
Causal Shapley values

Anonymous Author(s)

Affiliation

Address

email

Abstract

Shapley values have become one of the most popular model-agnostic methods within explainable artificial intelligence. They explain the output of a possibly complex and highly non-linear machine learning model by attributing the difference between the model output and an average, baseline output to the different features that are used as input to the model. Being based on game-theoretic principles, Shapley values uniquely satisfy several desirable properties. Shapley values are typically computed under the assumption that features are independent, which ignores any causal structure between the features and can lead to illogical explanations.

In this paper, we propose a novel framework for causal Shapley values that encompasses and sheds new light on recent previous work that aims to relax the independence assumption. Our framework is based on an active, causal interpretation of feature attribution. Replacing conventional conditioning by observation with conditioning by intervention, we call upon Pearl’s *do*-calculus to compute the Shapley values, without giving in on any of their desirable properties. We provide a practical implementation for deriving Shapley values based on causal chain graphs and illustrate causal Shapley values on several real-world examples.

Not completely happy yet. Suggestions welcome!

1 Introduction

Explainability. Attribution. Shapley. Conditioning or no conditioning.

2 A causal interpretation of Shapley values

Need more smooth talk.

We assume that we are given a set of features \mathbf{x} with corresponding model output $f(\mathbf{x})$. We compare this output with the average output

$$f_0 = \mathbb{E}f(\mathbf{X}) = \int d\mathbf{X} P(\mathbf{X})f(\mathbf{X}) ,$$

with expectation taken over some (for now assumed to be known) probability distribution $P(\mathbf{X})$. We would like to attribute the difference between $f(\mathbf{x})$ and f_0 in a sensible way to the different features $i \in N$ with $N = \{1, \dots, n\}$ and n the number of features. That is, we would like to write

$$f(\mathbf{x}) = f_0 + \sum_{i=1}^n \phi_i , \tag{1}$$

where we will refer to ϕ_i as the contribution of feature i to the output $f(\mathbf{x})$. Equation (1) is referred to as the efficiency property [9], which appears to be a sensible desideratum for any attribution method and we therefore take here as our starting point.

30 We can think of (at least) two different interpretations on how to go from knowing none of the feature
31 values in f_0 to knowing all feature values in $f(\mathbf{x})$.

32 **Passive.** We interpret the feature vector \mathbf{x} as a passive observation. Feature values come in one after
33 the other and the contribution of feature i should reflect the difference in expected value of
34 $f(\mathbf{X})$ after and before *observing* its feature value x_i .

35 **Active.** We interpret the feature vector \mathbf{x} as the result of an action. Feature values are imposed one
36 after the other and the contribution of feature i relates to the difference in expected value of
37 $f(\mathbf{X})$ after and before *setting* its value to x_i .

38 Following the above reasoning, the contribution of each feature depends on the order π in which the
39 feature values arrive or are imposed. We write the contribution of feature i given the permutation π
40 as

$$\phi_i(\pi) = v(\{j : j \preceq_\pi i\}) - v(\{j : j \prec_\pi i\}), \quad (2)$$

41 with $j \prec_\pi i$ if j precedes i in the permutation π and where we define the value function

$$v(S) = \mathbb{E}[f(\mathbf{X}) | op(\mathbf{x}_S)] = \int d\mathbf{X}_{\bar{S}} P(\mathbf{X}_{\bar{S}} | op(\mathbf{X}_S = \mathbf{x}_S)) f(\mathbf{X}_{\bar{S}}, \mathbf{x}_S). \quad (3)$$

42 Here S is the subset of indices of features with known ‘in-coalition’ feature values \mathbf{x}_S . To compute the
43 expectation, we still need to average over the ‘out-of-coalition’ feature values $\mathbf{X}_{\bar{S}}$ with $\bar{S} = N \setminus S$, the
44 complement of S . The operator $op()$ specifies how the distribution of the ‘out-of-coalition’ features
45 $\mathbf{X}_{\bar{S}}$ depends on the ‘in-coalition’ feature values \mathbf{x}_S . To arrive at the passive interpretation, we set
46 $op()$ to conventional conditioning by observation, yielding $P(\mathbf{X}_{\bar{S}} | \mathbf{x}_S)$. For the active interpretation,
47 we need to condition by intervention, for which we resort to Pearl’s *do*-calculus [6] and write
48 $P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S))$. **Do we need a longer introduction/explanation of *do*-calculus?** A third option
49 is to ignore the feature values \mathbf{x}_S and just take the unconditional, marginal distribution $P(\mathbf{X}_{\bar{S}})$. We
50 will refer to the corresponding Shapley values as conditional, causal, and marginal, respectively.

51 It is easy to check that the efficiency property (1) holds for any permutation π . So, for any distribution
52 over permutations $w(\pi)$ with $\sum_\pi w(\pi) = 1$, the contributions

$$\phi_i = \sum_\pi w(\pi) \phi_i(\pi)$$

53 still satisfy (1). An obvious choice would be to take a uniform distribution $w(\pi) = 1/n!$. We then
54 arrive at the standard definition of Shapley values:

$$\phi_i = \sum_{S \subseteq N \setminus i} \frac{|S|!(n - |S| - 1)!}{n!} [v(S \cup i) - v(S)],$$

55 where we use shorthand i for the singleton $\{i\}$. Besides efficiency, these Shapley values uniquely
56 satisfy three other desirable properties [9].

