

Information Decomposition Based on Cooperative Game Theory

Nihat Ay^{1,2,3}, Daniel Polani⁴, Nathaniel Virgo^{1,5}

September 27, 2019

¹Max Planck Institute for Mathematics in the Sciences, Leipzig, Germany

²University of Leipzig, Leipzig, Germany

³Santa Fe Institute, Santa Fe, NM, USA

⁴University of Hertfordshire, UK

⁵Earth-Life Science Institute (ELSI), Tokyo, Japan

Abstract

We offer a new approach to the *information decomposition* problem in information theory: given a ‘target’ random variable co-distributed with multiple ‘source’ variables, how can we decompose the mutual information into a sum of non-negative terms that quantify the contributions of each random variable, not only individually but also in combination? We derive our composition from cooperative game theory. It can be seen as assigning a “fair share” of the mutual information to each combination of the source variables. Our decomposition is based on a different lattice from the usual ‘partial information decomposition’ (PID) approach, and as a consequence our decomposition has a smaller number of terms: it has analogs of the synergy and unique information terms, but lacks terms corresponding to redundancy. Because of this, it is able to obey equivalents of the axioms known as ‘local positivity’ and ‘identity’, which cannot be simultaneously satisfied by a PID measure.

Keywords: partial information decomposition, information geometry, cooperative game theory

1 Introduction

We are interested in understanding the flow of information in dynamical systems. Several tools for this have been developed over recent decades. These include the transfer entropy (Schreiber, 2000; Lizier, 2014), which assumes a dynamical systems framework and that the system consists of identifiable components which can be tracked through time, and the causal information flow (Ay and Polani, 2008), which is defined in terms of general causal Bayesian networks and the interventional calculus (Pearl, 2009). These techniques provide ways to quantify the amount of influence that one part of the network has on another.

However, there is a growing awareness that it would be useful to quantify causal influences in a more fine-grained way than offered by current techniques.

Even in one of the simplest cases, in which multiple random variables exhibit a causal influence on a single ‘target’ variable, it would be desirable to have more detailed understanding in information theoretic terms, not only of how each variable affects the target *individually*, but of how *multiple* causes interact in bringing about their effects.

Of the existing approaches to this question, perhaps the best known is the Partial Information Decomposition (PID) framework, due to Williams and Beer (2010). The PID framework proposes that the mutual information between several ‘source’ variables and a single ‘target’ can be decomposed into a sum of several terms. In the case of two sources, these terms are (i) the information that the two sources provide redundantly about the target (known as redundant information, shared information or common information); (ii) the information provided uniquely by each of the two sources, and (iii) the synergistic or complementary information, which can only be obtained by knowing both of the sources simultaneously.

However, the axioms proposed by Williams and Beer do not completely determine these quantities. As a result, many PID measures have been proposed in the literature, each satisfying different additional properties beyond the ones given by Williams and Beer. Several approaches have been proposed. Among these are several that are based on information geometry (Harder et al., 2013; Bertschinger et al., 2014; Perrone and Ay, 2016; Olbrich et al., 2015; Griffith and Koch, 2014; James et al., 2019), which we build upon here.

Generalizing towards the case of three or more input variables has turned out to be more problematic under the PID framework. One of the most intuitive additional axioms proposed is known as the identity axiom, proposed by Harder et al. (2013), but it was shown by Rauh et al. (2014) that no measure can exist that obeys both Williams and Beer’s axioms (including “local positivity”) and the identity axiom. Because of this, there are a number of proposed PID measures that relax either the identity axiom or the local positivity axiom of Williams and Beer, or both. Such approaches include (Ince, 2017; Finn and Lizier, 2018; Kolchinsky, 2019). Another promising class of approaches involve changing to a slightly different perspective, for example, by considering the full joint distribution between multiple random variables, rather than singling out a single variable as the target (Rosas et al., 2016; James and Crutchfield, 2017). In the present paper, we present a different decomposition of the mutual information between a set of sources and a target. Our decomposition obeys analogs of both the local positivity and identity axioms, but it has a smaller number of terms than the partial information decomposition.

A related, but different, approach to multivariate information can be found in a family of measures that attempt quantify the complexity of a set of random variables, often also divided into input and output variables. These include Amari’s *hierarchical decomposition* (Amari, 2001), as well as Ay’s measure of complexity (Ay, 2015), and several measures that have arisen in the context of Integrated Information Theory (IIT), such as (Oizumi et al., 2016). This family of measures is reviewed in (Amari et al., 2016), which describes their relationships in terms of information geometry.

In this paper we are interested in a similar setup to the PID framework, in which several random variables, X_1, X_2, \dots, X_n , which we term *input random variables*, exhibit causal influences on a *target* random variable Y . This results in a joint distribution between X_1, X_2, \dots, X_n and Y . The mutual information

and conditional mutual information can be used to quantify the influence of each cause individually, as well as the conditional influence that one input variable has, once the value of another has been taken into account. However, in general it would be desirable to decompose the relationships between the causal influences more finely than the traditional conditional mutual information makes possible.

The promise of the PID approach was that it would offer a ready-made or at least preferred solution to this question. A PID measure would have allowed us to quantify not only the overall influence of X_1 upon Y , but also the extent to which it has a *unique* causal influence, which could be interpreted as distinct from that of the other causes; additionally, it would also allow synergistic or redundant causal influences to be quantified. This could be done simply by applying the PID to the joint distribution $X_1 X_2 \dots X_n Y$. This is, broadly speaking, the approach taken by Lizier et al. (2014). However, the lack of a non-negative PID measure for three or more input variables makes it difficult to interpret the decomposition in the case of three or more causes.

Here we propose a different approach, with slightly more modest goals than PID, in that we do not attempt to quantify redundancy. Instead, we provide a decomposition of the mutual information $I(X_1, X_2, \dots, X_n; Y)$ into a sum of terms corresponding to every possible subset of the causes (e.g. $\{X_1\}$, $\{X_1, X_2\}$, etc.). These terms resemble the unique information and synergy terms in the PID framework. We show that the terms represent, in a well-defined sense, a uniquely “fair apportionment” of the total mutual information into the contribution provided by each subset of sources. A set such as $\{X_1, X_2\}$ will make a contribution of zero if it provides no new information beyond that which is already provided by its subsets. How to achieve this will be made precise below. In this sense our measure plays a similar role to that of synergy in the PID lattice. However, our measure does not attempt to quantify redundancy, and as such it is not a solution to the PID problem. Because of this, we are able to give a non-negative decomposition for an arbitrary number of inputs, which obeys an analog of the identity axiom for PID measures.

To state our problem more precisely: we consider random variables X_1, \dots, X_n , the *input variables*, and Y , the *output variable*. We restrict ourselves to the case where these variables have finite state sets \mathbf{X}_i , $i = 1, \dots, n$, and \mathbf{Y} , but we expect our measure to generalise well to cases such as Gaussian models in which the state spaces are continuous.

We write V for the set of all input variables. Our goal is a decomposition of the mutual information, $I(X_1, \dots, X_n; Y)$ into a sum of terms corresponding to every subset of $\{X_1, \dots, X_n\}$. We refer to a set of input variables as a *predictor*. In other words, we will write

$$I(X_1, \dots, X_n; Y) = \sum_{A \in 2^V} I_A(X_1, \dots, X_n; Y) ,$$

where we write 2^V for the power set of V . For a given predictor A , the term $I_A(X_1, \dots, X_n; Y)$ shall indicate the proportion of the total mutual information that is contributed by A , beyond what is already provided by its subsets. How to do this will be made precise below. We call $I_A(X_1, \dots, X_n; Y)$ the *information contribution* of A to Y .

To construct our measure of information contribution we proceed in two steps. We begin by defining the mutual information provided by certain *sets*

of predictors, i.e. sets of sets of input variables. We do this via a sublattice of the lattice of probability distributions that James et al. (2019) termed the “constraint lattice.” The same lattice has appeared in the literature previously, within the topic of reconstructability analysis (Zwick, 2004). Having established the information contribution of each set of predictors, we then assign a contribution to each individual predictor by a method that involves summing over the maximal chains of the constraint lattice.

We then show that this procedure of summing over maximal chains can be derived using cooperative game theory. We can conceptualise our measure in terms of a cooperative game, in which each set of predictors is thought of as a coalition of players. Each coalition is assigned a score corresponding to the information they jointly provide about the target. Our measure can then be derived via a known generalisation of the Shapley value due to Faigle and Kern (1992), which assigns a score to each individual player (i.e. predictor) based on its average performance among all the coalitions in which it takes part, while respecting additional precedence constraints.

Since our measure is based on the constraint lattice, we review this concept in depth in section 2. We approach the constraint lattice from the perspective of information geometry and state its relationship to known results in that field. In section 3 we consider a sublattice of the constraint lattice which we term the *input lattice*, which allows us to define a quantity corresponding to the information that a set of predictors provides about the target. From this we derive our measure by summing over the maximal chains of the input lattice. After proving some properties of our information contribution measure and giving some examples (sections 5 and 6), we then make the connection to cooperative game theory in section 7, proving that our measure is equivalent to the generalised Shapley value of Faigle and Kern (1992).

2 Background: the constraint lattice

We begin by defining the so-called “constraint lattice” of James et al. (2019), which has also been defined previously in the context of reconstructability analysis (Zwick, 2004). This section serves to summarise previous work and to establish notation for the following sections.

2.1 Lattices of simplicial complexes

Suppose we have a set W of co-distributed random variables, $W = \{Z_1, Z_2, \dots, Z_m\}$. Subsets of W may also be considered as random variables. For example, $\{Z_1, Z_2\}$, which we also write $Z_1 Z_2$, can be thought of as a random variable whose sample space is the Cartesian product of the sample spaces of Z_1 and Z_2 . The constraint lattice is defined in terms of members of the power set 2^W .

For reasons that will be explained below, we want to put some restrictions on which members of 2^W are permitted as elements of the lattice. This may be done in terms of two different concepts, *antichains* or *simplicial complexes*. It is standard in the literature to define the constraint lattice in terms of antichains. However, in making the connection to cooperative game theory it will be more convenient to talk in terms of simplicial complexes instead. For this reason, we

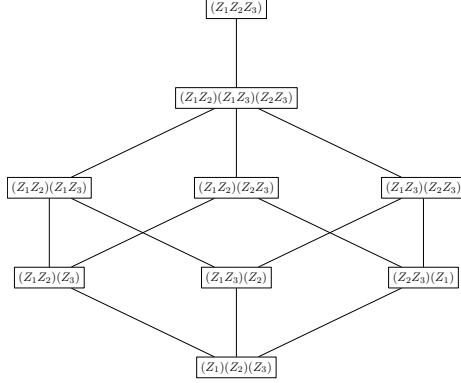


Figure 1: The Hasse diagram for the constraint lattice, as defined by (Zwick, 2004; James et al., 2019), for three random variables, $W = \{Z_1, Z_2, Z_3\}$.

define both terms here, but define the constraint lattice in terms of simplicial complexes.

A set of sets \mathcal{S} is called an antichain if for every $A \in \mathcal{S}$, no subset $B \subset A$ is a member of \mathcal{S} , i.e.

$$A \in \mathcal{S}, B \subset A \Rightarrow B \notin \mathcal{S}. \quad (1)$$

Similarly, a set of sets \mathcal{S} is called a simplicial complex if for every $A \in \mathcal{S}$, every subset $B \subset A$ is a member of \mathcal{S} ,

$$A \in \mathcal{S}, B \subset A \Rightarrow B \in \mathcal{S}. \quad (2)$$

Both of these concepts can be defined more generally in terms of arbitrary partial orders, but here we need only their definitions in terms of sets. In some contexts the empty set is considered a simplicial complex, but for most of this paper we consider only non-empty simplicial complexes.

We note that one can convert an antichain \mathcal{S} into a simplicial complex by adding every subset of every member of \mathcal{S} , and one can restore the antichain by removing every element that is a subset of some other element. This gives a one-to-one correspondence between antichains and simplicial complexes, which allows us to use the two concepts somewhat interchangeably.

