{X}$  and  $\text{Ch}(X_a)$  is not connected by bidirected paths. If identifiable, it’s given as

$$P(\mathbf{V}jdo(\mathbf{X})) = \frac{P(\mathbf{V})}{Q[C(\mathbf{X})]} \prod_{k=1}^c \sum_{\mathbf{x}_k} Q[C(\mathbf{X}_k)], \quad (\text{D.2})$$

where  $\{C_k\}_{k=1}^c$  is a  $C$ -partition of  $\mathbf{X}$  in  $G$ , and  $Q[C(\mathbf{X})] := \prod_{V_i \in C(\mathbf{X})} P(V_i|pre(V_i))$  is a  $C$ -factor of a  $C$ -component of  $\mathbf{X}$  ( $C(\mathbf{X})$ ), and  $Q[C(\mathbf{X}_k)] := \prod_{V_i \in C(\mathbf{X}_k)} P(V_i|pre(V_i))$  is a  $C$ -factor of a  $C$ -component of  $\mathbf{X}_k$  ( $C(\mathbf{X}_k)$ ).

*Proof.* In the proof, for a vector  $\mathbf{W}$ , we will use  $\text{De}(\mathbf{W})$  to denote a set of descendants of  $W_i \in \mathbf{W}$  in  $G$ .

Suppose  $\delta X_a \not\subseteq \mathbf{X}$  is not connected with  $\text{Ch}(X_a)$  by bidirected paths. We first show that  $P(\mathbf{V}jdo(\mathbf{X}_1))$  (for any  $\mathbf{X}_1 \subseteq \mathbf{X}$  a  $C$ -component in  $G(\mathbf{X})$ ) is identifiable and given as

$$P(\mathbf{V}jdo(\mathbf{X}_1)) = \frac{P(\mathbf{V})}{Q[C(\mathbf{X}_1)]} \sum_{\mathbf{x}_1} Q[C(\mathbf{X}_1)]. \quad (\text{D.3})$$

By the result of (Jaber et al., 2018, Lemma 1), it suffices to show that  $\mathbf{X}_1 = \text{De}(\mathbf{X}_1)_{G(C(\mathbf{X}_1))}$ . We show this by contradiction. Suppose  $V_a \in \text{De}(\mathbf{X}_1)_{G(C(\mathbf{X}_1))}$  such that  $V_a \notin \mathbf{X}_1$ . Since  $V_a \in G(C(\mathbf{X}_1))$ ,  $V_a$  is connected with  $\mathbf{X}_1$  by bidirected paths. Since  $V_a$  is a descendent of some  $X_a \in \mathbf{X}_1$  in  $G(C(\mathbf{X}_1))$ , this means that  $V_b \in \text{Ch}(X_a)$  is also in the  $G(C(\mathbf{X}_1))$ . This means that  $V_b$  and  $X_a$  is connected by a bidirected path, which is a contradiction of the given condition. Therefore,  $\mathbf{X}_1 = \text{De}(\mathbf{X}_1)_{G(C(\mathbf{X}_1))}$ , and Eq. (D.3) holds.

Now, consider a following inductive hypothesis for  $i = 1, 2, \dots, c$ :

$$Q[\mathbf{V}n\mathbf{X}^{(i)}] = \frac{Q[\mathbf{V}n\mathbf{X}^{(i-1)}]}{Q[C(\mathbf{X}_i)]} \sum_{\mathbf{x}_i} Q[C(\mathbf{X}_i)]. \quad (\text{D.4})$$

As shown in the above, it holds for  $i = 1$ . Suppose it holds for some  $i - 1 \leq i - 1$  for  $i \geq 2$ . Then, we first note that  $\mathbf{X}_i = \text{De}(\mathbf{X}_i)_{G(C(\mathbf{X}_i))_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}}$ . To witness, consider the contradiction – for some  $X_a \in \mathbf{X}_i$  there exists  $V_a \in \text{De}(\mathbf{X}_i)_{G(C(\mathbf{X}_i))_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}}$  s.t.  $V_a \notin \mathbf{X}_i$ . First,  $V_a$  is connected with  $X_a$  by bidirected paths since  $V_a \in G(C(\mathbf{X}_i))_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}$ . Also,  $V_a$  is a descendent of  $X_a$ , this means that a child of  $X_a$  is also in  $G(C(\mathbf{X}_i))$ , connected by bidirected paths. This is a contradiction. Therefore,  $\mathbf{X}_i = \text{De}(\mathbf{X}_i)_{G(C(\mathbf{X}_i))_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}}$ .

Now, we show that  $C(\mathbf{X}_i)_{G(\mathbf{V}n\mathbf{X}^{(i-1)})} = C(\mathbf{X}_i)_G$ . We start from an obvious observation –  $C(\mathbf{X}_i)_{G(\mathbf{V}n\mathbf{X}^{(i-1)})} \subseteq C(\mathbf{X}_i)_G$ . We now prove  $C(\mathbf{X}_i)_G \subseteq C(\mathbf{X}_i)_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}$ . For some  $V_a \in C(\mathbf{X}_i)_G$ , suppose  $V_a \notin C(\mathbf{X}_i)_{G(\mathbf{V}n\mathbf{X}^{(i-1)})}$ . This means that bidirected paths connecting  $V_a$  to some nodes in  $\mathbf{X}_1 \subseteq \mathbf{X}_i$  must be via other nodes in  $\mathbf{X}_2 \subseteq \mathbf{X}^{(i-1)}$ . This means that  $V_a, \mathbf{X}_2, \mathbf{X}_1$  are connected by bidirected paths. However, given that  $\mathbf{X}_2 \subseteq \mathbf{X}^{(i-1)}$  and  $\mathbf{X}_1 \subseteq \mathbf{X}_i$ , this is a contradiction, because they are in distinct  $C$ -partitions. Therefore,  $C(\mathbf{X}_i)_{G(\mathbf{V}n\mathbf{X}^{(i-1)})} = C(\mathbf{X}_i)_G$ .

Then, Eq. (D.4) holds. By unfolding it,

$$Q[\mathbf{V}n\mathbf{X}^{(i)}] = \frac{P(\mathbf{V})}{\prod_{k=1}^i Q[C(\mathbf{X}_k)]} \prod_{k=1}^c \sum_{\mathbf{x}_k} Q[C(\mathbf{X}_k)].$$

We note  $Q[C(\mathbf{X}^{(i)})] = \prod_{k=1}^i Q[C(\mathbf{X}_k)]$ , since

$$Q[C(\mathbf{X}^{(i)})] = \prod_{V_i \in C(\mathbf{X}^{(i)})} P(v_i|pre(v_i)) = \prod_{k=1}^i \prod_{V_i \in C(\mathbf{X}_k)} P(v_i|pre(v_i)) = \prod_{k=1}^i Q[C(\mathbf{X}_k)].$$