57 **Linearity:** for two value functions v_1 and v_2 , we have $\phi_i(\alpha_1 v_1 + \alpha_2 v_2) = \alpha_1 \phi_i(v_1) + \alpha_2 \phi_i(v_2)$.
58 This guarantees that the Shapley value of a linear ensemble of models is a linear combination
59 of the Shapley values of the individual models.

60 **Null player (dummy):** if $v(S \cup i) = v(S)$ for all $S \subseteq N \setminus i$, then $\phi_i = 0$. A feature that never
61 contributes to the model output receives zero Shapley value.

62 **Symmetry:** if $v(S \cup i) = v(S \cup j)$ for all $S \subseteq N \setminus \{i, j\}$, then $\phi_i = \phi_j$. Symmetry holds for
63 marginal, conditional, and causal Shapley values.

64 Efficiency, linearity, and null player still hold for a non-uniform distribution of permutations, but
65 symmetry is then typically lost.

66 **Part of the next two paragraphs probably to introduction.**

67 Janzing et al. [4] also argued for an active, interventional interpretation of Shapley values. They make
68 a case for using the marginal instead of the (observational) conditional distribution to compute the
69 Shapley values when dependencies between features are due to confounding. This follows directly
70 from our reasoning, since in models with no causal links between the features and any dependen-
71 cies only due to confounding, conditioning by intervention reduces to the marginal distribution:
72 $P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) = P(\mathbf{X}_{\bar{S}})$ for any subset S .



Figure 1: Two causal models. In both, X_1 causes X_2 and X_3 . In Model A the excess correlation between X_2 and X_3 is induced by a common confounder Z , in Model B by selection bias.

Frye et al. [3] introduce asymmetric Shapley values as a way to incorporate causal information. Instead of taking a uniform distribution over all possible permutations, these asymmetric Shapley values only consider those permutations that are consistent with the causal structure between the features, i.e., are such that a known causal ancestor always precedes its descendants. Frye et al. apply conventional conditioning by observation to make sure that the resulting explanations respect the data manifold. This makes the approach somewhat of a mix between an active (incorporating causal structure) and a passive approach (conditioning by observation). In this paper, we suggest a more direct approach to incorporate causality, by replacing conditioning by observation with conditioning by intervention. This can, but does not have to be combined with the idea to make the permutations match the causal structure. We will illustrate the differences in the example below.

3 Illustration

For illustration, we consider two causal models in Figure 1. They have a different causal structure, but the same dependency structure (all features are dependent) and we assume that the probability distribution $P(\mathbf{X})$ is exactly the same for Model A and Model B. Our estimate of the output Y is a linear function of the features:

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^3 \beta_i x_i.$$

Too much detail in this section: move parts to supplement? If so, which parts? Combining (2) and (3), we obtain, after some rewriting

$$\phi_i(\pi) = \beta_i (x_i - \mathbb{E}[X_i | op(\mathbf{x}_{j:j \prec_{\pi} i})]) + \sum_{k \succ_{\pi} i} \beta_k (\mathbb{E}[X_k | op(\mathbf{x}_{j:j \preceq_{\pi} i})] - \mathbb{E}[X_k | op(\mathbf{x}_{j:j \prec_{\pi} i})]).$$

For marginal Shapley values only the first term before the sum remains, yielding

$$\phi_i = \phi_i(\pi) = \beta_i (x_i - \mathbb{E}[X_i]),$$

as also derived in [1].

Analytically computing the conditional Shapley values is tedious, but conceptually straightforward. To write the equations in a compact form, we define $\bar{x}_{k|S} = \mathbb{E}[X_k | \mathbf{x}_S]$ and combine all expectations in a single matrix $\bar{\mathcal{X}}$:

$$\bar{\mathcal{X}} = \begin{pmatrix} x_1 & x_2 & x_3 \\ \bar{x}_1 & \bar{x}_2 & \bar{x}_3 \\ \bar{x}_{1|2} & \bar{x}_{2|1} & \bar{x}_{3|1} \\ \bar{x}_{1|3} & \bar{x}_{2|3} & \bar{x}_{3|2} \\ \bar{x}_{1|2,3} & \bar{x}_{2|1,3} & \bar{x}_{3|1,2} \end{pmatrix}.$$

Any vector ϕ with the Shapley values of the three features can be written as a linear combination of these expectations times the regression coefficients β . With definition of the matrix

$$\mathcal{C}^{\text{conditional}} = \frac{1}{6} \left(\begin{array}{cccc|cccccccc} 6 & -2 & -1 & -1 & -2 & 0 & -2 & 2 & -1 & 1 & 0 & -2 & 2 & -1 & 1 \\ 0 & -2 & 2 & -1 & 1 & 6 & -2 & -1 & -1 & -2 & 0 & -2 & -1 & 2 & 1 \\ 0 & -2 & -1 & 2 & 1 & 0 & -2 & -1 & 2 & 1 & 6 & -2 & -1 & -1 & -2 \end{array} \right),$$

expectation	model A	model B
$\hat{x}_{1 2}$	\bar{x}_1	
$\hat{x}_{1 3}$	\bar{x}_1	
$\hat{x}_{1 2,3}$	\bar{x}_1	
$\hat{x}_{2 1}$	$\bar{x}_{2 1}$	
$\hat{x}_{2 3}$	\bar{x}_2	$\bar{x}_{2 3}$
$\hat{x}_{2 1,3}$	$\bar{x}_{2 1}$	$\bar{x}_{2 1,3}$
$\hat{x}_{3 1}$	$\bar{x}_{3 1}$	
$\hat{x}_{3 2}$	\bar{x}_3	$\bar{x}_{3 2}$
$\hat{x}_{3 1,2}$	$\bar{x}_{3 1}$	$\bar{x}_{3 1,2}$