We define the constraint lattice in terms of simplicial complexes whose members *cover* W , meaning those simplicial complexes \mathcal{S} for which each element of W appears at least once in one of the members of \mathcal{S} . That is, \mathcal{S} covers W if $\bigcup_{A \in \mathcal{S}} A = W$. Such a simplicial complex is termed a *simplicial complex cover* of W .

The following partial order may be defined on simplicial complexes:

$$\mathcal{S} \leq \mathcal{S}' \text{ if and only if } \forall A \in \mathcal{S}, \exists B \in \mathcal{S}': A \subseteq B. \quad (3)$$

The constraint lattice is composed of the simplicial complexes that cover W , with this partial order.

The resulting lattice is illustrated in fig. 1. In the figures and elsewhere, we use the following shortcut notation for simplicial complexes: we take the corresponding antichain, write its elements as lists surrounded by parentheses,

and concatenate them. For example, the notation $(Z_1 Z_2)(Z_3)$ refers to the simplicial complex $\{\{Z_1, Z_2\}, \{Z_1\}, \{Z_2\}, \{Z_3\}, \emptyset\}$.

It is helpful to introduce some definitions from lattice theory. We write $\mathcal{S} < \mathcal{T}$ if $\mathcal{S} \leq \mathcal{T}$ and $\mathcal{S} \neq \mathcal{T}$. We say that a lattice element \mathcal{T} *covers* an element \mathcal{S} , written $\mathcal{S} \prec \mathcal{T}$, if $\mathcal{S} < \mathcal{T}$ and there exists no element \mathcal{U} such that $\mathcal{S} < \mathcal{U} < \mathcal{T}$. A sequence of elements $\mathcal{S}_1, \dots, \mathcal{S}_k$ is called a *chain* if $\mathcal{S}_1 < \mathcal{S}_2 < \dots < \mathcal{S}_k$. If we have in addition that $\mathcal{S}_1 \prec \mathcal{S}_2 \prec \dots \prec \mathcal{S}_k$, then it is called a *maximal chain*.

In fig. 1, the relationship $\mathcal{S} \prec \mathcal{T}$ is indicated by drawing \mathcal{T} above \mathcal{S} and connecting the elements with an edge. The resulting graph is called the Hasse diagram of the lattice. The maximal chains are the directed paths from the bottom node in fig. 1 to the top node.

2.2 Constraints and split distributions

Let $p = p(Z_1, \dots, Z_m)$ be the joint probability distribution of the members of W . We call this the *true distribution*. Following (James et al., 2019) and (Zwick, 2004), we now wish to associate with each simplicial complex cover \mathcal{S} of W a joint distribution $p_{\mathcal{S}} = p_{\mathcal{S}}(Z_1, \dots, Z_m)$. In the spirit of (Ay, 2015; Oizumi et al., 2016; Amari et al., 2016) we term these *split distributions*. Each split distribution captures only some of the correlations present in the true distribution, and we can think of the remaining correlations as being split apart, or forced to be as small as possible.

Specifically, each split distribution $p_{\mathcal{S}}$ is constructed so that it captures the correlations associated with the members of \mathcal{S} , in the sense that $p_{\mathcal{S}}(A) = p(A)$, for every $A \in \mathcal{S}$. This defines a family of distributions, and from this family we choose the one with the maximum entropy. Intuitively, the maximum entropy distribution is the least correlated one in the family, so it excludes any additional correlations aside from those specified by \mathcal{S} .

In the remainder of this section, we define the split distributions more rigorously, alongside some related objects, and we point out an important property, which follows from the so-called Pythagorean theorem of information geometry. This section is largely a review of previous work, and makes a connection between the constraint lattice of (James et al., 2019; Zwick, 2004) and the language of information geometry (Ay et al., 2017, chapter 2).

Let Δ be the set of all joint probability distributions of the random variables in W . For a simplicial complex cover \mathcal{S} , let

$$M_{\mathcal{S}} = \{q \in \Delta : \forall A \in \mathcal{S}, q(A) = p(A)\}.$$

That is, $M_{\mathcal{S}}$ is the set of all probability distributions for which the members of \mathcal{S} have the same marginal distributions as in the true distribution p . Note that if the constraint $q(A) = p(A)$ holds for some $A \subseteq W$, then it will also automatically hold for $B \subseteq A$. This is the reason for considering simplicial complexes. $M_{\mathcal{S}}$ is a mixture family, and we have that $\mathcal{S} \leq \mathcal{T} \implies M_{\mathcal{S}} \supseteq M_{\mathcal{T}}$.

We can now define the split distribution $p_{\mathcal{S}}$ as

$$p_{\mathcal{S}} = \operatorname{argmax}_{q \in M_{\mathcal{S}}} H(q). \quad (4)$$

Equivalently, we can instead define the split distributions in terms of the Kullback-Leibler divergence, as we will see below. This has the advantage that it is likely

to generalise to cases such as Gaussian models in which the state space is not discrete.

There is another interpretation of the split distributions, which is interesting to note. In addition to the mixture family $M_{\mathcal{S}}$, we can also define an exponential family corresponding to a given node in the constraint lattice. This can be seen as a family of *split models*, i.e. probability distributions in which some kinds of correlation are forced to be absent. The split distribution $p_{\mathcal{S}}$ can be seen as the closest member of this exponential family to the true distribution.

To see this, we define the exponential family

$$E_{\mathcal{S}} = \left\{ q \in \Delta : q(z_1, \dots, z_m) = \prod_{A \in \mathcal{S}} \mu_A(z_1, \dots, z_m), \text{ for some set of functions } \mu_A \right\}, \quad (5)$$

where the functions μ_A have the additional requirements that $\mu_A(z_1, \dots, z_m)$ depends only on z_i for $Z_i \in A$, and $\mu_A(z_1, \dots, z_m) > 0$. $E_{\mathcal{S}}$ is an exponential family, and we have $\mathcal{S} < \mathcal{T} \implies E_{\mathcal{S}} \subseteq E_{\mathcal{T}}$.

Finally, we let $\bar{E}_{\mathcal{S}}$ be the topological closure of the set $E_{\mathcal{S}}$, meaning that $\bar{E}_{\mathcal{S}}$ contains every member of $E_{\mathcal{S}}$, and in addition also contains all the limit points of sequences in $E_{\mathcal{S}}$. The difference is that $E_{\mathcal{S}}$ does not contain distributions with zero-probability outcomes, whereas the closure $\bar{E}_{\mathcal{S}}$ does.

Note that we cannot obtain $\bar{E}_{\mathcal{S}}$ by simply relaxing the condition that $\mu_A(z_1, \dots, z_m) > 0$. This is because although every member of $E_{\mathcal{S}}$ must factorise according to eq. (5), the limit points on the boundary of the simplex can fail to factorise in the same way. An example of this is given by (Lauritzen, 1996, Example 3.10). These limit points must be included in order to make sure the split distribution is always defined.

It is a known result in information geometry (Ay et al., 2017, Theorem 2.8) that for any \mathcal{S} , the sets $M_{\mathcal{S}}$ and $\bar{E}_{\mathcal{S}}$ intersect at a single point. In fact this point is the split distribution $p_{\mathcal{S}}$. With the Kullback-Leibler divergence

$$D_{\text{KL}}(q||p) = \sum_{z_1, \dots, z_m} q(z_1, \dots, z_m) \log \frac{q(z_1, \dots, z_m)}{p(z_1, \dots, z_m)},$$

we can equivalently characterise $p_{\mathcal{S}}$ by

$$p_{\mathcal{S}} = \operatorname{argmin}_{q \in M_{\mathcal{S}}} D_{\text{KL}}(q||u), \quad (6)$$

where u denotes the uniform distribution. This directly follows from (4). A further equivalent characterisation of $p_{\mathcal{S}}$ is given by

$$p_{\mathcal{S}} = \operatorname{argmin}_{q \in E_{\mathcal{S}}} D_{\text{KL}}(p||q). \quad (7)$$

In the terminology of information geometry, eq. (6) is an I-projection (information projection) and eq. (7) is an rI-projection (reverse I-projection). The classical theory of these information projections has been greatly extended by Csiszár and Matúš (2003, 2004).

We also have the so-called *Pythagorean theorem* of information geometry (Amari and Nagaoka, 2007), which in our notation says that for simplicial complex covers $\mathcal{S} < \mathcal{T} < \mathcal{U}$,

$$D_{\text{KL}}(p_{\mathcal{U}}||p_{\mathcal{S}}) = D_{\text{KL}}(p_{\mathcal{U}}||p_{\mathcal{T}}) + D_{\text{KL}}(p_{\mathcal{T}}||p_{\mathcal{S}}). \quad (8)$$

Equation (8) can be extended to any *chain* of elements in the constraint lattice $\mathcal{S}_1 < \mathcal{S}_2 < \dots < \mathcal{S}_k$, to give

$$D_{\text{KL}}(p_{\mathcal{S}_k} \| p_{\mathcal{S}_1}) = \sum_{i=2}^k D_{\text{KL}}(\mathcal{S}_i \| \mathcal{S}_{i-1}). \quad (9)$$

This will be crucial in defining our information contribution measure below.

Consider the top node in the constraint lattice, given by (Z_1, \dots, Z_m) , which we denote \top . We have $p_{\top} = p$. That is, the split distribution corresponding to \top is equal to the true distribution.

Since we are considering only simplicial complex covers of W , the bottom node of the lattice is given by $(Z_1) \dots (Z_m)$, which we denote \perp . We have $p_{\perp}(z_1, \dots, z_m) = p(z_1) \dots p(z_m)$. That is, its split distribution is given by the product of the marginal distributions for all the members of W .

Together with eq. (6), this allows us to interpret $p_{\mathcal{S}}$ as the distribution that is *as decorrelated as possible* (i.e. closest to the product distribution, in the Kullback-Leibler sense), subject to the constraint that the marginals of the members of \mathcal{S} match those of the true distribution. Alternatively, via eq. (7), we can see it as the distribution that is as close to the true distribution as possible, subject to the constraint that it lies in the closure of the exponential family $\bar{E}_{\mathcal{S}}$.

For a general antichain cover \mathcal{S} , the split distribution $p_{\mathcal{S}}$ may not have an analytical solution, and instead must be found numerically. One family of techniques for this is iterative scaling (Csiszár and Shields, 2004, chapter 5), which was used to calculate the examples below. Alternatively, one may solve eq. (6) as a numerical optimisation problem, starting from an element such as \perp with a known split distribution. This yields a convex optimisation problem with linear constraints, but it is not always well conditioned.

Finally, given an antichain cover \mathcal{S} , we define $I_{\mathcal{S}} := D_{\text{KL}}(p_{\top} \| p_{\mathcal{S}})$. This can be thought of as the amount of information that is present in the true distribution p_{\top} but is not present in $p_{\mathcal{S}}$. Note that due to the Pythagorean relation (eq. (8)) we have $D_{\text{KL}}(p_{\top} \| p_{\mathcal{S}}) = I_{\mathcal{S}} - I_{\mathcal{T}}$, for any antichain covers $\mathcal{S} \leq \mathcal{T}$. The quantity $I_{\mathcal{S}}$ turns out to be a useful generalisation of the mutual information, as shown in the following examples.

Example 2.1. Independence. Suppose $W = \{Z_1, Z_2\}$, and let $\mathcal{S} = (Z_1)(Z_2)$. Then $\bar{E}_{\mathcal{S}}$ is the set of distributions q that can be expressed as a product $q(z_1, z_2) = \mu_1(z_1)\mu_2(z_2)$, which we may also write $q(z_1, z_2) = q(z_1)q(z_2)$. So $\bar{E}_{\mathcal{S}}$ is the set of distributions for which Z_1 and Z_2 are independent. We have that $p_{\mathcal{S}} = p(z_1)p(z_2)$, and consequently, it is straightforward to show that in this example, $D_{\text{KL}}(p_{\top} \| p_{\mathcal{S}}) = I(X_1; X_2)$.