This completes the proof.  $\square$

Now back to witness Thm. D.2. Under the given condition that  $Y$  is not connected via bidirected paths to any nodes, the following holds: for any  $S \subseteq [n]$ ,

$$(Y \perp\!\!\!\perp \mathbf{V}_S | \mathbf{V}_{\bar{S}})_{G_{\mathbf{V}_S}}.$$



Therefore,

$$P(Y, \mathbf{V} \text{ jdo}(\mathbf{V}_S)) = P(Y \text{ jdo}(\mathbf{V}_S), \mathbf{V}_{\bar{S}}) Q[\mathbf{V}_S] = P(Y \text{ j} \mathbf{V}) Q[\mathbf{V}_S],$$

which implies that

$$E[Y \text{ jdo}(\mathbf{V}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Y \text{ j} \mathbf{V}] \frac{Q[\mathbf{V}]}{Q[C(\mathbf{V}_S)]} \prod_{k=1}^c \sum_{\mathbf{x}_k} Q[C(\mathbf{x}_k)].$$

This completes the proof. □

**Corollary D.2** (Restated Corollary 1). *In the Markovian case,  $E[Y \text{ jdo}(\mathbf{v}_S)]$  is given as*

$$E[Y \text{ jdo}(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Y \text{ j} \mathbf{v}_S, \mathbf{v}_{\bar{S}}] \prod_{i \in \bar{S}} P(v_i \text{ jpre}(v_i)).$$

*Proof.* In the Markovian case,  $C(\mathbf{W}) = \mathbf{W}$  for all  $\mathbf{W} \subseteq \mathbf{V}$ . Then,

$$P(Y, \mathbf{V}_{\bar{S}} \text{ jdo}(\mathbf{V}_S)) = \frac{P(\mathbf{V}, Y)}{Q[C(\mathbf{V}_S)]} \prod_{k=1}^c \sum_{\mathbf{x}_k} Q[C(\mathbf{x}_k)] = \frac{P(\mathbf{V}, Y)}{Q[\mathbf{V}_S]} = P(Y \text{ j} \mathbf{V}) \prod_{v_i \in \mathbf{V}_{\bar{S}}} P(V_i \text{ jpre}(V_i)).$$

This completes the proof. □

**Corollary D.2** (Restated Corollary 2). *In the Direct-cause case,  $E[Y \text{ jdo}(\mathbf{v}_S)]$  is given as*

$$E[Y \text{ jdo}(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} E[Y \text{ j} \mathbf{v}_S, \mathbf{v}_{\bar{S}}] P(\mathbf{v}_{\bar{S}}).$$

*Proof.* In the Direct-cause case,  $Q[\mathbf{W}] = P(\mathbf{W})$  for all  $\mathbf{W} \subseteq \mathbf{V}$  since there are no causal paths between a pair of variables in  $\mathbf{V}$ . Therefore,  $Q[\mathbf{V} \cap \mathbf{V}_S] = P(\mathbf{V} \cap \mathbf{V}_S) = P(\mathbf{V}_{\bar{S}})$ , which completes the proof. □

### D.3. Proofs from Section 5

**Lemma D.1** (Restated Lemma 1). *Let  $S = \{m_1, \dots, m_s\} \subseteq [n]$  denote an index set for  $\mathbf{V}_S$ . Let*

$$\begin{aligned} \omega_k^S &:= \prod_{r=1}^k \mathbb{1}_{V_{m_r}}(V_{m_r}) / h_r^S, \text{ for } k = s, \dots, 1; \\ \omega^S &:= \mathbb{1}_{\mathbf{v}_S}(\mathbf{V}_S) / h^S, \end{aligned}$$

where  $h_r^S := P(V_{m_r} \text{ jpre}(V_{m_r}))$  and  $h^S := P(\mathbf{V}_S \text{ j} \mathbf{V}_{\bar{S}})$ . Then,  $E[Y \text{ jdo}(\mathbf{v}_S)] = E[Y \omega]$  where  $\omega = \omega_k^S$  for the Markovian case, and  $\omega = \omega^S$  for the Direct-cause case.

*Proof.* For the Markovian case,

$$E[Y \text{ jdo}(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} \left( \frac{P(\mathbf{v})}{\prod_{r=1}^s P(v_{m_r} \text{ jpre}(v_{m_r}))} \right) = E \left[ \frac{\mathbb{1}_{V_{m_r}}(V_{m_r})}{h_r^S} \right].$$

For the Direct-cause case,

$$E[Y \text{ jdo}(\mathbf{v}_S)] = \sum_{\mathbf{v}_{\bar{S}}} \frac{P(\mathbf{v})}{P(\mathbf{v}_S \text{ j} \mathbf{v}_{\bar{S}})} = E \left[ \frac{\mathbb{1}_{\mathbf{v}_S}(\mathbf{V}_S)}{P(\mathbf{V}_S \text{ j} \mathbf{V}_{\bar{S}})} \right].$$

□

**Lemma D.2** (Restated Lemma 2). Let  $S := \{m_1, \dots, m_s\} \subseteq [n]$  denote an index set for  $\mathbf{V}_S$ . Let  $\theta_{S,1}^S := Y$ . For  $k = s, s-1, \dots, 1$ ,

$$\begin{aligned}\theta_{k,2}^S &:= E[\theta_{k,1}^S | V_{m_k}, \text{pre}(V_{m_k})] \\ \theta_{k-1,1}^S &:= E[\theta_{k,1}^S | v_{m_k}, \text{pre}(V_{m_k})], \\ \theta_a^S &:= E[Y | \mathbf{v}_S, \mathbf{V}_{\bar{S}}], \\ \theta_b^S &:= E[Y | \mathbf{V}_S, \mathbf{V}_{\bar{S}}].\end{aligned}$$

Then,  $E[Y | do(\mathbf{v}_S)] = E[\theta]$  where  $\theta = \theta_{0,1}^S$  for the Markovian case, and  $\theta = \theta_a^S$  for the Direct-cause case.

*Proof.* For the Markovian case, we will prove the following, which implies the result.

**Lemma S.3.** Suppose  $\mathbf{V}^\theta = fYg[\mathbf{V}]$  where  $\mathbf{V}^\theta$  is an ordered set. Assume that  $Y$  is the last variable in the given order. Let  $\mathbf{V}_S := \{V_{m_1}, \dots, V_{m_s}\} \subseteq \mathbf{V}$  (where  $\{m_1, \dots, m_s\} \subseteq [n]$ ) denote a set of discrete variables. Let  $\mathbf{V}_{\bar{S}} := \mathbf{V} \setminus \mathbf{V}_S$ . For each  $k = 2, \dots, s$ , let  $\mathbf{V}_{\cdot k} := \{V_j \in \mathbf{V}_{\bar{S}} : V_{m_k} \prec V_j \prec V_{m_k}\} \subseteq \mathbf{V}_{\bar{S}}$ . Let  $\mathbf{V}_{\cdot 1} := \{V_j \in \mathbf{V}_{\bar{S}} : V_j \prec V_{m_1}\} \subseteq \mathbf{V}_{\bar{S}}$  and  $\mathbf{V}_{\cdot s+1} := \{V_j \in \mathbf{V}_{\bar{S}} : V_{m_s} \prec V_j\} \subseteq \mathbf{V}_{\bar{S}}$ .

