

# Hierarchical Quantification of Synergy in Channels

Paolo Perrone  
Nihat Ay

SFI WORKING PAPER: 2015-12-050

SFI Working Papers contain accounts of scientific work of the author(s) and do not necessarily represent the views of the Santa Fe Institute. We accept papers intended for publication in peer-reviewed journals or proceedings volumes, but not papers that have already appeared in print. Except for papers by our external faculty, papers must be based on work done at SFI, inspired by an invited visit to or collaboration at SFI, or funded by an SFI grant.

©NOTICE: This working paper is included by permission of the contributing author(s) as a means to ensure timely distribution of the scholarly and technical work on a non-commercial basis. Copyright and all rights therein are maintained by the author(s). It is understood that all persons copying this information will adhere to the terms and constraints invoked by each author's copyright. These works may be reposted only with the explicit permission of the copyright holder.

[www.santafe.edu](http://www.santafe.edu)



SANTA FE INSTITUTE

# HIERARCHICAL QUANTIFICATION OF SYNERGY IN CHANNELS

Paolo Perrone<sup>\*1</sup> and Nihat Ay<sup>1,2,3</sup>

<sup>1</sup>Max Planck Institute for Mathematics in the Sciences, Inselstrasse 22, 04103  
Leipzig, Germany

<sup>2</sup>Faculty of Mathematics and Computer Science, University of Leipzig, PF 100920,  
04109 Leipzig, Germany

<sup>3</sup>Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

## Abstract

The decomposition of channel information into synergies of different order is an open, active problem in the theory of complex systems. Most approaches to the problem are based on information theory, and propose decompositions of mutual information between inputs and outputs in several ways, none of which is generally accepted yet.

We propose a new point of view on the topic. We model a multi-input channel as a Markov kernel. We can project the channel onto a series of exponential families which form a hierarchical structure. This is carried out with tools from information geometry, in a way analogous to the projections of probability distributions introduced by Amari. A Pythagorean relation leads naturally to a decomposition of the mutual information between inputs and outputs into terms which represent single node information; pairwise interactions; and in general  $n$ -node interactions.

The synergy measures introduced in this paper can be easily evaluated by an iterative scaling algorithm, which is a standard procedure in information geometry.

**Keywords:** Synergy, Redundancy, Hierarchy, Projections, Divergences, Interactions, Iterative Scaling, Information Geometry.

## 1 Introduction

In complex systems like biological networks, for example neural networks, a basic principle is that their functioning is based on the correlation and interaction of their different parts. While correlation between two sources is well understood, and can be quantified by Shannon's mutual information (see for example [12]),

---

<sup>\*</sup>Correspondence: perrone@mis.mpg.de

there is still no generally accepted theory for interactions of three nodes or more. If we label one of the nodes as “output”, the problem is equivalent to determine how much two (or more) input nodes interact to yield the output. This concept is known in common language as “synergy”, which means “working together”, or performing a task that would not be feasible by the single parts separately.

There are a number of important works which address the topic, but the problem is still considered open. The first generalization of mutual information was *interaction information* (introduced in [1]), defined for three nodes in terms of the joint and marginal entropies:

$$I(X : Y : Z) = -H(X, Y, Z) + H(X, Y) + H(X, Z) + H(Y, Z) - H(X) - H(Y) - H(Z). \quad (1)$$

Interaction information is defined symmetrically on the joint distribution, but most approaches interpret it by looking at a channel, rather than a joint distribution,  $(X, Y) \rightarrow Z$ . For example, we can rewrite (1) equivalently in terms of mutual information (choosing  $Z$  as “output”):

$$I(X : Y : Z) = I(X, Y : Z) - I(X : Z) - I(Y : Z), \quad (2)$$

where we see that it can mean intuitively “how much the whole  $(X, Y)$  gives more (or less) information about  $Z$  than the sum of the parts separately”. Another expression, again equivalent, is:

$$I(X : Y : Z) = I(X : Y|Z) - I(X : Y), \quad (3)$$

which we can interpret as “how much conditioning over  $Z$  changes the correlation between  $X$  and  $Y$ ” (see [2]). Unlike mutual information, interaction information carries a sign:

- $I > 0$ : *synergy*. Conditioning on one node *increases* the correlation between the remaining nodes. Or, the whole gives more information than the sum of the parts. Example: XOR function.
- $I < 0$ : *redundancy*. Conditioning on one node *decreases*, or *explains away* the correlation between the remaining nodes. Or, the whole gives less information than the sum of the parts. Example:  $X = Y = Z$ .
- $I = 0$ : *3-independence*. Conditioning on one node has no effect on the correlation between the remaining nodes. Or, the whole gives the same amount of information as the parts separately. The nodes can nevertheless still be conditionally dependent. Example: independent nodes.<sup>1</sup>

As argued in [3], [4], and [5], however, this is not the whole picture. There are systems which exhibit both synergetic and redundant behavior, and interaction information only quantifies the *difference* of synergy and redundancy, with a priori no way to tell the two apart. In a system with highly correlated

---

<sup>1</sup>For an example in which  $I = 0$  but the nodes are not independent, see [4].

inputs, for example, the synergy would remain unseen, as it would be cancelled by the redundancy. Moreover, this picture breaks down for more than three nodes. Another problem, pointed out in [3] and [8], is that redundancy (as for example in  $X = Y = Z$ ) can be described in terms of pairwise interactions, not triple, while synergy (as in the XOR function) is purely threewise. Therefore,  $I$  compares and mixes information quantities of different nature.

A detailed explanation of the problem for two inputs is presented in [4] and it yields a decomposition (“Partial Information Decomposition, PID) as follows: there exist two non-negative quantities, *Synergy* and *Redundancy*, such that

$$I(X, Y : Z) = I(X : Z) + I(Y : Z) + Syn - Red , \quad (4)$$

or equivalently:

$$I(X : Y : Z) = Syn - Red . \quad (5)$$

Moreover, they define *unique information* for the inputs  $X$  and  $X_2$  as:

$$UI(X) = I(X : Z) - Red , \quad (6)$$

$$UI(Y) = I(Y : Z) - Red , \quad (7)$$

so that the total mutual information is decomposed positively:

$$I(X, Y : Z) = UI(X) + UI(Y) + Red + Syn . \quad (8)$$

What these quantities intuitively mean is:

- Redundancy – information available in both inputs;
- Unique information – information available only in one of the inputs;
- Synergy – information available only when both inputs are present, arising purely from their interaction.

In this formulation, if one finds a measure of synergy, one can automatically define compatible measures of redundancy and unique information (and viceversa), provided that the measure of synergy is always larger or equal to  $I(X : Y : Z)$ , and that the resulting measure of redundancy is less or equal than  $I(X : Z)$  and  $I(Y : Z)$ . Synergy, redundancy, and unique information are defined on a channel, and choosing a different channel with the same joint distribution (e.g.  $(Y, Z) \rightarrow X$ ) may yield a different decomposition.

In [5] is presented an overview of (previous) measures of synergy, and their shortcomings in standard examples. In the same paper is then presented a newer measure for synergy, defined equivalently in [6] as:

$$CI(X, Y; Z) := I(X, Y : Z) - \min_{p^* \in \wedge} I_{p^*}(X, Y : Z) , \quad (9)$$

where  $\wedge$  is the space of distributions with prescribed marginals:

$$\wedge = \{q \in P(X, Y, Z) \mid q(X, Z) = p(X, Z), q(Y, Z) = p(Y, Z)\} . \quad (10)$$

This measure satisfies interesting properties (proven in [5] and [6]), which make it compatible with Williams and Beer’s PID, and with the intuition in most examples. However, it was proven in [7] that such an approach can *not* work in the desired way for more than three nodes (two inputs).

Our approach uses information geometry [13], extending previous work on hierarchical decompositions [8] and complexity [9]. (Compare the related approach on information decomposition pursued in [10].) The main tools of the present paper are KL-projections, and the Pythagorean relation that they satisfy. This allows (as in [8]) to form hierarchies of interactions of different orders in a geometrical way. In the present problem, we decompose mutual information between inputs and outputs of a channel  $k$ , for two inputs, as:

$$I(X, Y : Z) = d_1(k) + d_2(k) , \quad (11)$$

where  $d_2$  quantifies synergy (as in equation (8)), and  $d_1$  integrates all the lower order terms (*UI, Red*), quantifying the so-called *union information* (see [5]). One may want to use this measure of synergy to form a complete decomposition analogous to (8), but this does not work, as in general it is not true that  $d_2 \leq I(X : Y : Z)$ . For this reason, we keep the decomposition more coarse, and we do not divide union information into unique and redundant.