Example 2.2. Conditional independence. Suppose $W = \{Z_1, Z_2, Z_3\}$, and let $\mathcal{S} = (Z_1 Z_3)(Z_2 Z_3)$. Then $\bar{E}_{\mathcal{S}}$ is the set of distributions q that can be expressed as a product $q(z_1, z_2, z_3) = \mu_1(z_1, z_3)\mu_2(z_2, z_3)$. These are the distributions for which $q(z_1, z_2, z_3) = q(z_3)q(z_1|z_3)q(z_2|z_3)$, i.e. for which $Z_1 \perp\!\!\!\perp_q Z_2 \mid Z_3$. So in this case $E_{\mathcal{S}}$ can be seen as a conditional independence constraint. It is straightforward to show that that $p_{\mathcal{S}}(z_1, z_2, z_3) = p(z_3)p(z_1|z_3)p(z_2|z_3)$, and consequently $D_{\text{KL}}(p_{\top} \| p_{\mathcal{S}}) = I(X_1; X_2 | X_3)$.

Example 2.3. Amari's triplewise information. Suppose $W = \{Z_1, Z_2, Z_3\}$, and let $\mathcal{S} = (Z_1 Z_2)(Z_1 Z_3)(Z_2 Z_3)$. Then $E_{\mathcal{S}}$ is the set of distributions q that can

be expressed as a product $q(z_1, z_2, z_3) = \mu_1(z_1, z_2)\mu_2(z_1, z_3)\mu_3(z_2, z_3)$. Unlike the previous two examples, there is no analytic expression for μ_1 , μ_2 and μ_3 in terms of the probabilities $q(z_1, z_2, z_3)$. However, Amari (2001) argued that \bar{E}_S can be interpreted as the set of distributions in which there are no three-way, or “triplewise” interactions between the variables Z_1 , Z_2 and Z_3 , beyond those that are implied by their pairwise interactions. The split distribution p_S can be calculated numerically as described above, in order to obtain the quantity $D_{\text{KL}}(p_{\top} \| p_S)$, which quantifies the amount of information present in the triplewise interactions. Amari (2001) gives a straightforward generalisation, allowing n -way interactions to be quantified, among n or more random variables. As an example of triplewise information, consider the case where Z_1 and Z_2 are uniformly distributed binary variables, and $Z_3 = Z_1 \text{ XOR } Z_2$. In this case, in the split distribution $p_{(Z_1 Z_2)(Z_1 Z_3)(Z_2 Z_3)}$ all three variables are independent. The split distribution has 8 equally likely outcomes while the true distribution has 4 equally likely outcomes, leading to a triplewise information of 1 bit.

3 The input lattice

The constraint lattice is defined in terms of an arbitrary set of random variables $W = \{Z_1, \dots, Z_m\}$. We are interested specifically in the case where W is composed of a set of input variables X_1, \dots, X_n and a target variable Y . We write V for the set of input variables, so $W = V \cup \{Y\}$.

We wish to decompose the mutual information $I(X_1, \dots, X_n; Y)$ into a sum of terms $I_A(X_1, \dots, X_n; Y)$, one for each subset A of the input variables. To do this, we start by noting that

$$I(X_1, \dots, X_n; Y) = D_{\text{KL}}(p_{(X_1, \dots, X_n, Y)} \| p_{(X_1, \dots, X_n)}(Y)) .$$

Because of this, we can use the constraint lattice to derive decompositions of the mutual information.

Consider the set of lattice elements \mathcal{S} such that $(X_1, \dots, X_n)(Y) \leq \mathcal{S}$. This set forms a sublattice of the constraint lattice, i.e. a lattice under the same partial order. We call this sublattice the input lattice. The input lattice is highlighted in red in fig. 2, left.

Each element of the input lattice may be associated with a simplicial complex over the input variables only. That is, a non-empty set \mathfrak{S} of subsets of V , with the condition that every subset of a member must also be a member. (Unlike the elements of the constraint lattice, \mathfrak{S} need not cover V .) We use a Fraktur font for simplicial complexes over the input variables only, to distinguish them from simplicial complex covers of W . Their relationship to the input sublattice can be seen by noting that if $(X_1 \dots X_n)(Y) \leq \mathcal{S}$ then \mathcal{S} must include the element $\{X_1, \dots, X_n\}$ and its subsets. In addition \mathcal{S} contains elements of the form $A \cup \{Y\}$, where A is a subset of the input variables. These sets of input variables must by themselves form a simplicial complex, in order for \mathcal{S} to be a simplicial complex. This is the simplicial complex \mathfrak{S} over the input variables corresponding to \mathcal{S} .

Formally, given a simplicial complex \mathfrak{S} over the input variables, the corresponding member of the constraint lattice is given by

$$\sigma(\mathfrak{S}) = (X_1 \dots X_n) \cup \{A \cup \{Y\} : A \in \mathfrak{S}\} .$$

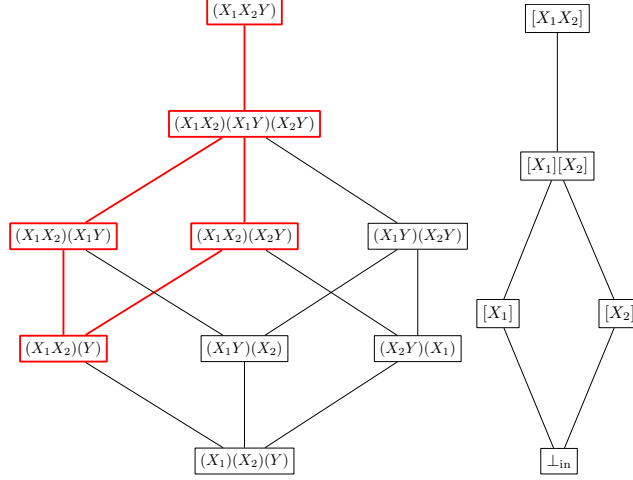


Figure 2: (Left) the Hasse diagram for the constraint lattice for $W = \{X_1, X_2, Y\}$. Highlighted in red bold is the sublattice that we call the input lattice, which provides decompositions of $I(X_1, X_2, \dots, X_n; Y)$. (Right) the input lattice alone, with the nodes labelled using simplicial complexes over the input random variables, rather than all random variables. The square brackets indicate simplicial complexes over the inputs, rather than over all random variables. The two lattices are related by the mapping σ , defined in the text.

For any \mathfrak{S} , we have that $\sigma(\mathfrak{S})$ is a simplicial complex cover of W , and $\sigma(\mathfrak{S}) \geq (X_1 \dots X_n)(Y)$. In fact, σ is an order-preserving invertible map from the lattice of simplicial complexes \mathfrak{S} over X to the sublattice of simplicial complex covers of W given by $(X_1 \dots X_n)(Y) \leq \mathfrak{S}$. This allows us to think of the elements of the input lattice as corresponding to simplicial complexes over the input variables.

The mapping is illustrated in fig. 2. When writing simplicial complexes over the input variables explicitly we use square brackets, in order to distinguish them from simplicial complexes over W . So for example, the simplicial complex $[X_1][X_2]$ over the input variables corresponds to the simplicial complex $(X_1 X_2)(X_1 Y)(X_2 Y)$ over W . We write the bottom node of the input lattice as \perp_{in} , which is equal to $\{\emptyset\}$ when considered as a simplicial complex over the input variables, or $(X_1 \dots X_n)(Y)$ when considered as a simplicial complex cover of W .

Every chain in this sublattice provides a decomposition of $I(X_1, \dots, X_n; Y)$ into a sum of non-negative terms. An example of such a decomposition is the chain rule for mutual information,

$$I(X_1, X_2; Y) = I(X_1; Y) + I(X_2; Y|X_1),$$

which can be derived by applying the Pythagorean theorem to the (not maximal) chain

$$\perp_{\text{in}} < [X_2] < [X_1 X_2].$$

This corresponds to

$$(X_1 X_2)(Y) < (X_1 X_2)(X_1 Y) < (X_1 X_2 Y)$$

when considered as elements of the constraint lattice. Applying Equation 8, we have

$$D_{\text{KL}}(p_{(X_1 X_2 Y)} \| p_{(X_1 X_2)(Y)}) = D_{\text{KL}}(p_{(X_1 X_2)(X_1 Y)} \| p_{(X_1 X_2)(Y)}) + D_{\text{KL}}(p_{(X_1 X_2 Y)} \| p_{(X_1 X_2)(X_1 Y)}), \quad (10)$$

which corresponds term-by-term to the chain rule for mutual information. While this chain is not maximal, considering the maximal chain, however, yields a more fine-grained decomposition:

$$\perp_{\text{in}} < [X_2] < [X_1][X_2] < [X_1 X_2],$$

This, in turn, yields an information decomposition with three non-negative terms,

$$I(X_1, X_2; Y) = I(X_1; Y) + (I(X_2; Y|X_1) - I_3(X_1, X_2, Y)) + I_3(X_1, X_2, Y),$$

where $I_3(X_1, X_2, Y)$ is Amari’s triplewise information. (See example 2.3 above.)

In this way we can write $I(X_1, \dots, X_n; Y)$ as a sum of non-negative terms in many different ways. However, these decompositions in general treat the input variables asymmetrically. The decompositions are “path-dependent,” in the sense that they depend on which particular chain is chosen. In the next section we turn these path-dependent decompositions into a single path-independent one by suitably averaging over the maximal chains.

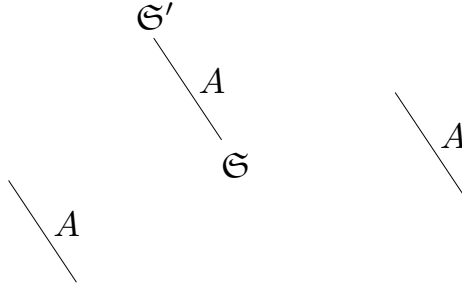
4 Defining the information contribution as a sum over chains

We now extend the decomposition along individual chains of the input lattice to a path-independent information decomposition: this decomposition will define a separate information contribution for each of the non-empty subsets A of V .

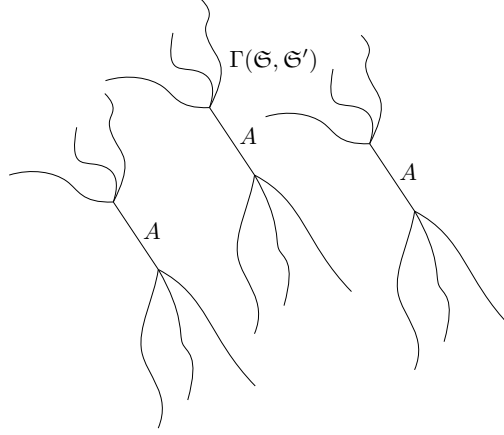
In order to do so, consider the set Γ of all maximal chains in the input lattice, that is, all directed paths from \perp_{in} to $[X_1, \dots, X_n]$. Consider a maximal chain $\gamma \in \Gamma$. For any index l in the chain, the collection $\gamma(l)$ of subsets forms a simplicial complex and for each transition from $\gamma(l)$ to $\gamma(l+1)$ a subset A of V is added to the simplicial complex $\gamma(l)$ until the topmost simplicial complex which ends up containing all subsets of V . In particular, the chain has the property that all non-empty subsets A of V are being added at some point along a chain γ .

In particular, this ensures that there is exactly one $l_\gamma(A)$ that satisfies the following condition: all simplicial complexes $\gamma(l)$, $0 \leq l < l_\gamma(A)$, do not contain A , and all simplicial complexes $\gamma(l)$, $l_\gamma(A) \leq l \leq 2^n - 1$, do contain A ; i.e. $l_\gamma(A)$ denotes the step in the chain γ at which A is added (note that the empty set \emptyset is necessarily contained in the first complex of each chain, i.e. $l_\gamma(\emptyset) = 0$).

Based on this, we now derive a decomposition of the mutual information between inputs and output “aligned” with respect to a particular subset A of inputs. For this purpose, consider the set \mathcal{E}_A of all edges $(\mathfrak{S}, \mathfrak{S}')$ where \mathfrak{S}' is obtained from \mathfrak{S} by adding A , i.e. where $\mathfrak{S}' = \mathfrak{S} \uplus \{A\}$:



We furthermore now subdivide the set Γ into classes of maximal chains, grouped by specific edges $(\mathfrak{S}, \mathfrak{S}')$. Denote by $\Gamma(\mathfrak{S}, \mathfrak{S}')$ the set of all maximal chains $\gamma \in \Gamma$ that contain this particular edge $(\mathfrak{S}, \mathfrak{S}')$:



Then, for any non-empty subset A , one has the following partition:

$$\Gamma = \bigsqcup_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \Gamma(\mathfrak{S}, \mathfrak{S}').$$

Every maximal chain is accounted for in this disjoint union, because for every maximal chain there is exactly one step (edge) at which the set A is added. This is illustrated in fig. 3 for the case of three input variables.