Table 1: Turning expectations under conditioning by intervention, $\hat{x}_{i|S} = \mathbb{E}[x_i | do(\mathbf{X}_S = \mathbf{x}_S)]$, into expectations under conventional conditioning by observation, $\bar{x}_{i|S} = \mathbb{E}[x_i | \mathbf{X}_S]$, for the two models in Figure 1. **Needed? Can put it next to Figure 1 to save space.**

the conditional Shapley values for any linear model with three variables and a uniform distribution over permutations can be written as

$$\phi^{\text{conditional}} = C^{\text{conditional}} \text{vec}(\text{diag}(\beta) \bar{\mathcal{X}}). \quad (4)$$

Skip this explanation? Vectorization stacks the columns on top of one another to end up with a 15-dimensional column vector. The vertical bars in the matrix $C^{\text{conditional}}$ indicate the three blocks, with the first 5 columns in the matrix mapping to the first column of $\bar{\mathcal{X}}$ with expectations of X_1 , the next 5 columns to the expectations of X_2 , and the final 5 columns to the expectations of X_3 .

Skip this paragraph? By summing every column of $C^{\text{conditional}}$, we can perform the sanity check that efficiency indeed holds. The first column of each block (which relates to the feature values themselves) adds up to 1, the second (corresponding to the marginal expectations) to -1 , and the other three columns to zero, as they should. Since Shapley values are constructed by always comparing two (possibly) different expectations, each row within each block sums up to zero.

Putting X_1 before X_2 and X_3 , and X_2 and X_3 on equal footing, asymmetric Shapley values only consider the two permutations where x_1 is observed before x_2 and x_3 , leading to (we divide by 6 to make it easier to compare with the other Shapley values)

$$C^{\text{asymmetric}} = \frac{1}{6} \left(\begin{array}{ccccc|ccccc|ccccc} 6 & -6 & 0 & 0 & 0 & 0 & -6 & 6 & 0 & 0 & 0 & 0 & -6 & 6 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 6 & 0 & -3 & 0 & -3 & 0 & 0 & 0 & 0 & -3 & 0 & 3 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -3 & 0 & 3 & 6 & 0 & -3 & 0 & 0 & -3 & 0 \end{array} \right).$$

Using the standard rules for *do*-calculus [6], we show in Table 1 how the expectations under conditioning by intervention reduce to expectations under conditioning by observation. Since in Model A the correlation between X_2 and X_3 is due to confounding, the interventional expectations and thus Shapley values simplify considerably:

$$C^{\text{causal,A}} = \frac{1}{6} \left(\begin{array}{ccccc|ccccc|ccccc} 6 & -6 & 0 & 0 & 0 & 0 & -3 & 3 & 0 & 0 & 0 & 0 & 0 & -3 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 6 & -3 & -3 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 6 & -3 & 3 & 0 & 0 & 0 & 0 \end{array} \right).$$

For Model B, on the other hand, features X_2 and X_3 do affect each other when intervened upon, which makes that compared to the conditional Shapley values only the first block changes:

$$C^{\text{causal,B}} = \frac{1}{6} \left(\begin{array}{ccccc|ccccc|ccccc} 6 & -6 & 0 & 0 & 0 & 0 & -2 & 2 & -1 & 1 & 0 & 0 & -2 & 2 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 6 & -2 & -1 & -1 & -2 & 0 & 0 & -2 & -1 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -2 & -1 & 2 & 1 & 6 & 0 & -2 & -1 & -1 & -2 & 0 \end{array} \right).$$

To make this more concrete, let us assume that $P(\mathbf{X})$ follows a multivariate normal distribution corresponding to the causal model

$$X_1 \sim \mathcal{N}(0; 1) \text{ and } (X_2, X_3 | X_1) \sim \mathcal{N}\left((\alpha_2 X_1, \alpha_3 X_1); \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right),$$

where for notational convenience, we chose zero means and unit variance for all noise variables. The first feature drives the second and third feature with coefficients α_2 and α_3 . The confounding in Model A or selection bias in Model B leads to correlation ρ on top of the correlation induced by the joint dependence on the first feature. Straightforward calculations yield

$$\bar{\mathcal{X}} = \begin{pmatrix} x_1 & x_2 & x_3 \\ 0 & 0 & 0 \\ \frac{\alpha_2}{\sigma_2^2} x_2 & \alpha_2 x_1 & \alpha_3 x_1 \\ \frac{\alpha_3}{\sigma_3^2} x_3 & \frac{\gamma}{\sigma_3^2} x_3 & \frac{\gamma}{\sigma_2^2} x_2 \\ \frac{\delta_{23} x_2 + \delta_{32} x_3}{1 - \rho^2 + \delta_{23} \alpha_2 + \delta_{32} \alpha_3} & (\alpha_2 - \rho \alpha_3) x_1 + \rho x_3 & (\alpha_3 - \rho \alpha_2) x_1 + \rho x_2 \end{pmatrix}.$$