Let  $g_S(P)$  denote a following functional (a.k.a.  $g$ -formula (Robins, 1986)).

$$g_S(P) := \int_{\mathbf{V}_{\bar{S}}} E[Y | \mathbf{v}] \prod_{i \in \bar{S}} P(v_i | \text{pre}(v_i)) d[\mathbf{v}_{\bar{S}}].$$

Let  $\theta_{S,1} := Y$ . For  $k = s, \dots, 1$ , and

$$\begin{aligned}\theta_{k,2} &:= E[\theta_{k,1} | V_{m_k}, \text{pre}(V_{m_k})] \\ \theta_{k-1,1} &:= E[\theta_{k,1} | v_{m_k}, \text{pre}(V_{m_k})].\end{aligned}$$

Then, the following holds:

$$g_S(P) = E[\theta_{0,1}].$$

*Proof.* Let

$$\begin{aligned}\mathbf{A}_k &:= \text{pre}(V_{m_k}) \\ \mathbf{B}_k &:= \{V_{\cdot k+1}, \mathbf{V}_{\cdot k+2}, \dots, \mathbf{V}_{\cdot s+1}\} \\ \mathbf{C}_k &:= \{V_{m_{k+1}}, V_{m_{k+2}}, \dots, V_{m_s}\}.\end{aligned}$$

For  $\mathbf{W} \subseteq \mathbf{V}$ ,

$$q(\mathbf{W}) := \begin{cases} \prod_{V_i \in \mathbf{W}} P(v_i | \text{pre}(v_i)) & \text{If } \mathbf{W} \subseteq \mathbf{V}_{\bar{S}}; \\ 1 & \text{If } \mathbf{W} \subseteq \mathbf{V}_S. \end{cases}$$

Then, it suffices to show that

$$\begin{aligned}\theta_{k,2} &= \int_{\mathbf{B}_k; \mathbf{C}_k} E[Y | V_{m_k}, \mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k] q(\mathbf{b}_k) \mathbb{1}_{\mathbf{c}_k}(\mathbf{C}_k) d[\mathbf{b}_k, \mathbf{c}_k] \\ \theta_{k-1,1} &= \int_{\mathbf{B}_k; \mathbf{C}_k} E[Y | \mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k] q(\mathbf{b}_k) \mathbb{1}_{\mathbf{c}_k}(\mathbf{C}_k) d[\mathbf{b}_k, \mathbf{c}_k],\end{aligned}$$

because witnessing  $E[\theta_{0,1}] = g_S(P)$  becomes trivial. Let  $\theta_{S,1} := Y$ . Then, it's easy to check that the above holds for  $\theta_{s,2}$  and  $\theta_{s-1,1}$ .

Suppose the above equation holds for  $k, k+1, \dots, s$ . Then, consider  $k-1$ . By the given definition,

$$\begin{aligned}\theta_{k-1,2} &:= E[\theta_{k-1,1} | V_{m_{k-1}}, \text{pre}(V_{m_{k-1}})] \\ \theta_{k-2,1} &:= E[\theta_{k-1,1} | v_{m_{k-1}}, \text{pre}(V_{m_{k-1}})].\end{aligned}$$

Then,

$$\begin{aligned}\theta_{k-1,2} &= \mathbb{E} \left[ \int_{B_k; C_{k-1}} \mathbb{E} [Yj\mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_{k-1}] q(\mathbf{b}_k) \mathbb{1}_{\mathbf{c}_{k-1}}(\mathbf{C}_{k-1}) d[\mathbf{b}_k, \mathbf{c}_{k-1}] \middle| V_{m_{k-1}}, \mathbf{A}_{k-1} \right] \\ &= \int_{B_{k-1}; C_{k-1}} \mathbb{E} [YjV_{m_{k-1}}, \mathbf{A}_{k-1}, \mathbf{b}_{k-1}, \mathbf{c}_{k-1}] q(\mathbf{b}_{k-1}) \mathbb{1}_{\mathbf{c}_{k-1}}(\mathbf{C}_{k-1}) d[\mathbf{b}_{k-1}, \mathbf{c}_{k-1}],\end{aligned}$$

and

$$\begin{aligned}\theta_{k-2,1} &= \mathbb{E} \left[ \int_{B_k; C_k} \mathbb{E} [Yj\mathbf{A}_k, \mathbf{b}_k, \mathbf{c}_k] q(\mathbf{b}_k) \mathbb{1}_{\mathbf{c}_k}(\mathbf{C}_k) d[\mathbf{b}_k, \mathbf{c}_k] \middle| v_{m_{k-1}}, \mathbf{A}_{k-1} \right] \\ &= \int_{B_{k-1}; C_{k-2}} \mathbb{E} [Yj\mathbf{A}_{k-1}, \mathbf{b}_{k-1}, \mathbf{c}_{k-2}] q(\mathbf{b}_{k-1}) \mathbb{1}_{\mathbf{c}_{k-2}}(\mathbf{C}_{k-2}) d[\mathbf{b}_{k-1}, \mathbf{c}_{k-2}].\end{aligned}$$

Therefore,

$$\theta_{0,1} = \int_{B_1; C_0} \mathbb{E} [Yj\mathbf{A}_1, \mathbf{b}_1, \mathbf{c}_0] q(\mathbf{b}_1) \mathbb{1}_{\mathbf{c}_0}(\mathbf{C}_0) d[\mathbf{b}_1, \mathbf{c}_0],$$

which gives the equality  $\mathbb{E} [\theta_{0,1}] = g_S(P)$ .  $\square$

For the Direct-cause case,

$$\mathbb{E} [\theta_a^S] = \sum_{\mathbf{v}_S} \mathbb{E} [Yj\mathbf{v}] P(\mathbf{v}_S) = \mathbb{E} [Yjdo(\mathbf{v}_S)],$$

which completes the proof.  $\square$

**Lemma D.3** (Restated Lemma 3). *Let*

$$\eta^S := \begin{cases} f\theta_{0,1}^S g [ f\theta_{k,1}^S, \theta_{k,2}^S g_{k=1}^S [ fh_r^S g_{r=1}^S \text{ (Markovian)} \\ f\theta_a^S, \theta_b^S, h^S g \text{ (Direct-cause)}, \end{cases}$$

defined in Defs. (4, 5) above, and

$$V_S(\mathbf{V}^0; \eta^S) := \begin{cases} \theta_{0,1}^S + \sum_{k=1}^s \omega_k^S (\theta_{k,1}^S \quad \theta_{k,2}^S) \text{ (Markovian)} \\ \theta_a^S + \omega^S (Y \quad \theta_b^S) \text{ (Direct-cause)}, \end{cases}$$

where  $\omega_k^S := \prod_{r=1}^k \mathbb{1}_{V_{m_r}}(V_{m_r})/h_r^S$  and  $\omega^S := \mathbb{1}_{\mathbf{v}_S}(\mathbf{V}_S)/h^S$ . Then,  $\mathbb{E} [Yjdo(\mathbf{v}_S)] = \mathbb{E} [V_S(\mathbf{V}^0; \eta^S)]$ .