We now consider a probability weighting over the maximal chains, that is, a

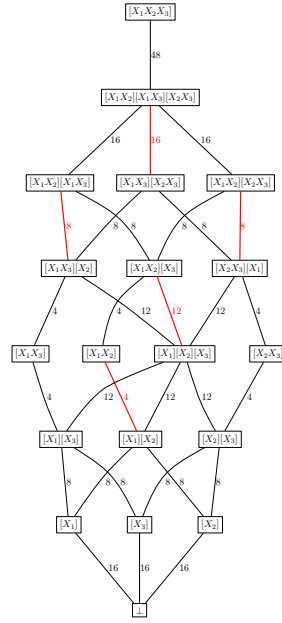


Figure 3: The input lattice for three inputs. Each edge is labelled with the total number of maximal chains that pass through that that edge. The edges where the subset $\{X_1, X_2\}$ appears for the first time are highlighted in red. Each maximal chain passes through exactly one of these edges. The contribution of $\{X_1, X_2\}$ to the total information is calculated by averaging over these edges, weighted by their path counts (the numbers in red.) In this lattice there are 48 maximal chains in total.

set of weights $\mu(\gamma)$ such that $\sum_{\gamma \in \Gamma} \mu(\gamma) = 1$. We obtain

$$I(X_1, \dots, X_n; Y) = D_{\text{KL}}(p_{[X_1, \dots, X_n]} \parallel p_{\perp_{\text{in}}}) \quad (11)$$

$$= \sum_{\gamma \in \Gamma} \mu(\gamma) \sum_{l=1}^{2^n-1} D_{\text{KL}}(p_{\gamma(l)} \parallel p_{\gamma(l-1)}) \quad (12)$$

$$= \sum_{\gamma \in \Gamma} \mu(\gamma) \sum_{\emptyset \neq A \subseteq V} D_{\text{KL}}(p_{\gamma(l_{\gamma}(A))} \parallel p_{\gamma(l_{\gamma}(A)-1)}) \quad (13)$$

$$= \sum_{\emptyset \neq A \subseteq V} \sum_{\gamma \in \Gamma} \mu(\gamma) D_{\text{KL}}(p_{\gamma(l_{\gamma}(A))} \parallel p_{\gamma(l_{\gamma}(A)-1)}) \quad (14)$$

$$= \sum_{\emptyset \neq A \subseteq V} \sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \underbrace{\left\{ \sum_{\gamma \in \Gamma(\mathfrak{S}, \mathfrak{S}')} \mu(\gamma) \right\}}_{=1} D_{\text{KL}}(p_{\mathfrak{S}'} \parallel p_{\mathfrak{S}}) \quad (15)$$

$$= \sum_{\emptyset \neq A \subseteq V} \sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \mu(\mathfrak{S}, \mathfrak{S}') D_{\text{KL}}(p_{\mathfrak{S}'} \parallel p_{\mathfrak{S}}) \quad (16)$$

Here, we used the short-hand notation $\mu(\mathfrak{S}, \mathfrak{S}')$ for $\mu(\Gamma(\mathfrak{S}, \mathfrak{S}')) = \sum_{\gamma \in \Gamma(\mathfrak{S}, \mathfrak{S}')} \mu(\gamma)$. The equality (12) follows because of the Pythagorean theorem (eq. (9)) and the normalization of the weights μ . Note, via (15), that the non-negative weights in the decomposition (16) satisfy the following condition:

$$\sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \mu(\mathfrak{S}, \mathfrak{S}') = 1. \quad (17)$$

This allows us to interpret

$$I_A^{(\mu)}(X_1, \dots, X_n; Y) := \sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \mu(\mathfrak{S}, \mathfrak{S}') D_{\text{KL}}(p_{\mathfrak{S}'} \parallel p_{\mathfrak{S}}) \quad (18)$$

as the mean information in A that is not contained in a proper subset of A .

This gives us a non-negative decomposition of $I(X_1, \dots, X_n; Y)$ into terms corresponding to each subset of the input variables, but note that this decomposition is dependent on the choice of weights μ .

A natural choice for the weights μ would be simply to choose the uniform distribution, i.e. $\mu(\gamma) = 1/|\Gamma|$ for all γ . It is not completely straightforward to justify the uniform distribution over maximal chains, because there is no obvious symmetry that transforms one maximal chain into another. Note, for example, that the connectivity of the node $[X_1][X_2][X_3]$ in fig. 3 is different from that of other nodes on the same level.

Nevertheless, we will now proceed with the uniform distribution as a reasonable intuitive choice. It will be shown in section 5 that choosing μ this way gives rise to a decomposition of $I(X_1, \dots, X_n; Y)$ that has some intuitively desirable properties. For the special choice of μ as the uniform distribution we will write $I_A^{(\mu)}(X_1, \dots, X_n; Y)$ simply as $I_A(X_1, \dots, X_n; Y)$. We denote this as the *information contribution* of A to Y . In section 7 we will then proceed to show that above originally merely intuitive choice of μ as uniform distribution finds a deeper justification in the theory of cooperative game theory.

For the practical calculation of I_A we first calculate the number $n_{(\mathfrak{S}, \mathfrak{S}')}$ of maximal chains that pass through each edge $(\mathfrak{S}, \mathfrak{S}')$ in the Hasse diagram of the input lattice. These numbers are shown in fig. 3, as well as the total number of maximal chains, $|\Gamma|$. For each node \mathfrak{S} in the lattice we calculate the distribution $p_{\mathfrak{S}}$ by iterative scaling (Csiszár and Shields, 2004, chapter 5), from which we obtain $D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\{\emptyset\}})$. We then find the set \mathcal{E}_A of edges in which a given predictor A is added for the first time in a maximal chain (for the example of $A = \{X_1, X_2\}$ this is shown in red in fig. 3). We then calculate our measure I_A from the Kullback-Leibler gains accrued on these paths by adding the set A of interest, weighted by the chain counts $n_{(\mathfrak{S}, \mathfrak{S}')}$ of the respective edges:

$$I_A(X_1, \dots, X_n; Y) = \frac{1}{|\Gamma|} \sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} n_{(\mathfrak{S}, \mathfrak{S}')} (D_{\text{KL}}(p_{\mathfrak{S}'} \| p_{\{\emptyset\}}) - D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\{\emptyset\}})) . \quad (19)$$

5 Properties of the information contribution

We now prove the following properties of the information contribution, as a decomposition of the mutual information.

Theorem 1. For $A \subseteq \{X_1, \dots, X_n\}$ we have

- I. $I_A(X_1, \dots, X_n; Y) \geq 0$ (*nonnegativity*)
- II. $\sum_{A \in 2^V} I_A(X_1, \dots, X_n; Y) = I(X_1, \dots, X_n; Y)$ (*completeness*)
- III. $I_A(X_1, \dots, X_n; Y)$ is invariant under permutations of X_1, \dots, X_n . (*symmetry*)
- IV. $I_A(X_1, \dots, X_n; (X_1, \dots, X_n)) = 0$ if $|A| > 1$. (*singleton*)
- V. if $X_i = (X'_i, X''_i)$ for all i , $Y = (Y', Y'')$, and

$$p(x_1, \dots, x_n, y) = p(x'_1, \dots, x'_n, y') p(x''_1, \dots, x''_n, y'') ,$$

then

$$I_A(X_1, \dots, X_n; Y) = I_{A'}(X'_1, \dots, X'_n; Y') + I_{A''}(X''_1, \dots, X''_n; Y'') ,$$

where $A' = \{X'_i : (X'_i, X''_i) \in A\}$ and $A'' = \{X''_i : (X'_i, X''_i) \in A\}$. (*additivity*)

As we discuss below, the singleton property is somewhat analogous to the identity axiom proposed by (Harder et al., 2013) for partial information decomposition measures, which effectively says that there should be no synergy terms if the output is simply an identical copy of the input.

Proof. (I) follows from the nonnegativity of the Kullback-Leibler divergence. (II) is proved in Section 4 above. (III) is true by construction, since the values of the Kullback-Leibler divergences do not depend on the order in which the input variables are considered, and the uniform distribution over maximal chains is invariant to reordering the input variables.

To prove (IV), write $Y = (\bar{X}_1, \dots, \bar{X}_n)$, where \bar{X}_i is considered to be a copy of X_i , in the sense that X_i and \bar{X}_i are separate random variables but we have

$$p(x_1, \dots, x_n, \bar{x}_1, \dots, \bar{x}_n) = \begin{cases} p(x_1, \dots, x_n) & \text{if } x_1 = \bar{x}_1, \dots, x_n = \bar{x}_n, \\ 0 & \text{otherwise,} \end{cases} \quad (20)$$

which implies that $p(x_i, \bar{x}_i) = \delta_{x_i, \bar{x}_i} p(x_i)$, for every i . We then have

$$p_{\perp_{\text{in}}}(X_1, \dots, X_n, \bar{X}_1, \dots, \bar{X}_n) = p(X_1) \dots p(X_n) p(\bar{X}_1, \dots, \bar{X}_n).$$

Consider now any edge $(\mathfrak{S} \setminus \{A\}, \mathfrak{S})$ in the Hasse diagram of the input lattice. There are two cases to consider:

- (i) $A = \{X_i\}$ for some i . In this case $\sigma(\mathfrak{S})$ contains the element $\{X_i, Y\}$. Therefore, from its definition, the marginal $p_{\mathfrak{S}}(X_i, Y)$ must match the true marginal $p(X_i, Y)$, which implies that $p_{\mathfrak{S}}(x_i, \bar{x}_i) = \delta_{x_i, \bar{x}_i} p(x_i)$. However, $\sigma(\mathfrak{S} \setminus \{A\})$ does not contain the element $\{X_i, Y\}$, and so the marginal $p_{\mathfrak{S} \setminus \{A\}}(x_i, \bar{x}_i)$ may in general differ from $\delta_{x_i, \bar{x}_i} p(x_i)$, and $D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\mathfrak{S} \setminus A})$ can be nonzero.
- (ii) $|A| > 1$. Consider first the case that $A = \{X_i, X_j\}$. Because \mathfrak{S} is a simplicial complex, we have that $\{X_i\} \in \mathfrak{S} \setminus \{A\}$ and $\{X_j\} \in \mathfrak{S} \setminus \{A\}$. Therefore $p_{\mathfrak{S} \setminus \{A\}}$ has to match the constraints $p_{\mathfrak{S} \setminus \{A\}}(x_i, \bar{x}_i) = \delta_{x_i, \bar{x}_i} p(x_i)$ and $p_{\mathfrak{S} \setminus \{A\}}(x_j, \bar{x}_j) = \delta_{x_j, \bar{x}_j} p(x_j)$. We also have, from eq. (20), that $p_{\mathfrak{S} \setminus \{A\}}(\bar{x}_i, \bar{x}_j) = p(\bar{x}_i, \bar{x}_j)$. From these constraints we have

$$p_{\mathfrak{S} \setminus A}(x_i, x_j, \bar{x}_i, \bar{x}_j) = p(\bar{x}_i, \bar{x}_j) \delta_{x_i, \bar{x}_i} \delta_{x_j, \bar{x}_j} = p(x_i, x_j, \bar{x}_i, \bar{x}_j).$$

Therefore $p_{\mathfrak{S} \setminus A}$ already meets the constraint that the marginals for X_i, X_j, Y match those of the true distribution and minimising $D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\mathfrak{S} \setminus A})$ subject to this constraint must result in zero. The proof of this is similar if $|A| > 2$.

Therefore every term in eq. (16) will be zero if $|A| > 1$, but in general they can be nonzero if $|A| = 1$.