with $\delta_{ij} = \alpha_i - \rho \alpha_j$, $\sigma_i^2 = 1 + \alpha_i^2$ (the marginal variance for feature i), and $\gamma = \rho + \alpha_2 \alpha_3$ (the total covariance between the second and third feature). Plugging this and the expressions for the different \mathcal{C} matrices into (4), we obtain the asymmetric Shapley values

$$\begin{aligned} \phi_1^{\text{asymmetric}} &= \beta_1 x_1 + (\alpha_2 \beta_2 + \alpha_3 \beta_3) x_1 \\ \phi_2^{\text{asymmetric}} &= \beta_2 x_2 - \alpha_2 \beta_2 x_1 + \frac{\rho}{2} (\alpha_3 \beta_2 - \alpha_2 \beta_3) x_1 + \frac{\rho}{2} (\beta_3 x_2 - \beta_2 x_3) \\ \phi_3^{\text{asymmetric}} &= \beta_3 x_3 - \alpha_3 \beta_3 x_1 + \frac{\rho}{2} (\alpha_2 \beta_3 - \alpha_3 \beta_2) x_1 + \frac{\rho}{2} (\beta_2 x_3 - \beta_3 x_2), \end{aligned}$$

and the causal Shapley values, for Model A,

$$\begin{aligned} \phi_1^{\text{causal,A}} &= \beta_1 x_1 + \frac{1}{2} (\alpha_2 \beta_2 + \alpha_3 \beta_3) x_1 \\ \phi_2^{\text{causal,A}} &= \beta_2 x_2 - \frac{1}{2} \alpha_2 \beta_2 x_1 \\ \phi_3^{\text{causal,A}} &= \beta_3 x_3 - \frac{1}{2} \alpha_3 \beta_3 x_1. \end{aligned}$$

and, for Model B, (really ugly... can we prevent this?)

$$\begin{aligned} \phi_1^{\text{causal,B}} &= \beta_1 x_1 + \frac{1}{2} (\alpha_2 \beta_2 + \alpha_3 \beta_3) x_1 - \frac{\rho}{6} (\alpha_3 \beta_2 + \alpha_2 \beta_3) x_1 - \frac{1}{6} \left(\frac{\alpha_3 \delta_{23}}{\sigma_3^2} \beta_2 x_3 + \frac{\alpha_2 \delta_{32}}{\sigma_2^2} \beta_3 x_2 \right) \\ \phi_2^{\text{causal,B}} &= \beta_2 x_2 - \frac{1}{2} \alpha_2 \beta_2 x_1 + \frac{\rho}{6} (2 \alpha_3 \beta_2 - \alpha_2 \beta_3) x_1 + \\ &\quad \frac{1}{6} \left(2 \frac{\alpha_2 \delta_{32}}{\sigma_2^2} \beta_3 x_2 - \frac{\alpha_3 \delta_{23}}{\sigma_3^2} \beta_2 x_3 \right) + \frac{\rho}{2} (\beta_3 x_2 - \beta_2 x_3) \\ \phi_3^{\text{causal,B}} &= \beta_3 x_3 - \frac{1}{2} \alpha_3 \beta_3 x_1 + \frac{\rho}{6} (2 \alpha_2 \beta_3 - \alpha_3 \beta_2) x_1 + \\ &\quad \frac{1}{6} \left(2 \frac{\alpha_3 \delta_{23}}{\sigma_3^2} \beta_2 x_3 - \frac{\alpha_2 \delta_{32}}{\sigma_2^2} \beta_3 x_2 \right) + \frac{\rho}{2} (\beta_2 x_3 - \beta_3 x_2). \end{aligned}$$

In this linear model, the asymmetric Shapley value for the first feature adds its indirect causal effects on the output through the second and third feature, $\alpha_2 \beta_2 x_1 + \alpha_3 \beta_3 x_1$, to its direct effect, $\beta_1 x_1$. The causal Shapley values for the first feature are somewhat more conservative: they essentially claim only half of the indirect effects through the other two features. **Move the rest of this paragraph elsewhere? Now overlap with direct/indirect in next section.** This is a direct consequence of taking a uniform distribution over all permutations: for any pair of features i and j , the feature value x_i is set before and after x_j for exactly half of the number of permutations. Which distribution over permutations to prefer, a uniform one or one that respects the causal structure, depends on the question the practitioner tries to answer and possibly on the application. For example, when a causal link represents a temporal relationship, it may make no sense to set a feature value before the values of all features preceding it in time have been set. In that case, it would be wise to consider a non-uniform distribution over permutations as in [3]. On the other hand, for causal models without temporal interpretation, e.g., describing presumed causal relationships between personal and biomedical variables related to Alzheimer [10] or between social and economic characteristics in census data [2], deviating from a uniform distribution over permutations (and hence sacrificing the symmetry property) seems unnecessary. With or without uniform distribution over permutations, applying *do*-calculus instead of conditioning by observation is a natural way to incorporate causal information.