*Proof.* For the Markovian case, it suffices to show that  $\mathbb{E} [\theta_{k,1}^S \quad \theta_{k,2}^S]$  for any  $k = 1, 2, \dots, s$ . This holds since

$$\mathbb{E} [\theta_{k,1}^S \quad \theta_{k,2}^S] = \mathbb{E} [\mathbb{E} [\theta_{k,1}^S \quad \theta_{k,2}^S | V_{m_k}, \text{pre}(V_{m_k})]] = \mathbb{E} [\theta_{k,2}^S \quad \theta_{k,2}^S] = 0.$$

Therefore,  $\mathbb{E} [V(\mathbf{V}^0; \eta^S)] = \mathbb{E} [\theta_{0,1}^S] = \mathbb{E} [Yjdo(\mathbf{v}_S)]$ , where the 2nd equality holds by Lemma 2.

For the Direct-cause case,  $\mathbb{E} [Y \quad \theta_b^S] = 0$  by the definition of  $\theta_b^S$ . Therefore,  $\mathbb{E} [V(\mathbf{V}^0; \eta^S)] = \mathbb{E} [\theta_a^S] = \mathbb{E} [Yjdo(\mathbf{v}_S)]$ .  $\square$

**Theorem D.3** (Restated Theorem 3). *Let  $f\pi_j g_{j=1}^M$  denote  $M$  randomly generated permutations of  $[n]$ . For the fixed index  $i$ , let  $S_{j,0} := \text{pre}_j(i)$  and  $S_{j,1} := f\hat{\eta} g [ S_{j,a}$ . Let  $f\hat{\eta}^{S_{j,0}}, \hat{\eta}^{S_{j,1}} g_{j=1}^M$  denote  $L_2$ -consistent estimates for all nuisances  $f\hat{\eta}^{S_{j,0}}, \eta^{S_{j,1}} g_{j=1}^M$  defined in Def. 6. Let  $R_{M,N} := O_P(M^{-1/2} + N^{-1/2})$ . Let  $e(\hat{g}) := k\hat{g} - g$  denote an error for a nuisance estimates for any  $\hat{g} \geq \hat{\eta}$  and  $g \geq \eta$ . For the do-Shapley estimators defined in Def. 7, suppose the estimators  $T^{\text{est}}(S)$  are bounded. Let  $\epsilon_{V_i}^{\text{est}} := \phi_{V_i}^{\text{est}} - \phi_{V_i}$  (where  $\text{est} \geq \text{fipw, reg, dmlg}$ ).*

Under the Markovian case,

$$\begin{aligned}\epsilon_{V_i}^{\text{ipw}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M e(\hat{\omega}_{S_j^{j,p}}^{S_j,p}) g, \\ \epsilon_{V_i}^{\text{reg}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M e(\hat{\theta}_{0;1}^{S_j,p}) g, \\ \epsilon_{V_i}^{\text{dml}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M \sum_{k=1}^{S_j} e(\hat{h}_k^{S_j,p}) e(\hat{\theta}_{k;2}^{S_j,p}) g.\end{aligned}$$

Under the Direct-cause case,

$$\begin{aligned}\epsilon_{V_i}^{\text{ipw}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M e(\hat{\omega}^{S_j,p}) g, \\ \epsilon_{V_i}^{\text{reg}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M e(\hat{\theta}_2^{S_j,p}) g, \\ \epsilon_{V_i}^{\text{dml}} &= R_{M;N} + O_P f \sum_{p \geq 2} \sum_{f_0;1} \sum_{g j=1}^M e(\hat{h}^{S_j,p}) e(\hat{\theta}_b^{S_j,p}) g.\end{aligned}$$

*Proof.* In the proof, we will use a notation  $E_{D \perp P} [f(\mathbf{V})]$  for  $f(\mathbf{V}) := E_D [f(\mathbf{V})] - E [f(\mathbf{V})]$ . We use  $N := jDj$ . Also, for any quantity  $A, B$ ,  $A \asymp B$  if there is a constant  $c$  s.t.  $A \leq cB$ . We first introduce a useful tool for analyzing errors of the proposed estimator.

**Lemma S.4.** *Let  $\eta_0$  denote some nuisance and  $\hat{\eta}$  denote its  $L_2$  consistent estimate. Let  $f(\mathbf{V}; \eta)$  denote an arbitrary function having a bounded second moment for any fixed  $\eta$ . Suppose samples used for constructing  $\hat{\eta}$  and for evaluating  $f(\mathbf{V}; \hat{\eta})$  are independent.*

$$E_D [f(\mathbf{V}; \hat{\eta})] - E [f(\mathbf{V}; \eta_0)] = O_P(N^{-1/2}) + E [f(\mathbf{V}; \hat{\eta}) - f(\mathbf{V}; \eta_0)].$$

*Proof.* We first note that

$$E_D [f(\mathbf{V}; \hat{\eta})] - E [f(\mathbf{V}; \eta_0)] = E_{D \perp P} [f(\mathbf{V}; \eta_0)] - E_{D \perp P} [f(\mathbf{V}; \hat{\eta}) - f(\mathbf{V}; \eta_0)] + E [f(\mathbf{V}; \hat{\eta}) - f(\mathbf{V}; \eta_0)].$$

First,  $E_{D \perp P} [f(\mathbf{V}; \eta_0)] = O_P(N^{-1/2})$  by the classical central limit theorem. Second,  $E_{D \perp P} [f(\mathbf{V}; \hat{\eta}) - f(\mathbf{V}; \eta_0)] = O_P(N^{-1/2})$  under given conditions by (Kennedy et al., 2020, Lemma 2).  $\square$

Now, we introduce an equivalent representation of the *do*-Shapley:

**Proposition S.2** ((Štrumbelj & Kononenko, 2014, Eq. (10))). *An equivalent representation of the *do*-Shapley in Eq. (2) is given as*

$$\tilde{\phi}_{V_i} := \frac{1}{n!} \sum_{\pi \in \Pi([n])} \{E[Y \text{do}(\mathbf{v}_{\text{pre}_\pi(i);i})] - E[Y \text{do}(\mathbf{v}_{\text{pre}_\pi(i)})]\},$$

where  $\Pi([n])$  is a set of all possible permutations of  $[n]$ ,  $\pi$  is an individual permutation in  $\Pi([n])$ ,  $\text{pre}_\pi(i) := \{k \in [n] \mid k < i \text{ in } \pi([n])\}$ .