To prove (V) we first note the following general additivity property of the Kullback-Leibler divergence. Let Z' and Z'' be two co-distributed random variables, let $p_0(z', z'') = p_0(z') p_0(z'')$ for each z', z'' in the sample spaces of Z', Z'' , that is, render the two random variables independent according to the distribution p_0 . Then let M be a mixture family defined by constraints that depend only on either Z' or Z'' . That is,

$$M = \left\{ q : \sum_{z'} q(z') f^{(i)}(z') = F^{(i)} \quad (i = 1, \dots, r), \right. \\ \left. \sum_{z''} q(z'') g^{(j)}(z'') = G^{(j)} \quad (j = 1, \dots, s) \right\}. \quad (21)$$

Calculating $\text{argmin}_{p \in M} D_{\text{KL}}(p \| p_0)$, introducing Lagrange multipliers in the usual way, gives us

$$p(z', z'') = p_0(z') p_0(z'') e^{\sum_i \lambda_i f^{(i)}(z') + \sum_j \eta_j g^{(j)}(z'') - \psi} = p(z') p(z''),$$

where $p(z') = p_0(z')e^{\sum_i \lambda_i f^{(i)}(z') - \psi'}$ and $p(z'') = p_0(z'')e^{\sum_j \eta_j g^{(j)}(z'') - \psi''}$. Note that these are the same distributions that would be obtained if the projection were performed on each of the marginals rather than the joint distribution. We have both that Z' and Z'' remain independent after projecting onto M , and also that $D_{\text{KL}}(p(Z', Z'') \parallel p_0(Z', Z'')) = D_{\text{KL}}(p(Z') \parallel p_0(Z')) + D_{\text{KL}}(p(Z'') \parallel p_0(Z''))$.

Now consider constructing a system of random variables $X_i = (X'_i, X''_i)$, $Y = (Y', Y'')$, according to the condition of property V. Each of the split distributions is defined as a projection from the product distribution onto a mixture family. By construction, all of these mixture families satisfy eq. (21). Because of this, every term in eq. (18) can be written as a sum of the corresponding terms for the systems $\{X'_1, \dots, X'_n, Y'\}$ and $\{X''_1, \dots, X''_n, Y''\}$. The additivity property follows from this. \square

We note that these properties do not uniquely determine the information contribution measure. In particular, one could choose a different measure μ over the maximal chains besides the uniform measure; there are in general many such measures that would yield an information measure satisfying theorem 1. To see this, note that properties I, II, IV and V do not depend on the choice of measure μ , and hence don't constrain it. Property III does restrict the choice of measure, but for more than two inputs the number of paths in the lattice is greater than the number of inputs, and consequently the symmetry axiom does not provide enough constraint to uniquely specify μ . However, as argued above, the uniform measure is a natural choice, and we will show below that its use can be more systematically justified from the perspective of cooperative game theory.

5.1 Comparison to partial information decomposition

As noted above, the information contribution I_A is not a partial information decomposition (PID) measure, because it decomposes the mutual information into a different number of terms than the latter. In the case of two input variables X_1 and X_2 , the PID has four terms (synergy, redundancy, and two unique terms), whereas the information contribution has only three, $I_{\{X_1\}}$, $I_{\{X_2\}}$ and $I_{\{X_1, X_2\}}$. The joint term, $I_{\{X_1, X_2\}}$, behaves somewhat similarly to a synergy term, and the two singleton contributions $I_{\{X_1\}}$ and $I_{\{X_2\}}$ are roughly analogous to the two unique information terms, but there is no term corresponding to shared/redundant information. For more than two inputs, the terms of a partial information measure can be expressed in terms of a lattice known as the redundancy lattice Williams and Beer (2010), which is different from the constraint lattice or the input lattice discussed above.

Within the PID framework, (Harder et al., 2013) introduced the *identity axiom*, which states that a measure of redundant information I_{\cap} , should satisfy

$$I_{\cap}(X_1, X_2; (X_1, X_2)) = I(X_1; X_2) .$$

This is equivalent to saying that the corresponding measure of synergy, I_{\cup} , should be zero in the case where the output is a copy of its two input variables:

$$I_{\cup}(X_1, X_2; (X_1, X_2)) = 0 . \quad (22)$$

It was proven by Rauh et al. (2014) that there can be no non-negative PID measure that satisfies all of Williams and Beer’s axioms together with the identity axiom. This can be achieved if we restrict ourselves to two input variables, but for three or more inputs there are distributions for which it cannot be achieved. (See Example RBOJ below.)

While our information contribution measure is not a PID measure, if we take the joint term $I_{\{X_1, X_2\}}$ to be analogous to a synergy term, then the singleton decomposition property (property IV), for two inputs, is roughly analogous to eq. (22). Therefore our measure obeys an analog of the identity axiom for PID measures, alongside analogs of the non-negativity and symmetry axioms for PID measures. This is possible only because the information contribution is not a PID measure, and hence does not have to obey the precise set of Williams-Beer lattice axioms.

It is also worth comparing our information decomposition measure with the framework proposed by James and Crutchfield (2017), which seeks a different kind of information decomposition from PID. In this framework, instead of decomposing the mutual information between a set of sources and a target, one instead wishes to decompose the joint entropy $H(Z_1, \dots, Z_n)$ of several jointly distributed random variables, into a sum of terms corresponding to each subset of the variables. Our framework sits somewhere between this approach and PID, since we have the distinction between the inputs and the target, but we decompose $I(X_1, \dots, X_n; Y)$ into a sum of terms corresponding to subsets of the inputs, in a similar manner to James and Crutchfield’s proposal.

6 Examples

We now explore a few examples of our information contribution measure (which we will also sometimes denote by Shapley information decomposition, as will be justified by the game-theoretic analysis in section 7 below). Here, we apply it to joint distributions between a target and two or three inputs. Note that our framework does not require any restrictions on these joint distributions. In particular, it is expressly not assumed that the inputs are independent of one another. Importantly, the measures will in general be affected by dependent inputs, which is a desirable property of such a measure, because it has been observed before that appropriate attributions of joint interactions should depend on input correlations (see the discussion on *source* vs. *mechanistic* redundancy in Harder et al., 2013).

We take most of our examples from the literature on partial information decomposition, in particular (Williams and Beer, 2010; Griffith and Koch, 2014; Harder et al., 2013; Bertschinger et al., 2014). These examples are relatively standardised, and give some intuition for how our measure compares to PID measures.

We first explore some basic examples with two predictors, which are presented in table 1. For each of these examples, the method attributes an amount of information contribution to the predictors $\{X_1\}$, $\{X_2\}$ and $\{X_1, X_2\}$. The numbers assigned to these sets are nonnegative and, together, they sum up to the mutual information $I(X_1, X_2; Y)$. This is in many ways similar to the partial information decomposition framework, but we note again that the information contribution decomposition has fewer terms than the partial information

decomposition (e.g. three rather than four in the two-input case).

In the example RDN in table 1, the two inputs share a single bit of information about the target. In the PID framework, this typically corresponds to one bit of shared or redundant information. However, the Shapley decomposition does not try to identify redundancy as a separate term, and instead assigns half a bit to each of the predictors. The joint predictor $\{X_1, X_2\}$ is assigned a zero contribution. This reflects the fact that once the correlations between Y and the two individual predictors are known the three-way correlations are already determined, and so learning them does not reveal any extra information.

In the second example, XOR, we have $Y = X_1 \oplus X_2$, where \oplus is the exclusive-or function. In this example, no contribution is assigned to the individual predictors X_1 and X_2 , but one bit is assigned to the joint predictor $\{X_1, X_2\}$. This can be seen as a kind of synergy measure — it says that all of the information that the predictors give about the target is found in the three-way correlations between X_1 , X_2 and Y , and none in the pairwise correlations between either predictor and the target. Interpreting this causally, it means that the causal influences of X_1 and X_2 on Y are strongly tied together. While the Shapley decomposition does not have a term corresponding to redundancy, we see that it characterises synergy in a rather intuitive way.

Our third and fourth examples are discussed in (Harder et al., 2013). The “two bit copy” operation plays an important role in the PID literature in the context of the identity axiom. The Shapley decomposition assigns one bit each to both of the predictors and none to the joint predictor, reflecting the fact that the two inputs each provide a different piece of information about the target. This can be compared to the PID framework, since it is usually seen as desirable for a PID measure to assign zero bits of synergy in this case. Note, however, that because our decomposition does not try to identify redundancy, it does not distinguish between this case and the case of RDN, where the information is also shared equally between the two predictors. The results for the AND distribution are similar, telling us that there is also no synergy in this case. This is because for AND the joint distribution can be inferred completely by knowing the marginals (X_1, Y) , (X_2, Y) and (X_1, X_2) , and consequently there is no triplewise information.

Our final two-predictor example is SYN RDN, which can be formed by combining the XOR example with an independent copy of the RDN example. The values assigned to the two predictors and the joint predictor are simply the sum of their values in the original two examples, which is a result of the additivity property (theorem 1, part V).

Table 2 shows the results for three input variables. In this case the method assigns an amount of information to every non-empty subset of $\{X_1, X_2, X_3\}$, representing the share of the mutual information provided by that set of inputs. The first example, PARITY, is a three-input analog of the XOR example, since $Y = X_1 \oplus X_2 \oplus X_3$. In this example it is not possible to infer anything about the value of Y until the values of all three inputs are known. Correspondingly, the method assigns all of the total mutual information (1 bit) to the predictor $\{X_1, X_2, X_3\}$ and none to the others.

Our second example, XORMULTICOAL (which we take from (Griffith and Koch, 2014)) has the property that knowing any single input gives no information about the target, but any pair of predictors completely determines it. This is reflected in the contributions assigned by the Shapley decomposition:

RDN					
X_1	X_2	Y	p	predictor	contribution (bits)
0	0	0	1/2	$\{X_2\}$	1/2
1	1	1	1/2	$\{X_1\}$	1/2
				$\{X_1, X_2\}$	0

XOR					
X_1	X_2	Y	p	predictor	contribution (bits)
0	0	0	1/4	$\{X_2\}$	0
0	1	1	1/4	$\{X_1\}$	0
1	0	1	1/4	$\{X_1, X_2\}$	1
1	1	0	1/4		

2 BIT COPY					
X_1	X_2	Y	p	predictor	contribution (bits)
0	0	0	1/4	$\{X_2\}$	1
0	1	1	1/4	$\{X_1\}$	1
1	0	2	1/4	$\{X_1, X_2\}$	0
1	1	3	1/4		

AND					
X_1	X_2	Y	p	predictor	contribution (bits)
0	0	0	1/4	$\{X_2\}$	0.40563765
0	1	0	1/4	$\{X_1\}$	0.40563762
1	0	0	1/4	$\{X_1, X_2\}$	0
1	1	1	1/4		

SYNRDN					
X_1	X_2	Y	p	predictor	contribution (bits)
0	0	0	1/8	$\{X_2\}$	1/2
0	1	1	1/8	$\{X_1\}$	1/2
1	0	1	1/8	$\{X_1, X_2\}$	1
1	1	0	1/8		
2	2	2	1/8		
2	3	3	1/8		
3	2	3	1/8		
3	3	2	1/8		

Table 1: Examples of the Shapley information decomposition for several simple two-predictor cases. For each example the joint distribution is shown on the left, and on the right we tabulate $I_{\{X_1\}}$, $I_{\{X_2\}}$ and $I_{\{X_1, X_2\}}$, the contributions made by the two singleton predictors $\{X_1\}$ and $\{X_2\}$ and the joint predictor $\{X_1, X_2\}$. These three values always sum to the mutual information $I(X_1, X_2; Y)$. All logarithms are taken to base 2, so that the numbers are in bits. The interpretation of these examples is given in the text.