The Shapley values for Model A are different from those for Model B, even though the observable probability distribution $P(\mathbf{X})$ is exactly the same. Those for Model A simplify a lot, because in this model any excess correlation between X_2 and X_3 beyond the correlation resulting from the common parent X_1 results from a confounder. This correlation vanishes when we intervene on either X_2 or X_3 . The contributions of the second and third feature are therefore just their direct effect, minus half of the indirect effect, which already has been attributed to the first feature. In Model B, on the other hand, for most expectations conditioning by intervention reduces to conditioning by observation on the same variables and does not further simplify to conditioning on less or even no variables as for Model A. Since the causal Shapley values consider all six permutations, in contrast to the asymmetric Shapley values which only consider two of them, expectations such as $\mathbb{E}[X_2|x_3]$ now also enter the equation, which considerably complicates the analytical expressions.

4 Causal chain graphs

Added a theorem, corollaries, a figure, and an algorithm. Still need to decide which ones to keep and then adapt the text accordingly by adding the appropriate references and removing repetitions.

Computing causal Shapley values not only requires knowledge of the probability distribution $P(\mathbf{X})$, but also of the underlying causal structure. And even then, there is no guarantee that any causal query is identifiable (see e.g., [7]). For example, if Model A or B also includes a causal link from X_2 to X_3 , even knowing the probability distribution $P(\mathbf{X})$ and the causal structure is insufficient: it is impossible to express, for example, $P(X_3|do(X_2 = x_2))$ in terms of $P(\mathbf{X})$, essentially because, without knowing the parameters of the causal model, there is no way to tell which part of the observed dependence between X_2 and X_3 is due to the causal link and which due to the confounding or selection bias.

Furthermore, and perhaps more importantly, requiring a practitioner to specify a complete causal structure, possibly even including some of its parameters, would be detrimental to the method's general applicability. We therefore follow the same line of reasoning as in [3] and assume that a practitioner may be able to specify a causal ordering, but not much more.

In the special case that a complete causal ordering of the features can be given and that all causal relationships are unconfounded, $P(\mathbf{X})$ satisfies the Markov properties associated with a directed acyclic graph (DAG) and can be written in the form

$$P(\mathbf{X}) = \prod_{j \in N} P(X_j | \mathbf{X}_{pa(j)}) ,$$

with $pa(j)$ the parents of node j . If no further conditional independences are assumed, the parents of j are all nodes that precede j in the causal ordering. For causal DAGs, we have the interventional formula [5]:

$$P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) = \prod_{j \in \bar{S}} P(X_j | \mathbf{X}_{pa(j) \cap \bar{S}}, \mathbf{x}_{pa(j) \cap S}) , \quad (5)$$

with $pa(j) \cap T$ the parents of j that are also part of subset T . The interventional formula can be used to answer any causal query of interest. We will often approximate the expectations needed to compute the Shapley values through sampling, which is particularly straightforward for causal DAGs under conditioning by intervention. Variables are sampled consecutively by following the causal ordering. The probability distribution for a feature then only depends on the values of its parents, which by then is either sampled or fixed. Since the intervention blocks the influence of all descendants, there is no need for an MCMC approach such as Gibbs sampling: the values of all features can be sampled in a single pass through the graph.

We may not always be willing or able to give a complete ordering between the individual variables, but rather a partial ordering as, for example, in Figure 1 where we have the partial ordering $(\{1\}, \{2, 3\})$: the first feature precedes the second and third feature in the causal ordering, with the second and third feature on equal footing, i.e., without specifying whether the second causes the third or vice versa. Here causal chain graphs [5] come to the rescue. A causal chain graph has directed and undirected edges. All features that are treated on an equal footing are linked together with undirected edges and become part of the same chain component. Edges between chain components are directed and represent causal relationships. The probability distribution $P(\mathbf{X})$ now factorizes as a “DAG of chain

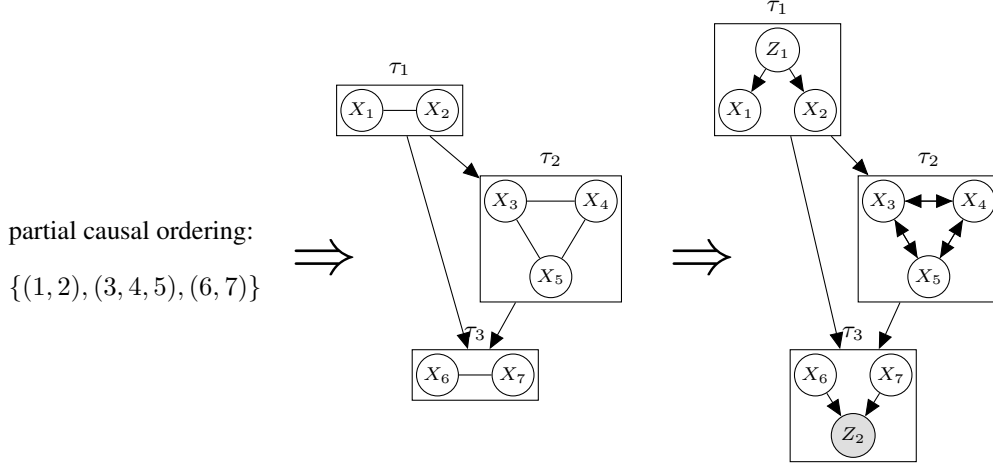


Figure 2: From partial ordering to causal chain graph. Features on equal footing are combined into a fully connected chain component. How to handle interventions within each component depends on the generative process that best explains the (surplus) dependencies. In this example, the dependency between X_1 and X_2 in chain component τ_1 is assumed to be the result of a common confounder. The surplus dependencies in τ_2 and τ_3 are assumed to be caused by mutual feedback and selection bias, respectively. **Attempt to illustrate the main ideas. Could be nice, but probably not enough space?**