This representation motivates a following Monte-Carlo-based approximation:

$$\tilde{\phi}_{V_i} := \frac{1}{M} \sum_{j=1}^M \{E[Y \text{do}(\mathbf{v}_{\text{pre}_{\pi_j}(i);i})] - E[Y \text{do}(\mathbf{v}_{\text{pre}_{\pi_j}(i)})]\}, \quad (\text{D.5})$$

where  $M$  is the number of randomly generated permutation of  $[n]$  and  $\pi_j$  denotes  $k$ th permutation. Convergence of  $\tilde{\phi}_{V_i}$  is guaranteed by the following result:

**Lemma S.5.**

$$\tilde{\phi}_{V_i} - \phi_{V_i} = O_P(M^{-1/2}). \quad (\text{D.6})$$

*Proof.* Let  $Z(\sigma) := E[Yjdo(\mathbf{v}_{i, \text{pre}_{\sigma(k)}(i)})] - E[Yjdo(\mathbf{v}_{\text{pre}_{\sigma(k)}(i)})]$  denote a random variable where the randomness is over the permutation  $\sigma$ , where  $P(\sigma) = \frac{1}{n!}$ . Then,  $E_P[Z(\sigma)] = \phi_{V_i}$ . By the given assumption,  $Z(\sigma)$  and

$$\tilde{\phi}_{V_i} := \frac{1}{M} \sum_{k=1}^M Z(\sigma(k))$$

are bounded random variables. Let  $B$  denote such bound. Then, by (Lattimore & Szepesvári, 2020, Corollary 5.5),

$$\tilde{\phi}_{V_i} > \phi_{V_i} - \sqrt{\frac{2B^2 \log(1/\delta)}{M}} \quad \text{and} \quad \tilde{\phi}_{V_i} < \phi_{V_i} + \sqrt{\frac{2B^2 \log(1/\delta)}{M}}$$

in probability  $(1 - \delta)$ , which implies that  $\tilde{\phi}_{V_i}$  converges in  $\frac{P}{M}$  rate. This completes the proof.  $\square$

Let  $S_{j;a} := \text{pre}_j(i)$  and  $S_{j;b} := \text{fig}[\text{pre}_j(i)]$ . By Def. 7, Eqs. (D.5, D.6),

$$\phi_{V_i}^{\text{est}} - \phi_{V_i} = \phi_{V_i}^{\text{est}} + \tilde{\phi}_{V_i} - \tilde{\phi}_{V_i} + \phi_{V_i} \quad (\text{D.7})$$

$$= \frac{1}{M} \sum_{i=1}^J (\{T^{\text{est}}(S_{j;b}) - E[Yjdo(\mathbf{v}_{S_{j;b}})]\} + \{T^{\text{est}}(S_{j;a}) - E[Yjdo(\mathbf{v}_{S_{j;a}})]\}) + O_P(M^{-1/2}). \quad (\text{D.8})$$

Now, we analyze each of IPW, REG, DML estimators in Defs. (4,5,6).

**Lemma S.6 (Error analysis for IPW).** For any nonempty  $S \subseteq [n]$ ,

$$T^{\text{IPW}}(S) - E[Yjdo(\mathbf{v}_S)] = \begin{cases} O_P(N^{-1/2}) + O_P(\|\hat{\omega}_S^S - \omega_S^S\|) & (\text{Markovian}) \\ O_P(N^{-1/2}) + O_P(\|\hat{\omega}_S^S - \omega_S^S\|) & (\text{Direct-cause}), \end{cases} \quad (\text{D.9})$$

*Proof.* We will prove only for the Markovian case, since the exactly same proof is applied for the Direct-cause case. First,  $E[Yjdo(\mathbf{v}_S)] = E[Y\omega_S^S]$ .

From Lemma S.4, it suffices to show that  $E[Y\hat{\omega}_S^S - Y\omega_S^S] = O_P(\|\hat{\omega}_S^S - \omega_S^S\|)$ . It can be shown by

$$E[Y\hat{\omega}_S^S - Y\omega_S^S] \leq k_Y k \|\hat{\omega}_S^S - \omega_S^S\| \cdot \|\hat{\omega}_S^S - \omega_S^S\|,$$

where the first inequality by Cauchy-Schwarz inequality and the second by the boundness of  $Y$ .  $\square$

**Lemma S.7 (Error analysis for REG).** For any nonempty  $S \subseteq [n]$ ,

$$T^{\text{REG}}(S) - E[Yjdo(\mathbf{v}_S)] = \begin{cases} O_P(N^{-1/2}) + O_P\left(\left\|\begin{matrix} \hat{\theta}_{0,1}^S & \theta_{0,1}^S \end{matrix}\right\|\right) & (\text{Markovian}) \\ O_P(N^{-1/2}) + O_P\left(\left\|\begin{matrix} \hat{\theta}_a^S & \theta_a^S \end{matrix}\right\|\right) & (\text{Direct-cause}). \end{cases}$$

*Proof.* We will prove only for the Markovian case, since the exactly same proof is applied for the Direct-cause case. We note that  $E[\theta_{0,1}^S] = E[Yjdo(\mathbf{v}_S)]$  by Lemma 2. From Lemma S.4, it suffices to show that  $E\left[\begin{matrix} \hat{\theta}_{0,1}^S & \theta_{0,1}^S \end{matrix}\right] = O_P\left(\left\|\begin{matrix} \hat{\theta} & \theta \end{matrix}\right\|\right)$ . It holds by Cauchy-Schwarz inequality.  $\square$

**Lemma S.8 (Error analysis for DML).** For any nonempty  $S \subseteq [n]$ ,

$$T^{\text{DML}}(S) - E[Yjdo(\mathbf{v}_S)] = \begin{cases} O_P(N^{-1/2}) + \sum_{j=1}^S O_P\left(\left\|\begin{matrix} \hat{\theta}_{j,2}^S & \theta_{j,2}^S \end{matrix}\right\| \left\|\begin{matrix} \hat{h}_j^S & h_j^S \end{matrix}\right\|\right) & (\text{Markovian}) \\ O_P(N^{-1/2}) + O_P\left(\left\|\begin{matrix} \hat{\theta}_a^S & \theta_a^S \end{matrix}\right\| \left\|\begin{matrix} \hat{h}_S & h_S \end{matrix}\right\|\right) & (\text{Direct-cause}). \end{cases}$$

*Proof.* We note that  $E[V(\mathbf{V}^\theta, \eta^S)] = E[Ydo(\mathbf{v}_S)]$  by Lemma 3. From Lemma S.4, it suffices to show that

$$E[V(\mathbf{V}^\theta, \hat{\eta}^S) \mid V(\mathbf{V}^\theta, \eta^S)] = \begin{cases} \sum_{j=1}^S O_P \left( \left\| \hat{\theta}_{j,2}^S \quad \theta_{j,2}^S \right\| \left\| \hat{h}_j^S \quad h_j^S \right\| \right) & \text{(Markovian)} \\ O_P \left( \left\| \hat{\theta}_a^S \quad \theta_a^S \right\| \left\| \hat{h}_S \quad h_S \right\| \right) & \text{(Direct-cause).} \end{cases}$$