PARITY						
X_1	X_2	X_3	Y	p	predictor	contribution (bits)
0	0	0	0	1/8	$\{X_1\}$	0
0	0	1	1	1/8	$\{X_2\}$	0
0	1	0	1	1/8	$\{X_3\}$	0
0	1	1	0	1/8	$\{X_1, X_2\}$	0
1	0	0	1	1/8	$\{X_1, X_3\}$	0
1	0	1	0	1/8	$\{X_2, X_3\}$	0
1	1	0	0	1/8	$\{X_1, X_2, X_3\}$	1
1	1	1	1	1/8	total	1

XORMULTICOAL						
X_1	X_2	X_3	Y	p	predictor	contribution (bits)
4	0	4	0	1/8	$\{X_1\}$	0
0	2	2	0	1/8	$\{X_2\}$	0
1	1	0	0	1/8	$\{X_3\}$	0
5	3	6	0	1/8	$\{X_1, X_2\}$	1/3
5	1	4	1	1/8	$\{X_1, X_3\}$	1/3
1	3	2	1	1/8	$\{X_2, X_3\}$	1/3
0	0	0	1	1/8	$\{X_1, X_2, X_3\}$	0
4	2	6	1	1/8	total	1

RBOJ						
X_1	X_2	X_3	Y	p	predictor	contribution (bits)
0	0	0	0	1/4	$\{X_1\}$	2/3
0	1	1	1	1/4	$\{X_2\}$	2/3
1	0	1	2	1/4	$\{X_3\}$	2/3
1	1	0	3	1/4	$\{X_1, X_2\}$	0
					$\{X_1, X_3\}$	0
					$\{X_2, X_3\}$	0
					$\{X_1, X_2, X_3\}$	0
					total	2

THREE WAY AND						
X_1	X_2	X_3	Y	p	predictor	contribution (bits)
0	0	0	0	1/8	$\{X_1\}$	0.18118725
0	0	1	0	1/8	$\{X_2\}$	0.18118724
0	1	0	0	1/8	$\{X_3\}$	0.18118724
0	1	1	0	1/8	$\{X_1, X_2\}$	0
1	0	0	0	1/8	$\{X_1, X_3\}$	0
1	0	1	0	1/8	$\{X_2, X_3\}$	0
1	1	0	0	1/8	$\{X_1, X_2, X_3\}$	0
1	1	1	1	1/8	total	0.54356444

Table 2: Some examples of our measure, applied to joint distributions between a target and three inputs. The interpretation of these examples is given in the text.

the singleton predictors $\{X_1\}$, $\{X_2\}$ and $\{X_3\}$ each make no contribution to the total. Instead, the total one bit of mutual information is shared equally between the three two-input predictors, $\{X_1, X_2\}$, $\{X_1, X_3\}$ and $\{X_2, X_3\}$. The three-input predictor $\{X_1, X_2, X_3\}$ makes no contribution, because the target is already fully determined by knowing any of the pairwise predictors.

The third example, RBOJ, played an important role in the literature on PID, because it was used in (Rauh et al., 2014) to prove that no partial information decomposition is possible that obeys the so-called identity axiom, along with the axioms of Williams and Beer (2010) and local-positivity. In particular, no such decomposition is possible for this distribution. In this joint distribution, the inputs X_1 , X_2 and X_3 are related by the exclusive-or function, and the target Y is in a one-to-one relationship with its inputs. As a result, each input provides one bit (in the usual sense) of information about the target, and each pair of inputs provides two bits, which completely determine the target. Consequently, learning the third input adds no new information about the target, if the other two are already known. Because our decomposition is different from PID, it is able to assign non-negative values to each of the predictors. It shares out the total two bits of mutual information equally between the three singleton predictors, $\{X_1\}$, $\{X_2\}$ and $\{X_3\}$. This can be seen as a compromise between the fact that the contributions of each member of a pair of input variables are independent (similarly to the 2-bit copy) and that they, at the same time, need to be fairly allocated to three variables.

We finish with an example, THREE WAY AND, in which the decomposition is less intuitive. In this case, the target is 1 if and only if all three inputs are 1. Similarly to the AND example, our measure divides the information contributions between the three singleton predictors, assigning none to the two- or three-input predictors. The reason for this is similar to the AND example. Because of this, from the perspective of our measure, this example looks similar to the RBOJ example.

7 Cooperative game theory and weighted path summation

In section 4 we defined the information contribution $I_A(X_1, \dots, X_n; Y)$ based on a uniform weighting of the maximal chains in the input lattice. In this section we return to the question of how this uniform distribution would be justified.

To do so, we use the notion of the *Shapley value* (Shapley, 1953) from cooperative game theory. Informally, the idea of the Shapley value is that one has a set of players $N = \{A_i, i = 1 \dots |N|\}$. Subsets of the players are called *coalitions*, and each coalition is assigned a *total score* (we will use this term interchangeably with *payoff*), which is to be interpreted as how well that set of players could do at some task, without the participation of the remaining non-coalition players. Given this data, the problem is to assign a score to each individual player, such that the scores of each individual player sum up to the total score. The players' scores should reflect their "fair" contribution in achieving the total score.

For this assignment of scores to be uniquely characterized, Shapley postulates that the scores assigned to players should be a linear function of the coalitions' scores, a notion of relevance (explained below) and a notion of sym-

Shapley Theory	Information in Simplicial Complexes
player A, B, C	set (simplex) A, B, C
coalition	simplicial complex
empty coalition \emptyset	empty ¹ simplicial complex $\mathfrak{S}^{(0)} = \{\}$
coalition of all players \mathfrak{N}	all subsets of $\{1, \dots, n\} = [n]$, i.e. $2^{[n]}$
value of a coalition \mathfrak{S}	$D_{\text{KL}}(p_{\mathfrak{S}} \ p_{\{\emptyset\}})$ (0 if $\mathfrak{S} = \{\}$)
Shapley value $\phi_A(v)$	information contribution $I_A(X_1, \dots, X_n; Y)$

Table 3: Correspondence between our quantities and coalitional game theory. Note that the empty coalition/simplicial complex is included at the bottom of the lattice.

metry amongst the players (where players whose contribution to value cannot be distinguished via a symmetrical exchange of players should attain the same Shapley value). The basic Shapley value assumes that all subsets of N are possible as coalitions. Since Shapley’s original work, many generalizations of the Shapley value have been developed (Bilbao, 1998; Bilbao and Edelman, 2000; Grabisch and Michel, 2009; Faigle and Grabisch, 2012; Ulrich Faigle and Michel Grabisch, 2013).

The purpose of our measure $I_A(X_1, \dots, X_n; Y)$ is to assign to each predictor A a unique share of the mutual information $I(X_1, \dots, X_n; Y)$. Equation (19) calculates this as a linear function of the quantities $D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\{\emptyset\}})$, which can be thought of as the information provided by \mathfrak{S} , which is a set of predictors. This is closely reminiscent to the task of the Shapley value to identify contribution of a particular player when the values of all valid coalitions of players are known.

In fact, we can apply cooperative game theory directly to our problem, by treating sets A of input variables (i.e. predictors) as players in a cooperative game, in which the score of a coalition \mathfrak{S} is given by $D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\{\emptyset\}})$, and hence the total score is $I(X_1, \dots, X_n; Y)$. The only complication is that not every set of players forms a viable coalition, because \mathfrak{S} is constrained to be a simplicial complex. We thus need a formulation which permits us to restrict the possible coalitions to the ones imposed by the simplicial partial order \leq on the set \mathcal{D} of all simplicial complexes. This restriction also necessitates a modification of the symmetry axiom of the original Shapley value to guarantee that the generalized Shapley allocation becomes uniquely determined.

Concretely, here we argue that the specific quantity in Eq. (18) can be interpreted precisely as the generalized Shapley value under *precedence constraints* in the sense of Faigle and Kern (1992).

We will use similar notation for cooperative games to the notation we have used this far for information quantities. We use the symbols A, B, C, \dots for players, and similarly $\mathfrak{S}, \mathfrak{S}', \dots$ for coalitions, \mathcal{D} for the set of all feasible coalitions, to keep a coherent notation. Finally, let \mathfrak{N} denote the set of all players. Table 3 gives the relationship between game-theoretic quantities and the quantities defined in previous sections. To simplify the exposition and render it coherent with respect to existing literature on cooperative game theory, we additionally include the empty coalition (simplicial complex) below the coalition $\{\emptyset\}$ and assign to it the value 0. This will not affect any of the results on the input lattice.

¹In previous chapters, we considered $\{\emptyset\}$ as bottom node of the input lattice. Here, we

In what follows, we introduce Faigle and Kern's extension of the Shapley value, and then show that applying it to this 'information game' is indeed equivalent to eq. (18) with μ taken as the uniform distribution over maximal chains, resulting in eq. (19). This demonstrates that our measure obeys the axioms of Linearity, Carrier and Hierarchical strength, described below, which are used to derive Faigle and Kern's result.

7.1 Shapley Value under Precedence Constraints

We now proceed to define the (generalized) Shapley value under precedence constraints as defined in (Faigle and Kern, 1992). For brevity, when we henceforth say "Shapley value", we will refer to this variant unless stated otherwise.

Let \mathfrak{N} be a finite partially ordered set of players, where for $A, B \in \mathfrak{N}$ the relation $B \leq A$ enforces that, if $A \in \mathfrak{S}$ for any coalition $\mathfrak{S} \subseteq \mathfrak{N}$, one also has $B \in \mathfrak{S}$ (compare with (2)). Under this constraint, not necessarily every subset of \mathfrak{N} is a valid coalition. Let \mathcal{D} be the set of all valid coalitions. \mathcal{D} is closed under intersection and union operations, but not necessarily under the complement operation.

A *cooperative game* on \mathfrak{N} is now a function

$$v : \mathcal{D} \rightarrow \mathbb{R} \quad (23)$$

such that $v(\emptyset) = 0$. Consider the vector space Υ of all cooperative games on \mathfrak{N} , then the Shapley value is defined as a function

$$\Phi : \Upsilon \rightarrow \mathbb{R}^{\mathfrak{N}} \quad (24)$$

which defines, for each player A from \mathfrak{N} , their share $\Phi_A(v)$ for the game v .

Several axioms are postulated for the Shapley value.

Axiom 1 (Linearity). For all $c \in \mathbb{R}, v, w \in \Upsilon$, demand

$$\begin{aligned} \Phi(cv) &= c\Phi(v) \\ \Phi(v+w) &= \Phi(v) + \Phi(w) \end{aligned}$$

Axiom 2 (Carrier). Call a coalition $\mathfrak{U} \in \mathcal{D}$ a *carrier* of $v \in \Upsilon$ if $v(\mathfrak{S}) = v(\mathfrak{S} \cap \mathfrak{U})$ for all $\mathfrak{S} \in \mathcal{D}$. Then, if \mathfrak{U} is a carrier of v , we have

$$\sum_{A \in \mathfrak{U}} \Phi_A(v) = v(\mathfrak{U}) . \quad (25)$$

The carrier axiom needs a brief explanation. It unifies two intuitive axioms that are sometimes used instead, the *dummy axiom* (a player that does not affect the value (or payoff) of any coalition attains a Shapley value of 0) and the *efficiency axiom* (the sum of the Shapley values of all players sums up to the total payoff of the whole set of players).

The third axiom of the traditional Shapley value postulates that players whose contribution to coalition payoffs are equivalent with respect to a symmetric permutation will also receive the same Shapley allocation. This axiom

also include the node below $\{\emptyset\}$, interpreted as coalition, namely, the *empty coalition* in the lattice.

cannot be used in our case, because it requires all subsets of \mathfrak{N} to be feasible coalitions. To obtain a unique characterization of the generalized Shapley value discussed here, a stronger requirement needs to be imposed. There are several axiom sets which are equivalent on the ordered coalition games discussed here (see introduction of section 7 above). We follow (Faigle and Kern, 1992) in choosing the formulation via hierarchical strength.

We need a number of definitions. Call an injective map

$$\pi : \mathfrak{N} \rightarrow \{1, 2, \dots, |\mathfrak{N}|\}$$

a (*feasible*) *ranking* of the players in \mathfrak{N} if for all $A, B \in \mathfrak{N}$ we have that $A < B$ (i.e. $A \leq B$ and $A \neq B$) implies $\pi(A) < \pi(B)$.

The ranking π of \mathfrak{N} induces a ranking $\pi_{\mathfrak{S}} : \mathfrak{S} \rightarrow \{1, 2, \dots, |\mathfrak{S}|\}$ on all coalitions $\mathfrak{S} \in \mathcal{D}$ via $\pi_{\mathfrak{S}}(A) < \pi_{\mathfrak{S}}(B)$ if and only if $\pi(A) < \pi(B)$ for all $A, B \in \mathfrak{S}$. Note that only ordering, but not numbering are equivalent in π and $\pi_{\mathfrak{S}}$.