193 components”:

$$P(\mathbf{X}) = \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_\tau | \mathbf{X}_{pa(\tau)}) ,$$

194 with each τ corresponding to a chain component, consisting of all features that are treated on an
 195 equal footing.

196 How to compute the effect of an intervention now depends on the interpretation of the generative
 197 process leading to the (surplus) dependencies between features within each component. If we assume
 198 that these are the consequence of marginalizing out a common confounder, as in Model A in Figure 1,
 199 intervention on a particular feature will break the dependency with the other features. We will refer
 200 to the set of chain components for which this applies as $\mathcal{T}_{\text{confounding}}$. Another possible interpretation is
 201 that the undirected part corresponds to the equilibrium distribution of a dynamic process resulting
 202 from interactions between the variables within a component [5]. In this case, setting the value of a
 203 feature does affect the distribution of the variables within the same component. The same applies to
 204 the case of selection bias, as in Model B in Figure 1.

205 **Theorem 1.** *For causal chain graphs, we have the interventional formula*

$$P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) = \prod_{\tau \in \mathcal{T}_{\text{confounding}}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}) \times \prod_{\tau \in \overline{\mathcal{T}_{\text{confounding}}}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}, \mathbf{x}_{\tau \cap S}) .$$

Proof.

$$\begin{aligned} P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) &\stackrel{(1)}{=} \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, do(\mathbf{X}_S = \mathbf{x}_S)) \\ &\stackrel{(3)}{=} \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, do(\mathbf{X}_{pa(\tau) \cap S} = \mathbf{x}_{pa(\tau) \cap S}), do(\mathbf{X}_{\tau \cap S} = \mathbf{x}_{\tau \cap S})) \\ &\stackrel{(2)}{=} \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}, do(\mathbf{X}_{\tau \cap S} = \mathbf{x}_{\tau \cap S})) , \end{aligned}$$

206 where the number above each equal sign refers to the standard *do*-calculus rule from [7] that is
 207 applied. For a chain component with dependencies induced by a common confounder, rule (3) applies

208 once more and yields

$$P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}),$$

209 whereas for a chain component with dependencies induced by selection bias or mutual interactions,
210 rule (2) again applies:

$$P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}, \mathbf{x}_{\tau \cap S}).$$

211 **Proof probably needs to be extended, which is fine when it goes to the supplement anyway.** \square

Algorithm 1 Compute the value function $v(S)$ under conditioning by intervention. **Include in main text? Seems rather obvious...**

```

1: function VALUEFUNCTION( $S$ )
2:   for  $k \leftarrow 1$  to  $N_{\text{samples}}$  do
3:     for all  $j \leftarrow 1$  to  $|\mathcal{T}|$  do ▷ run over all chain components in their causal order
4:       if confounding( $\tau_j$ ) then
5:         for all  $i \in \tau_j \cap \bar{S}$  do
6:           Sample  $\tilde{x}_i^{(k)} \sim P(X_i | \tilde{\mathbf{x}}_{pa(\tau_j) \cap \bar{S}}^{(k)}, \mathbf{x}_{pa(\tau_j) \cap S})$  ▷ can be drawn independently
7:         end for
8:       else
9:         Sample  $\tilde{\mathbf{x}}_{\tau_j \cap \bar{S}}^{(k)} \sim P(\mathbf{X}_{\tau_j \cap \bar{S}} | \tilde{\mathbf{x}}_{pa(\tau_j) \cap \bar{S}}^{(k)}, \mathbf{x}_{pa(\tau_j) \cap S}, \mathbf{x}_{\tau_j \cap S})$  ▷ e.g., Gibbs sampling
10:      end if
11:    end for
12:  end for
13:   $v \leftarrow \frac{1}{N_{\text{samples}}} \sum_{k=1}^{N_{\text{samples}}} f(\mathbf{x}_S, \tilde{\mathbf{x}}_{\bar{S}}^{(k)})$ 
14:  return  $v$ 
15: end function

```

212 Theorem 1 connects to observations made and algorithms proposed in recent papers.

213 **Corollary 2.** *With all features combined in a single component and all dependencies induced by*
214 *confounding, as in [4], causal Shapley values are equivalent to marginal Shapley values.*

215 *Proof.* We immediately get $P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) = P(\mathbf{X}_{\bar{S}})$ for all subsets S , i.e., as if all features
216 are independent. \square

217 **Corollary 3.** *With all features combined in a single component and all dependencies induced by*
218 *selection bias or mutual interactions, causal Shapley values are equivalent to conditional Shapley*
219 *values as proposed in [1].*