First, consider the Markovian case. We omit the superscript  $S$ . Consider a following quantity: For  $j = 1, 2, \dots, s$ ,

$$Q_j := \theta_{j-1,1} + \sum_{k=j}^s \omega_{j:k}(\theta_{k,1} \quad \theta_{k,2}),$$

where  $\omega_{j:k} := \prod_{r=j}^k \frac{\mathbb{1}_{V_{m_r}}(V_{m_r})}{\hat{\pi}_r^S}$ . Let  $Q_{s+1} := Y$  and  $\omega_{j+1,1} = 0$ . We note that  $Q_1 = V(\mathbf{V}^\theta; \eta^S)$ , and  $E[Q_1] = E[Ydo(\mathbf{v}_S)]$ . Also, the following holds, by the definition of  $\theta_{k-1,1}, \theta_{k,2}$ :

$$\begin{aligned} E[\theta_{k-1,1}] &= E[\mathbb{1}_{V_{m_k}}(V_{m_k})\theta_{k,2}], \\ E[\hat{\theta}_{k-1,1}] &= E[\mathbb{1}_{V_{m_k}}(V_{m_k})\hat{\theta}_{k,2}]. \end{aligned}$$

First, we note that  $Q_j$  can be written in a recursion as follow: For  $j = 1, 2, \dots, s$ ,

$$Q_j = \theta_{j-1,1} + \omega_{j,j}(Q_{j+1} \quad \theta_{j,2}).$$

To witness, consider the followings:

$$\begin{aligned} Q_j &= \theta_{j-1,1} + \omega_{j,j}(\theta_{j,1} \quad \theta_{j,2}) + \omega_{j,j+1}(\theta_{j+1,1} \quad \theta_{j+1,2}) + \omega_{j,j+2}(\theta_{j+2,1} \quad \theta_{j+2,2}) + \\ Q_{j+1} &= \theta_{j,1} + \omega_{j+1,j+1}(\theta_{j+1,1} \quad \theta_{j+1,2}) + \omega_{j+1,j+2}(\theta_{j+2,1} \quad \theta_{j+2,2}) + \\ \omega_{j,j}Q_{j+1} &= \omega_{j,j}\theta_{j,1} + \omega_{j,j+1}(\theta_{j+1,1} \quad \theta_{j+1,2}) + \omega_{j,j+2}(\theta_{j+2,1} \quad \theta_{j+2,2}) + \dots \end{aligned}$$

Then,

$$\begin{aligned} Q_j &= \omega_{j,j}Q_{j+1} \quad \omega_{j,j}\theta_{j,1} + \theta_{j-1,1} + \omega_{j,j}(\theta_{j,1} \quad \theta_{j,2}) \\ &= \theta_{j-1,1} + \omega_{j,j}(Q_{j+1} \quad \theta_{j,2}). \end{aligned}$$

Finally, we will witness the following holds:

$$E[\hat{Q}_j \quad Q_j] = E[\hat{Q}_j \quad \theta_{j-1,1}] = \sum_{k=j}^s O_P \left( \left\| \theta_{k,2} \quad \hat{\theta}_{k,2} \right\| \left\| \hat{h}_k \quad h_k \right\| \right).$$

We will prove this by using an inductive hypothesis. First, at  $j = s$ ,

$$\begin{aligned} E[\hat{Q}_s \quad Q_s] &= E[\hat{Q}_s \quad \theta_{s-1,1}] = E[\hat{\theta}_{s-1,1} + \hat{\omega}_{s,s}(Y \quad \hat{\theta}_{s,2}) \quad \theta_{s-1,1}] \\ &= E\left[\hat{\theta}_{s-1,1} + \frac{\mathbb{1}_{V_{m_s}}(V_{m_s})}{\hat{\pi}_s}(Y \quad \hat{\theta}_{s,2}) \quad \theta_{s-1,1}\right] \\ &= E\left[\mathbb{1}_{V_{m_s}}(V_{m_s})(\hat{\theta}_{s,2} \quad \theta_{s,2}) + \frac{\mathbb{1}_{V_{m_s}}(V_{m_s})}{\hat{\pi}_s}(\theta_{s,2} \quad \hat{\theta}_{s,2})\right] \\ &= O_P \left( \left\| \theta_{s,2} \quad \hat{\theta}_{s,2} \right\| k\hat{\pi}_s \quad \pi_s k \right). \end{aligned}$$

For any  $j = s-1, \dots, 1$ ,

$$\begin{aligned}
 \mathbb{E} [\hat{Q}_j - Q_j] &= \mathbb{E} [\hat{Q}_j - \theta_{j-1,1}] = \mathbb{E} [\hat{\theta}_{j-1,1} - \theta_{j-1,1} + \hat{\omega}_{j:j} (\hat{Q}_{j+1} - \hat{\theta}_{j,2})] \\
 &= \mathbb{E} [\hat{\theta}_{j-1,1} - \theta_{j-1,1} + \hat{\omega}_{j:j} (\hat{Q}_{j+1} - \theta_{j,1}) + \hat{\omega}_{j:j} (\theta_{j,1} - \hat{\theta}_{j,2})] \\
 &= \mathbb{E} [\hat{\omega}_{j:j} (\hat{Q}_{j+1} - \theta_{j,1})] + \mathbb{E} [\mathbb{1}_{V_{m_j}}(V_{m_j}) (\hat{\theta}_{j,2} - \theta_{j,2}) + \hat{\omega}_{j:j} (\theta_{j,1} - \hat{\theta}_{j,2})] \\
 &= \mathbb{E} [\hat{\omega}_{j:j} (\hat{Q}_{j+1} - \theta_{j,1})] + \mathbb{E} \left[ \frac{1}{\hat{P}(V_{m_j} | \mathbf{W}_{m_j})} \{ \theta_{j,2} - \hat{\theta}_{j,2} \} \{ \hat{h}_j - h_j \} \right] \\
 &= \mathbb{E} [(\hat{Q}_{j+1} - \theta_{j,1})] + \mathbb{E} [\{ \theta_{j,2} - \hat{\theta}_{j,2} \} \{ \hat{h}_j - h_j \}] \\
 &= \mathbb{E} [(\hat{Q}_{j+1} - \theta_{j,1})] + \|\theta_{j,2} - \hat{\theta}_{j,2}\| \|\hat{h}_j - h_j\|.
 \end{aligned}$$

If we assume  $\mathbb{E} [\hat{Q}_r - \theta_{r-1,1}] = \sum_{k=r}^s O_P(\|\theta_{k,2} - \hat{\theta}_{k,2}\| \|\hat{h}_k - h_k\|)$  for  $r = j+1, \dots, s$ , then it's easy to witness that it holds for  $r = j$ , too. Therefore, by an induction, the equality holds for all  $r = 1, 2, \dots, 1$ . This completes the proof for Markovian case.