We say that player $C \in \mathfrak{S}$ is \mathfrak{S} -maximal in the ranking π if $\pi_{\mathfrak{S}}(C) = |\mathfrak{S}|$ which is the same as saying that $\pi(C) = \max_{A \in \mathfrak{S}} \pi(A)$ or that there is no player A in coalition \mathfrak{S} with $C < A$.

We are now ready to express the concept of hierarchical strength: the *hierarchical strength* $h_{\mathfrak{S}}(C)$ of the player C in \mathfrak{S} is defined as the proportion of (total) rankings π in which C is \mathfrak{S} -maximal. Formally,

$$h_{\mathfrak{S}}(C) := \frac{1}{|\mathcal{R}(\mathfrak{N})|} |\{\pi \in \mathcal{R}(\mathfrak{N}) \mid C \text{ is } \mathfrak{S}\text{-maximal for } \pi\}| \quad (26)$$

where $\mathcal{R}(\mathfrak{N})$ is the set of all rankings for the set \mathfrak{N} of players.

Define now a particular fundamental game type, the *inclusion game* over \mathfrak{S} , called $\zeta_{\mathfrak{S}}$ via:

$$\zeta_{\mathfrak{S}}(\mathfrak{T}) := \begin{cases} 1 & \text{if } \mathfrak{S} \subseteq \mathfrak{T} \\ 0 & \text{otherwise.} \end{cases} \quad (27)$$

In other words, the payoff of the game is 1 if the tested coalition \mathfrak{T} encompasses a given reference coalition \mathfrak{S} and vanishes otherwise. We mention without proof that these games form a basis of Υ and thus it is sufficient to define the Shapley value over all inclusion games on \mathfrak{N} . Inclusion games form a preferred set of coalitions in the study of coalitional games since they are closely related to elementary games (games where only a specific coalition achieves a non-zero payoff) and have an intuitive interpretation.

Finally, we can now define

Axiom 3 (Hierarchical Strength (Equivalence)). For any $\mathfrak{S} \in \mathcal{D}$, $A, B \in \mathfrak{S}$, we demand:

$$h_{\mathfrak{S}}(A)\Phi_B(\zeta_{\mathfrak{S}}) = h_{\mathfrak{S}}(B)\Phi_A(\zeta_{\mathfrak{S}}) \quad (28)$$

Informally, the Shapley value of a player B in a coalition \mathfrak{S} for the inclusion game is weighted against that of another player A in the same coalition via their hierarchical strength. All else being equal, the Shapley values of the two players relate to each other as their hierarchical strengths — a larger value of the hierarchical strength corresponds to a larger Shapley value, i.e. larger allocation of payoff.

Faigle and Kern (1992) note that the hierarchical strength emphasizes the given player being *on top* of its respective coalition in the given ranking rather than, say, considering its average rank. This is insofar an intuitive choice for the generalized Shapley value, since it is only the top-ranked player in a coalition which determines whether that particular coalition is formed at all. In other words, it is a measuring in how many rankings (relative to the total number of rankings) that particular player has the power to decide whether the given coalition will be formed or not.

It turns out this has a straightforward reinterpretation and generalization in the context of Markovian coalition processes (Ulrich Faigle and Michel Grabisch, 2012). In addition, there are many other formulations equivalent with it (see references mentioned in the introduction to the present chapter section 7). We opted for the formulation via hierarchical strength since it is the most widely established generalization of the symmetry axiom for the classical Shapley value in the literature.

We state here without proof that the unique payoff allocation is given by

$$\Phi_C(v) = \frac{1}{|\mathcal{R}(\mathfrak{N})|} \sum_{\substack{\mathfrak{T} \in \mathcal{D} \\ C \in \mathfrak{T}^+}} |\mathcal{R}(\mathfrak{T} \setminus \{C\})| \cdot |\mathcal{R}(\mathfrak{N} \setminus \mathfrak{T})| (v(\mathfrak{T}) - v(\mathfrak{T} \setminus \{C\})) \quad (29)$$

where we trivially assume $|\mathcal{R}(\emptyset)| = 1$ and where $C \in \mathfrak{T}^+$ sums over all coalitions \mathfrak{T} for which C is maximal. The (generalized) Shapley value of C is given by the marginal contribution of C to all coalitions \mathfrak{T} for which it is maximal, weighted by the proportion of rankings for which this is the case.

We now show that our definition of information contribution of a simplex (eq. (19)) is equivalent to the generalized Shapley value under precedence constraints if the value is the mutual information between input variables X_1, \dots, X_n and output variable Y . Thus, the information contribution has a natural interpretation in the context of game-theoretic payoff allocation.

7.2 Equivalence of Generalized Shapley Value and the Sum over Maximal Chains

Theorem 2. We now prove that, under the identifications of table 3, the information contribution of a set A is identical with its Shapley value under precedence constraints, with A interpreted as a player. More precisely:

$$\begin{aligned} \sum_{(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A} \mu(\mathfrak{S}, \mathfrak{S}') D_{\text{KL}}(p_{\mathfrak{S}'} \| p_{\mathfrak{S}}) = \\ \sum_{\substack{\mathfrak{S} \in \mathcal{D} \\ A \in \mathfrak{S}^+}} \frac{|\mathcal{R}(\mathfrak{S} \setminus \{A\})| \cdot |\mathcal{R}(2^{[n]} \setminus \mathfrak{S})|}{|\mathcal{R}(2^{[n]})|} [D_{\text{KL}}(p_{\mathfrak{S}} \| p_{\{\emptyset\}}) - D_{\text{KL}}(p_{\mathfrak{S} \setminus \{A\}} \| p_{\{\emptyset\}})] , \quad (30) \end{aligned}$$

where the weighting of the lattice chains μ is chosen as the uniform distribution, $\mu(\gamma) = 1/|\Gamma|$. To permit consistency between lattice and Shapley model, we furthermore define the bracketed term on the right side to be 0 for $\mathfrak{S} = \emptyset$.

Proof. Consider $\mathfrak{N} = 2^{[n]}$. Identify the elements $A \in \mathfrak{N}$, i.e. the subsets of $[n]$, with the players in a Shapley coalition game with partial ordering defined via

the subset relation, i.e. via

$$B \leq A : \Longleftrightarrow B \subseteq A.$$

Per definition, the partial order-compatible coalitions \mathfrak{S} then precisely constitute the simplicial complexes of \mathfrak{N} .

We now show that, under these identifications, the (feasible) rankings of players define precisely the maximal chains over simplicial complexes. In other words, there is a one-to-one correspondence between the rankings of the ordered coalition game and the maximal chains over its corresponding simplicial complexes.

The intuition of the proof is as follows: in the Hasse diagram for the input lattice section 4, each maximal chain is formed by successively adding each of the predictors, one at a time, in such a way that each step of the chain remains a simplicial complex. We will demonstrate that the orders in which the predictors are added in a maximal chain correspond precisely to the feasible rankings of the predictors interpreted as Shapley players. We now proceed to show this formally.

We first show well-definedness, i.e. that each ranking defines a maximal chain. Let π , a (feasible) ranking over the set of players \mathfrak{N} , be given (we remind that each player is a subset of $[n]$). Define the sequence

$$\mathfrak{S}_0 := \emptyset, \mathfrak{S}_1, \mathfrak{S}_2, \dots, \mathfrak{S}_{|\mathfrak{N}|} \quad (31)$$

where for $k = 1, \dots, |\mathfrak{N}|$

$$\mathfrak{S}_k := \{A \in \mathfrak{N} \mid \pi(A) \leq k\} \quad (32)$$

$$= \pi^{-1}(\{1, \dots, k\}). \quad (33)$$

We now need to show now that this sequence $(\mathfrak{S}_k)_{k=0, \dots, |\mathfrak{N}|}$ is, first, a chain of simplicial complexes (equivalently, feasible coalitions) and, second, maximal.

If $k = 0$, then $\mathfrak{S}_k = \emptyset$ is trivially a simplicial complex. Else, let $1 \leq k \leq |\mathfrak{N}|$. Consider now $A \in \mathfrak{S}_k$, and any $B \in \mathfrak{N}$ with $B \subseteq A$. We have $\pi(B) \leq \pi(A) \in \{1, \dots, k\}$ per ranking property, and thus $\pi(B) \in \{1, \dots, k\}$, and thus $B \in \mathfrak{S}_k$ and \mathfrak{S}_k is a simplicial complex.

From (33) it follows that, for $k \leq l$, $\mathfrak{S}_k \subseteq \mathfrak{S}_l$. Therefore, if $A \in \mathfrak{S}_k$, also $A \in \mathfrak{S}_l$ and thus $\mathfrak{S}_k \leq \mathfrak{S}_l$ and the $(\mathfrak{S}_k)_k$ form a chain.

This chain is maximal. To show this, consider successive simplicial complexes $\mathfrak{S}_k, \mathfrak{S}_{k+1}$, $k = 0, \dots, |\mathfrak{N}| - 1$ in the sequence. Consider $\tilde{\mathfrak{S}}$ such that $\mathfrak{S}_k \leq \tilde{\mathfrak{S}} \leq \mathfrak{S}_{k+1}$ according to the natural partial order \leq on simplicial complexes. If $\mathfrak{S}_k \neq \tilde{\mathfrak{S}}$, then there exists a $B \in \tilde{\mathfrak{S}} \setminus \mathfrak{S}_k$ and, since $\tilde{\mathfrak{S}} \leq \mathfrak{S}_{k+1}$, one has $B \subseteq C$ for some $C \in \mathfrak{S}_{k+1}$. This means that $\pi(B) \leq \pi(C)$. Since $B \notin \mathfrak{S}_k$, also $\pi(B) \notin \{1, \dots, k\}$, so, per construction of \mathfrak{S}_{k+1} , necessarily $\pi(B) = k + 1$ and $B = \pi^{-1}(k + 1) = C \in \mathfrak{S}_{k+1}$. It follows that $\tilde{\mathfrak{S}}$ must be either \mathfrak{S}_k or \mathfrak{S}_{k+1} , thus, $\mathfrak{S}_k \prec \mathfrak{S}_{k+1}$ and the chain is maximal. This shows that the mapping from rankings to maximal chains is well-defined.

We show now that mapping rankings to maximal chains (31) via (33) is injective. For this, consider two rankings $\pi \neq \rho$. We have to show that they induce different maximal chains.

Consider B with $\pi(B) \neq \rho(B)$. Assume, without loss of generality, $\pi(B) < \rho(B)$. If we consider the chain $(\mathfrak{S}_k^\pi)_k$ induced by π (and analogously $(\mathfrak{S}_k^\rho)_k$ for ρ), then observe that the chain can be written in the form of inclusion chain as

$$\emptyset \subseteq \mathfrak{S}_0^\pi \subseteq \mathfrak{S}_1^\pi \subseteq \cdots \subseteq \mathfrak{S}_{\pi(B)}^\pi \subseteq \cdots \subseteq \mathfrak{S}_{|\mathfrak{N}|}^\pi = \mathfrak{N}. \quad (34)$$

\uparrow first time where B appears in $(\mathfrak{S}_k^\pi)_k$

In this chain, the first simplicial complex to contain B is the one with index $\pi(B)$. Under the same consideration for the chain induced by ρ , the first member of the chain to contain B is the one with index $\rho(B)$. However, $\pi(B) < \rho(B)$ and therefore the chains must differ and assigning chains to rankings via (31) is injective.

Show now surjectivity: for each maximal chain, there is a ranking that produces it. Let

$$\emptyset = \mathfrak{S}_0 \subseteq \mathfrak{S}_1 \subseteq \cdots \subseteq \mathfrak{S}_{|\mathfrak{N}|} = \mathfrak{N} \quad (35)$$

be a maximal chain. We show now that each step adds exactly one $C \in \mathfrak{N}$. Assume none of the steps in the sequence is trivial, i.e. we always have $\mathfrak{S}_j \subsetneq \mathfrak{S}_{j+1}$. All \mathfrak{S}_k are at the same time simplicial complexes as well as — equivalently — order-compatible coalitions. Choose $C \in \mathfrak{S}_{j+1} \setminus \mathfrak{S}_j$ minimal (i.e. such that for any $B \in \mathfrak{S}_{j+1} \setminus \mathfrak{S}_j$ with $B \subseteq C$, we have $B = C$).