220 *Proof.* Now $P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)) = P(\mathbf{X}_{\bar{S}} | \mathbf{x}_S)$ for all subsets S , which boils down to conven-
221 tional conditioning by observation. \square

222 **Corollary 4.** *When we only take into account permutations that match the causal ordering and*
223 *assume that all dependencies within chain components are induced by selection bias or mutual*
224 *interactions, the resulting asymmetric causal Shapley values are equivalent to the asymmetric*
225 *conditional Shapley values as defined in [3].*

226 *Proof.* Following [3], asymmetric Shapley values only include those permutations π for which
227 $i \prec_{\pi} j$ for all known ancestors i of descendants j in the causal graph. For those permutations, we are
228 guaranteed to have $\tau \prec_{\mathcal{CG}} \tau'$ for all $\tau, \tau' \in \mathcal{T}$ such that $\tau \cap S \neq \emptyset, \tau' \cap \bar{S} \neq \emptyset$. That is, the chain
229 components that contain features from S are never later in the causal ordering of the chain graph \mathcal{CG}
230 than those that contain features from \bar{S} . We then have

$$\begin{aligned}
P(\mathbf{X}_{\bar{S}} | \mathbf{x}_S) &= \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_S) \\
&= \prod_{\tau \in \mathcal{T}} P(\mathbf{X}_{\tau \cap \bar{S}} | \mathbf{X}_{pa(\tau) \cap \bar{S}}, \mathbf{x}_{pa(\tau) \cap S}, \mathbf{x}_{\tau \cap S}) = P(\mathbf{X}_{\bar{S}} | do(\mathbf{X}_S = \mathbf{x}_S)),
\end{aligned}$$

231 where in the last step we used the fact that $\mathcal{T}_{\text{confounding}} = \emptyset$. \square

Conditioning by intervention in causal chain graphs boils down to conditioning by observation, where features within components that are later in the causal ordering are excluded from the conditioning set. The asymmetric Shapley values in [3], since still based on conditioning by observation, need to choose a non-uniform distribution over permutations to achieve the same. Our analysis shows that this restriction to permutations that are consistent with the causal ordering is not needed and be considered a separate choice, orthogonal to conditioning by observation or by intervention.

With a causal interpretation of Shapley values, we can decompose our explanation to reflect the contribution of direct and indirect effects. Let us consider the contribution $\phi_i(\pi)$ of a particular permutation π and feature i . With $\underline{S} = \{j : j \prec_\pi i\}$ and $\bar{S} = \{j : j \succ_\pi i\}$, we can write:

$$\begin{aligned}\phi_i(\pi) &= \mathbb{E}[f(\mathbf{X}_{\bar{S}}, \mathbf{x}_{\underline{S} \cup i}) | do(\mathbf{X}_{\underline{S} \cup i} = \mathbf{x}_{\underline{S} \cup i})] - \mathbb{E}[f(\mathbf{X}_{\bar{S} \cup i}, \mathbf{x}_{\underline{S}}) | do(\mathbf{X}_{\underline{S}} = \mathbf{x}_{\underline{S}})] && \text{(total effect)} \\ &= \mathbb{E}[f(\mathbf{X}_{\bar{S}}, \mathbf{x}_{\underline{S} \cup i}) | do(\mathbf{X}_{\underline{S}} = \mathbf{x}_{\underline{S}})] - \mathbb{E}[f(\mathbf{X}_{\bar{S} \cup i}, \mathbf{x}_{\underline{S}}) | do(\mathbf{X}_{\underline{S}} = \mathbf{x}_{\underline{S}})] + && \text{(direct effect)} \\ &\quad \mathbb{E}[f(\mathbf{X}_{\bar{S}}, \mathbf{x}_{\underline{S} \cup i}) | do(\mathbf{X}_{\underline{S} \cup i} = \mathbf{x}_{\underline{S} \cup i})] - \mathbb{E}[f(\mathbf{X}_{\bar{S}}, \mathbf{x}_{\underline{S} \cup i}) | do(\mathbf{X}_{\underline{S}} = \mathbf{x}_{\underline{S}})] && \text{(indirect effect)}\end{aligned}$$

The direct effect measures the expected change in model output when the stochastic feature X_i is replaced by its feature value x_i , without changing the distribution of the other ‘out-of-coalition’ features. The indirect effect measures the difference in expectation when the distribution of the other ‘out-of-coalition’ features changes due to the additional intervention $do(X_i = x_i)$. For any pair of features (i, j) with $j \succ_{CG} i$, the indirect effect of feature i through feature j is added to the Shapley value for i , if and only if $j \succ_\pi i$. This goes at the expense of the direct effect of feature j , essentially because when feature j is set to its value, the direct effect is computed relative to the expectation conditioned upon intervention given a set of features, including x_i . Asymmetric (causal) Shapley values only incorporate permutations π with $j \succ_\pi i$ if $j \succ_{CG} i$. With symmetric causal Shapley values, $j \succ_\pi i$ in half of the permutations π : in the other half feature j is intervened upon before feature i and there is no indirect effect to be accounted for. This makes the indirect effect in asymmetric Shapley values roughly a factor two times the indirect effect in symmetric Shapley values.