For Direct-cause case,

$$\begin{aligned}
 \mathbb{E} [V(\mathbf{V}^0; \eta^S) - V(\mathbf{V}; \hat{\eta}^S)] &= \mathbb{E} \left[ \frac{\mathbb{1}_{V_S}(\mathbf{V}_S)}{\hat{h}^S} (Y - \hat{\theta}_a) + \hat{\theta}_b - \theta_b \right] \\
 &= \mathbb{E} \left[ \frac{\mathbb{1}_{V_S}(\mathbf{V}_S)}{\hat{h}^S} (\theta_a - \hat{\theta}_a) + \hat{\theta}_b - \theta_b \right] \\
 &= \mathbb{E} \left[ \frac{h^S}{\hat{h}^S} (\theta_b - \hat{\theta}_b) + \hat{\theta}_b - \theta_b \right] \\
 &= \mathbb{E} \left[ \frac{1}{\hat{h}^S} (\theta_b - \hat{\theta}_b) (h^S - \hat{h}^S) \right] \\
 &= \mathbb{E} \left[ \frac{1}{\hat{h}^S} (\theta_b - \hat{\theta}_b) (h^S - \hat{h}^S) \right] \\
 &= \|\theta_b - \hat{\theta}_b\| \|h^S - \hat{h}^S\|.
 \end{aligned}$$

□

By combining Lemmas (S.4,S.5,S.6,S.7,S.8), we complete the proof of Theorem D.3. □

## E. Additional Experimental Details From Section 6

### E.1. Data Generating Processes

Here, we present the structural causal model for the data generating processes used for the data generating process used in Section 6.

We first note that  $U \sim \text{Bernoulli}(0.4)$ ,  $U_{V_1} \sim \text{Bernoulli}(0.8)$ ,  $U_{V_3} \sim \text{Bernoulli}(0.4)$ ,  $U_{V_2} \sim \text{Bernoulli}(0.3)$ , and  $U_Y \sim \text{Normal}(0, 1)$ . The SCM that induced the graph in Fig. 2a is

$$\begin{aligned}
 V_1 &= U_{V_1} - U \\
 V_3 &= U_{V_3} - U \\
 V_2 &= (V_1 \wedge V_3) - U_{V_2} \\
 Y &= 3V_1 + 0.4V_2 + V_3 + U_Y.
 \end{aligned}$$

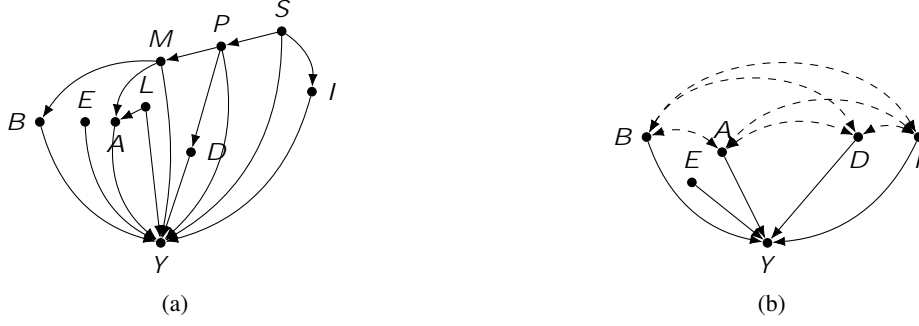


Figure F.4: **(a)** Causal graph for Example 1, taken from Lundberg (2021) **(b)** Variables  $fS, P, M, Lg$  are hidden. These graphs are used for Appendix F.

The SCM that induced the graph in Fig. 2b is

$$\begin{aligned} V_1 &= U_{V_1} \\ V_3 &= U_{V_3} \\ V_2 &= (V_1 \wedge V_3) \_ U_{V_2} \\ Y &= 3V_1 + 0.4V_2 + V_3 + U_Y. \end{aligned}$$

The SCM that induced the graph in Fig. 2c is

$$\begin{aligned} V_1 &= U_{V_1} \quad U \\ V_3 &= U_{V_3} \_ U \\ V_2 &= U_{V_2} \\ Y &= 3V_1 + 0.4V_2 + V_3 + U_Y. \end{aligned}$$

## F. Additional Experiments

In this section, we consider a different data generation process based on Example 1.

**Experimental Setup.** We use synthetic datasets based on: (a) Example 1 for which the corresponding causal graph Fig. F.4a is Markovian, and (b) the graph in Fig. F.4b which matches with Direct-cause case. These two graphs share the same data generating process since the graph in Fig. F.4b is generated from the graph in Fig. F.4a by omitting a set of variables. Details of the data generating process are provided in Appendix E. Throughout the simulation, we denote  $f\phi_{V_i} g_{i=1}^n$  as the ground-truth  $do$ -Shapley values.

**Comparison Between Estimators.** We compare the three estimators (IPW, REG, DML), denoted by  $f\phi_{V_i}^{ipw}, \phi_{V_i}^{reg}, \phi_{V_i}^{dml} g$  respectively, for scenarios depicted in graphs in Figs. (F.4a, F.4b). For all estimators, nuisances are estimated using gradient boosting model called XGBoost (Chen & Guestrin, 2016).

Let  $\phi_{V_i, k}^{est} \geq f\phi_{V_i, k}^{dml}, \phi_{V_i, k}^{ipw}, \phi_{V_i, k}^{reg} g$  denote an estimated importance of the  $i$ th feature of  $j$ th samples (i.e.,  $V_{i:k} \geq \mathbf{V}(k) \geq D$ ). As in Section 6, we assess the quality of the estimator by computing the  $L_1$  error as

$$L_1(\text{est}, k) := (1/n) \sum_{i=1}^n |\phi_{V_i, k}^{est} - \phi_{V_i, k}|,$$

(where  $n$  is the number of features). We ran the simulation for 100 randomly samples; i.e.,  $k \geq f1, 2, \dots, 100g$ , and with sample size  $N := jDj \geq f100, 1000, 5000, 10000g$  to observe convergence behaviors of estimators. We fix  $M = 100$ .

**Data Generating Processes.** Here, we present the structural causal model for the data generating processes used for the data generating process, where the qualitative graphical description is provided as causal graphs in Fig. F.4a. We will denote  $V_0$  : sales calls,  $V_1$  : interaction,  $V_2$  : economic factors,  $V_3$  : last upgrade,  $V_4$  : product needs,  $V_5$  : discounts provided,  $V_6$  : monthly usage,  $V_7$  : Ad spend,  $V_8$  : bugs reported,  $Y$  : customer retention (target variable).



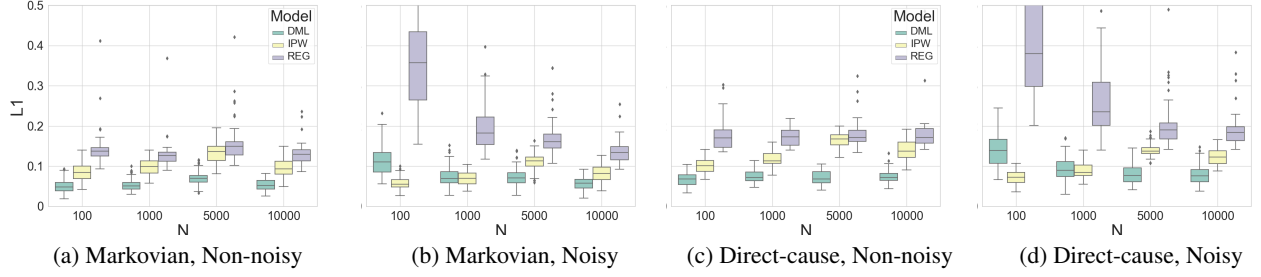


Figure F.5: The L1-error plots for the scenario in Section F.