Since $C \in \mathfrak{S}_{j+1}$, for any $B \subseteq C$, we have $B \in \mathfrak{S}_{j+1}$. It follows that either $B \in \mathfrak{S}_j$ or $B \in \mathfrak{S}_{j+1} \setminus \mathfrak{S}_j$; in the latter case, however, because of minimality of C in $\mathfrak{S}_{j+1} \setminus \mathfrak{S}_j$, it follows $B = C$. Thus $\mathfrak{S}_j \cup \{C\}$ is a simplicial complex, and because of maximality of the chain, it must be identical to \mathfrak{S}_{j+1} . In summary, in each step of the maximal chain precisely one simplicial complex is added.

Finally, given a maximal chain

$$\emptyset \prec \mathfrak{S}_1 \prec \mathfrak{S}_2 \prec \cdots \prec \mathfrak{S}_{|\mathfrak{N}|} = \mathfrak{N}, \quad (36)$$

define for every $j = 1, \dots, |\mathfrak{N}|$ the inverse ranking $\pi^{-1}(j)$ to map onto the unique set (player) in $\mathfrak{S}_j \setminus \mathfrak{S}_{j-1}$. The maximal chain (36) is induced by the ranking π ; we have thus shown the mapping (33) of rankings to maximal chains to be surjective (for every maximal chain there is a ranking that is mapped to it). With the injectivity shown earlier, this mapping is thus bijective. In short, we have shown that to each maximal chain corresponds one and only one feasible ranking.

Consider now $(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A$, i.e. an edge where $\mathfrak{S}' = \mathfrak{S} \cup \{A\}$ is obtained by adding A to \mathfrak{S} . For the set $\Gamma(\mathfrak{S}, \mathfrak{S}')$ of (maximal) chains $(\mathfrak{S}_k)_k$ who pass through this edge, i.e. for which $\mathfrak{S}_j = \mathfrak{S}$ and $\mathfrak{S}_{j+1} = \mathfrak{S}'$ for some j , one has, in analogy to the derivation above, a one-to-one map to the pairs of the rankings over $\mathfrak{S}' \setminus \{A\}$ and those over $\mathfrak{N} \setminus \mathfrak{S}'$:

$$\mathcal{R}(\mathfrak{S}' \setminus \{A\}) \times \mathcal{R}(\mathfrak{N} \setminus \mathfrak{S}'). \quad (37)$$

This is seen by replacing the full ranking with two subrankings, one over the lower sublattice with \mathfrak{S} as top element instead of \mathfrak{N} and one over the upper one which has \mathfrak{S}' as bottom element replacing \mathfrak{S}_0 . It follows that we have

$$|\Gamma(\mathfrak{S}, \mathfrak{S}')| = |\mathcal{R}(\mathfrak{S}' \setminus \{A\})| \cdot |\mathcal{R}(\mathfrak{N} \setminus \mathfrak{S}')|. \quad (38)$$

Consider a particular edge $(\mathfrak{S}, \mathfrak{S}') \in \mathcal{E}_A$. We note that this edge corresponds precisely to the simplicial complexes $\mathfrak{S}' \in \mathcal{D}$ where A is maximal in \mathfrak{S}' , i.e. $A \in \mathfrak{S}'^+$. We had earlier the short-hand notation $\mu(\mathfrak{S}, \mathfrak{S}') = \sum_{\gamma \in \Gamma(\mathfrak{S}, \mathfrak{S}')} \mu(\gamma)$ where $\Gamma(\mathfrak{S}, \mathfrak{S}')$ ranges over all chains containing a particular edge. If all chains/paths γ are equally weighted, their weight is given by

$$\frac{1}{|\Gamma|} = \frac{1}{|\mathcal{R}(\mathfrak{N})|} = \frac{1}{|\mathcal{R}(2^{[n]})|} \quad (39)$$

Finally, note that

$$D_{\text{KL}}(p_{\mathfrak{S}'} \| p_{\mathfrak{S}}) = D_{\text{KL}}(p_{\mathfrak{S}'} \| p_{\mathfrak{S}^{(0)}}) - D_{\text{KL}}(p_{\mathfrak{S}' \setminus \{A\}} \| p_{\mathfrak{S}^{(0)}}) \quad (40)$$

because of the Pythagorean relation (9). This completes the proof of (30). \square

Note that, when constructing the correspondence between the input lattice to the Shapley value, for the former we had the maximal chains start at $\{\emptyset\}$ rather than at \emptyset as bottom of the lattice. However, the property from theorem 2 continues to hold in this case, since the bottom step from \emptyset to $\{\emptyset\}$ is unique and does not affect the path counts.

8 Discussion

In the search for a partial information measure that allocates informational contributions to various input variable sets (i.e. predictors) we relinquished the demand to quantify redundancy and instead applied the Pythagorean decomposition to characterize the additional contribution of an input variable set as it is added on the relevant maximal chains. This “longitudinal” contribution is chain-dependent, though. To be able to talk about a contribution of an individual predictor, though, we need to express this contribution independently of the particular chain.

Intuitively, this can be done by assigning a probability distribution over the chains and averaging a predictor’s contribution over all these chains; most naturally, the equidistribution could be chosen for this purpose. A more justified reasoning for this choice can be derived by observing that the setup of information contribution precisely matches the situation of a coalition game where the value of a coalition is the contribution of that coalition to the overall “value”, i.e. information about the target variable; and that contribution can be fairly assigned via the Shapley value concept. Of course, with the natural precedence order of predictors, not all coalitions of predictors (i.e. players in the language of game theory) are viable. We needed to resort to the variant of the Shapley value under *precedence constraints* which, as it turns out, corresponds precisely to the averaging over all maximal chains of the input lattice, strengthening both the confidence in the appropriateness of the measure and the intuition behind it.

While the view of a predictor contribution stemming from averaging over chains (paths) through the lattice seems abstract and artificial, the Shapley value-based interpretation justifies its use. In fact, this perspective finds, again, additional justification from more recent coalition game theory in which coalitions are not considered as immutable, but can change as per a stochastic process

via local incentives (Ulrich Faigle and Michel Grabisch, 2012). In our context, this would correspond to a dynamically chosen path in an input lattice. At this stage, however, we are interested in the static contributions of the predictors; whether there will be an incentive to invoke a complex trajectory in the input lattice over which the contributions will be averaged, remains a question for the future.

Acknowledgement

DP would like to acknowledge support by H2020-641321 socSMCs FET Proactive project. NA and NV acknowledge the support of the Deutsche Forschungsgemeinschaft Priority Programme “The Active Self” (SPP 2134).

References

- Amari, S.-i. (2001). Information geometry on hierarchical decomposition of stochastic interactions. *IEEE Transaction on Information Theory*, 47:1701–1711.
- Amari, S.-i. and Nagaoka, H. (2007). *Methods of information geometry*, volume 191. American Mathematical Soc.
- Amari, S.-i., Tsuchiya, N., and Oizumi, M. (2016). Geometry of information integration. In *Information Geometry and its Applications IV*, pages 3–17. Springer.
- Ay, N. (2001/2015). Information geometry on complexity and stochastic interaction. *Entropy*, 17(4):2432–2458. Originally published in 2001 as MiS-Preprint 95/2001. Journal version published 2015. Preprint URL <https://www.mis.mpg.de/publications/preprints/2001/prepr2001-95.html>.
- Ay, N., Jost, J., Vân Lê, H., and Schwachhöfer, L. (2017). *Information Geometry*. Springer.
- Ay, N. and Polani, D. (2008). Information flows in causal networks. *Advances in complex systems*, 11(01):17–41.
- Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., and Ay, N. (2014). Quantifying unique information. *Entropy*, 16(4):2161–2183.
- Bilbao, J. M. (1998). Axioms for the Shapley value on convex geometries. *European Journal of Operational Research*, 110(2):368–376.
- Bilbao, J. M. and Edelman, P. H. (2000). The Shapley value on convex geometries. *Discrete Applied Mathematics*, 103(1-3):33–40.
- Csiszár, I. and Matúš, F. (2003). Information projections revisited. *IEEE Transactions on Information Theory*, 49(6):1474–1490.
- Csiszár, I. and Matúš, F. (2004). On information closures of exponential families: counterexample. *IEEE Transactions on Information Theory*, 50(5):922–924.

- Csiszár, I. and Shields, P. C. (2004). Information theory and statistics: A tutorial. *Foundations and Trends in Communications and Information Theory*, 1(4):417–528.
- Faigle, U. and Grabisch, M. (2012). Values for Markovian coalition processes. *Economic Theory*, 51(3):505–538.
- Faigle, U. and Kern, W. (1992). The Shapley value for cooperative games under precedence constraints. *International Journal of Game Theory*, 21(3):249–266.
- Finn, C. and Lizier, J. (2018). Pointwise partial information decomposition using the specificity and ambiguity lattices. *Entropy*, 20(4):297.
- Grabisch, F. L. and Michel (2009). Values on regular games under Kirchhoff’s laws. *Mathematical Social Sciences*, 58:322–340.
- Griffith, V. and Koch, C. (2014). Quantifying synergistic mutual information. In *Guided Self-Organization: Inception*, pages 159–190. Springer.
- Harder, M., Salge, C., and Polani, D. (2013). Bivariate measure of redundant information. *Phys. Rev. E*, 87:012130. <http://arxiv.org/abs/1207.2080>.
- Ince, R. A. (2017). The partial entropy decomposition: Decomposing multivariate entropy and mutual information via pointwise common surprisal. *arXiv preprint arXiv:1702.01591*.
- James, R. and Crutchfield, J. (2017). Multivariate dependence beyond shannon information. *Entropy*, 19(10):531.
- James, R. G., Emenheiser, J., and Crutchfield, J. P. (2019). Unique information via dependency constraints. *Journal of Physics A: Mathematical and Theoretical*, 52(1):014002.
- Kolchinsky, A. (2019). A novel approach to multivariate redundancy and synergy. *arXiv preprint arXiv:1908.08642*.
- Lauritzen, S. L. (1996). *Graphical Models*. Oxford Science Publications.
- Lizier, J. T. (2014). Jidt: An information-theoretic toolkit for studying the dynamics of complex systems. *Frontiers in Robotics and AI*, 1:11.
- Lizier, J. T., Prokopenko, M., and Zomaya, A. Y. (2014). A framework for the local information dynamics of distributed computation in complex systems. In *Guided self-organization: inception*, pages 115–158. Springer.
- Oizumi, M., Tsuchiya, N., and Amari, S.-i. (2016). Unified framework for information integration based on information geometry. *Proceedings of the National Academy of Sciences*, 113(51):14817–14822.
- Olbrich, E., Bertschinger, N., and Rauh, J. (2015). Information decomposition and synergy. *Entropy*, 17(5):3501–3517.
- Pearl, J. (2009). *Causality*. Cambridge university press.

- Perrone, P. and Ay, N. (2016). Hierarchical quantification of synergy in channels. *Frontiers in Robotics and AI*, 2:35.
- Rauh, J., Bertschinger, N., Olbrich, E., and Jost, J. (2014). Reconsidering unique information: Towards a multivariate information decomposition. In *2014 IEEE International Symposium on Information Theory*, pages 2232–2236. IEEE.
- Rosas, F., Ntranos, V., Ellison, C., Pollin, S., and Verhelst, M. (2016). Understanding interdependency through complex information sharing. *Entropy*, 18(2):38.
- Schreiber, T. (2000). Measuring information transfer. *Phys. Rev. Lett.*, 85:461–464.
- Shapley, L. S. (1953). A Value for n -person Games. In Kuhn, H. and Tucker, A., editors, *Annals of Mathematics Studies*, volume 28, pages 307–317. Princeton University Press.
- Ulrich Faigle and Michel Grabisch (2012). Values for Markovian coalition processes. *Economic Theory*, 51(3):505–538.
- Ulrich Faigle and Michel Grabisch (2013). A Concise Axiomatization of a Shapley-type Value for Stochastic Coalition Processes. *Economic Theory Bulletin*, pages 189–199.
- Williams, P. L. and Beer, R. D. (2010). Nonnegative decomposition of multivariate information. *arXiv preprint arXiv:1004.2515*.
- Zwisk, M. (2004). An overview of reconstructability analysis. *Kybernetes*, 33(5/6):877–905.