Turn the above into a Theorem? Bit hard to make concrete.

So, to be able to compute the expectations in the Shapley equations under an interventional interpretation, we need to specify (1) a partial order and (2) whether any dependencies between features that are treated on an equal footing are most likely the result of mutual interaction or selection, or of a common confounder. Based on this information, any expectation by intervention can be translated to an expectation by observation.

To compute these expectations, we can rely on the various methods that have been proposed to compute conditional Shapley values [1, 3]. Following [1], we will assume a multivariate Gaussian distribution for $P(\mathbf{X})$ that we estimate from the training data. Alternative proposals include assuming a Gaussian copula distribution, estimating from the empirical (conditional) distribution (both from [1]) and a variational autoencoder [3].

5 Experiments

From here on just rough text and ideas.

Show that it works. For now: example on bike rental. Do we predict more bike shares on a warm, but cloudy day in August because of the season or because of the weather? ADNI as another example? Tried German Credit Data, but hard to see differences between causal and conditioning, mainly because the features, such as gender and age, that can be considered causes of some of the others, hardly affect the prediction. Other suggestions? Currently using a relatively straightforward adaptation of the code of [1]. How to describe this? Do we need to publish the code? Do we need to show results for asymmetric Shapley values as well? If so, need to dig deeper into the code. Also: currently no code for handling discrete variables. Could connect to Ruifei’s Gaussian copula’s for mixed missing data, if needed?

6 Discussion

Whether to use conditional or causal Shapley values depends on interpretation: what happens if we set the features to their values compared to what happens to when we observe that features obtain

279 their values. Both interpretations are fine, but the latter enables us to leverage additional available
280 information regarding the causal structure.

281 Marginal Shapley values are perfectly fine when dependencies are purely the result of confounding,
282 as argued in [4]. However, when these dependencies are the result of causal relationships or even due
283 to selection bias or mutual feedback, then expectations under conditioning by intervention do not
284 simplify to marginal expectations.

285 We proposed a novel algorithm to compute causal Shapley values, based on causal chain graphs.
286 All that a practitioner needs to provide is a partial causal order (as for asymmetric Shapley values)
287 and a way to interpret dependencies between features that are on equal footing. Conditioning by
288 intervention becomes conditioning by observation, but only on ancestors and possibly, depending on
289 the interpretation, on features within the same component. Any existing approach for conditional
290 Shapley values can be easily adapted. If anything the expectations simplify, since there is no
291 conditioning by observation on the descendants.

292 Our analysis also provides new insights on asymmetric Shapley values: for all permutations that
293 are consistent with the causal ordering, conditioning by intervention boils down to conditioning by
294 observation. However, the current definition in [3] implicitly assumes that dependencies between
295 features that are on equal footing is due to mutual feedback or selection bias, not common con-
296 founding. Whether or not to only consider permutations that match the causal ordering depends
297 on the application and possibly one’s preference. This paper shows that it is unnecessary to forgo
298 symmetry in order to arrive at a causal interpretation of Shapley values. Roughly speaking, asym-
299 metric (causal/conditioning) Shapley values attribute all indirect effects of a causal variable through
300 a mediator on the output to the causal variable (is this the right term?), subtracting them from the
301 Shapley value of the mediator. Symmetric Shapley values are more conservative and attribute only
302 half to the causal variable.

303 Discuss non-manipulable causes as in [8]?

304 Mention that it’s easy to combine with any of TreeSHAP, KernelShap, and so on?

305 Compare with counterfactual explanations?

306 References

- 307 [1] Kjersti Aas, Martin Jullum, and Anders Løland. Explaining individual predictions when
308 features are dependent: More accurate approximations to Shapley values. *arXiv preprint*
309 *arXiv:1903.10464*, 2019.
- 310 [2] Silvia Chiappa. Path-specific counterfactual fairness. In *Proceedings of the AAAI Conference*
311 *on Artificial Intelligence*, volume 33, pages 7801–7808, 2019.
- 312 [3] Christopher Frye, Ilya Feige, and Colin Rowat. Asymmetric Shapley values: incorporating
313 causal knowledge into model-agnostic explainability. *arXiv preprint arXiv:1910.06358*, 2019.
- 314 [4] Dominik Janzing, Lenon Minorics, and Patrick Blöbaum. Feature relevance quantification in
315 explainable AI: A causality problem. *arXiv preprint arXiv:1910.13413*, 2019.
- 316 [5] Steffen L Lauritzen and Thomas S Richardson. Chain graph models and their causal in-
317 terpretations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
318 64(3):321–348, 2002.
- 319 [6] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- 320 [7] Judea Pearl. The do-calculus revisited. *arXiv preprint arXiv:1210.4852*, 2012.
- 321 [8] Judea Pearl. Does obesity shorten life? Or is it the soda? On non-manipulable causes. *Journal*
322 *of Causal Inference*, 6(2), 2018.
- 323 [9] Lloyd S Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2(28):307–
324 317, 1953.
- 325 [10] Xinpeng Shen, Sisi Ma, Prashanthi Vemuri, and Gyorgy Simon. Challenges and opportunities
326 with causal discovery algorithms: Application to Alzheimer’s pathophysiology. *Scientific*
327 *reports*, 10(1):1–12, 2020.