$V_0 \sim P(U_{V_0})$ ,  $V_2 \sim P(U_{V_2})$ ,  $V_3 \sim P(U_{V_3})$  where  $P(U_{V_0})$  is a Uniform distribution ranging over 0 to 4;  $P(U_{V_2})$  is a Uniform distribution ranging over 0 to 1; and  $P(U_{V_3})$  is a Uniform distribution ranging over 0 to 20. For the rest of variables,

$$\begin{aligned} V_1 &= V_0 + U_{V_1} \\ V_4 &= 0.1 V_0 + U_{V_4} \\ V_5 &= 0.5(1 - \text{logit}(V_4)) + 0.5 U_{V_5} \\ V_6 &= \text{logit}(0.3 V_4 + U_{V_6}) \\ V_7 &= V_6 + U_{V_6} + \mathbb{1}_{V_3 < 1}(V_3) + \mathbb{1}_{V_3 < 2}(V_3) \\ V_8 &= U_{V_8}(V_6), \end{aligned}$$

where  $U_{V_1}$  is a Poisson random variable with the parameter 0.2,  $U_4 \sim \text{Normal}(0, 1)$ ,  $U_{V_5}$  is a Uniform variable ranging (0, 1),  $U_{V_6} \sim \text{normal}(0, 1)$ ,  $U_{V_6}$  is an Uniform variable ranging (0.9, 0.99), and  $U_{V_8}(V_6)$  is a Poisson random variable such that its parameter follows  $2V_6$ , and  $\mathbb{1}_{V_a < c}(V_a)$  is an indicator function for the variable  $V_a$  for the event  $V_a < c$  for some constant  $c$ . Finally,

$$Y^\theta = 0.9V_4^\theta + 0.8V_6^\theta - 0.2V_2^\theta + 0.05V_5^\theta - 0.015(1 - V_8^\theta) + 0.2V_0^\theta + 0.3V_1 + 0.5(V_3 + 0.25) + 0.6V_7 - U_Y - 0.45,$$

where  $fV_4^\theta, V_6^\theta, V_2^\theta, V_5^\theta, V_8^\theta, V_0^\theta g$  are random variables from the normal distribution where the variance is 1 and their means are  $fV_4, V_6, V_2, V_5, V_8, V_0 g$ . Finally,  $Y = \text{logit}(7Y)$ .

For the Case 2, we drop the variable  $V_0, V_4, V_6$ , and we used

$$Y^\theta = -0.2V_2^\theta + 0.05V_5^\theta - 0.015(1 - V_8^\theta) + 0.3V_1 + 0.5/(V_3/4 + 0.25) + 0.6V_7 - U_Y - 0.45.$$

We also recommend checking the code `data_generator_1.py`, `data_generator_2.py` for the detailed configurations of the data generating processes.

**Experimental Results.** For the non-noisy setting, the L1-error plots for *f*Markovian, Direct-cause*g* cases are presented in Figs. (F.5a, F.5c) respectively. The DML-based estimator  $f\hat{\phi}_{V_i}^{\text{dml}} g_{I=1}^\theta$  outperforms  $(f\hat{\phi}_{V_i}^{\text{ipw}}, \hat{\phi}_{V_i}^{\text{reg}} g_{I=1}^\theta)$  for all  $N \geq f100, 1000, 5000, 10000 g$ , and it achieves the smallest variance compared to other estimators. This result corroborates with the robustness property of the DML-based estimator (see Remark 4). The L1-error plots for the noisy setting for *f*Markovian, Direct-cause*g* cases are presented in Figs. (F.5b, F.5d) respectively. In this case, the DML-based estimator  $f\hat{\phi}_{V_i}^{\text{dml}} g_{I=1}^\theta$  exhibits the debiasedness property against the converging noise, while other estimators converge much slower.

**Contrasting with the ICC Approach (Janzing et al., 2020a).** We contrast the *do*-Shapley with the ICC approach (Janzing et al., 2020a). The *do*-Shapley measures the feature importance based on the total effect of variables, while the ICC measures based on their *intrinsic* effects. It is not possible to quantitatively compare these two contrasting definitions.

We compute the feature importance as proposed in (Molnar, 2020, Chap. 9.6.5), where the importance of the  $j$ th feature is defined as:

$$I_j := \frac{1}{jD^j} \sum_{\mathbf{v}_{(j)} \in 2^{D^j}} |\phi_j(\mathbf{v}_{(j)})|,$$

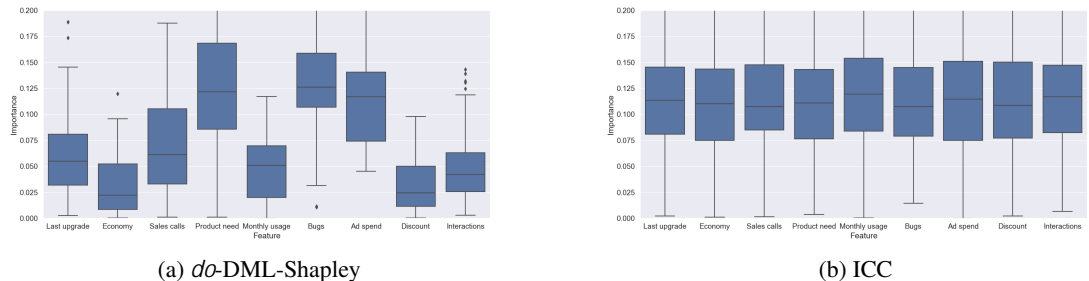


Figure F.6: Feature importance plots for the *do*-DML-Shapley and the ICC approaches.

	S	P	I	M	D	L	E	A	B
<i>do</i> -DML	0.07	0.12	0.05	0.05	0.03	0.06	0.03	0.12	0.13
ICC	0.12	0.13	0.12	0.13	0.13	0.12	0.12	0.12	0.12

Table 3: Average of feature importances produced by the *do*-DML-Shapley and ICC approaches.

where  $\phi_i$  is the Shapley value, and  $D^\emptyset \subseteq D$  is a subset of samples.

In our experiments, we randomly selected 100 samples and compare the feature importance using the *do*-DML-Shapley ( $\phi_{V_i}^{dml}$ ) and the ICC approach, denoted  $\phi_{V_i}^{icc}$ . The average of the estimated importance of each features described in Example 1 is presented in Table 3. In Fig. F.6, we present the bar-plot for both the *do*-DML-Shapley and the ICC approaches using the observations  $f|\phi_i^{dml}(\mathbf{V}_{(j)})|g_{\mathbf{V}_{(j)} \setminus D^\emptyset}$  and  $f|\phi_i^{icc}(\mathbf{V}_{(j)})|g_{\mathbf{V}_{(j)} \setminus D^\emptyset}$ .

In our experiments, the *do*-Shapley approach gives that the production needs (*P*) has the largest total effect where *P* is in fact the variable with largest coefficient (0.9). in our data generating process, whereas ICC approach gives that all variables have similar intrinsic effects.