For more inputs  $X_1, \dots, X_N$ , the decomposition generalizes to:

$$I(X_1, \dots, X_N : Z) = d_1(k) + \dots + d_N(k) = \sum_{i=1}^N d_i(k) , \quad (12)$$

where higher orders of synergy appear.

Until now, there seems to be no way of rewriting the decomposition of [5] and [6] in a way consistent with information geometry, and more in general, Williams and Beer’s PID seems hard to write as a geometric decomposition. A comparison between  $d_2$  and the measure  $CI$  of [5] and [6] is presented in Section 5. There we show that  $d_2 \leq CI$ , and we argue, with a numerical example, that  $CI$  overestimates synergy at least in one case.

For a small number of inputs ( $\lesssim 5$ ), our quantities are easily computable with the standard algorithms of information geometry (like iterative scaling [11]). This allowed to get precise quantities for all the examples considered.

## 1.1 Technical Definitions

We consider a set of  $N$  input nodes  $V = \{1, \dots, N\}$ , taking values in the sets  $X_1, \dots, X_N$ , and an output node, taking values in the set  $Y$ . We write the input globally as  $X := X_1 \times \dots \times X_N$ . For example, in biology  $Y$  can be the phenotype, and  $X$  can be a collection of genes determining  $Y$ . We denote by  $F(Y)$  the set of real functions on  $Y$ , and with  $P(X)$  the set of probability measures on  $X$ .

We can model the channel from  $X$  to  $Y$  as a Markov kernel (called also stochastic kernel, transition kernel, or stochastic map)  $k$ , that assigns to each  $x \in X$  a probability measure on  $Y$  (for a detailed treatment, see [12]). Here

we will consider only finite systems, so we can think of a channel simply as a transition matrix (or stochastic matrix), whose rows sum to one.

$$k(x; y) \geq 0 \quad \forall x, y; \quad \sum_y k(x; y) = 1 \quad \forall x. \quad (13)$$

The space of channels from  $X$  to  $Y$  will be denoted by  $K(X; Y)$ . We will denote by  $X$  and  $Y$  also the corresponding random variables, whenever this does not lead to confusion.

Conditional probabilities define channels: if  $p(X, Y) \in P(X, Y)$  and the marginal  $p(X)$  is strictly positive, then  $p(Y|X) \in K(X; Y)$  is a well-defined channel. Viceversa, if  $k \in K(X; Y)$ , given  $p \in P(X)$  we can form a well-defined joint probability:

$$pk(x, y) := p(x) k(x; y) \quad \forall x, y. \quad (14)$$

An “input distribution”  $p \in P(X)$  is crucial also to extend the notion of divergence from probability distributions to channels. The most natural way of doing it is the following.

**Definition 1.** Let  $p \in P(X)$ , let  $k, m \in K(X; Y)$ . Then:

$$D_p(k||m) := \sum_{x,y} p(x) k(x; y) \log \frac{k(x; y)}{m(x; y)}. \quad (15)$$

Defined this way,  $D_p$  is affine in  $p$ . It is worth noticing that  $D_p$  is in general *not* equal to  $D(k_*p||m_*p)$ . Moreover, it has an important compatibility property. Let  $p, q$  be joint probability distributions on  $X \times Y$ , and let  $D$  be the KL-divergence. Then:

$$D(p(X, Y)||q(X, Y)) = D(p(X)||q(X)) + D_{p(X)}(p(Y|X)||q(Y|X)). \quad (16)$$

We will now illustrate our geometric ideas in channels with one, two, and three input nodes, then we present some examples. The general case will be addressed in Section 4.

## 2 Geometric Idea of Synergy

**Mutual information as motivation.** It is a well-known fact in information theory that Shannon’s mutual information can be written as a KL-divergence:

$$I_p(X : Y) = D(p(X, Y)||p(X)p(Y)). \quad (17)$$

From the point of view of information geometry, this can be interpreted as a “distance” between the real distribution and a product distribution that has exactly the same marginals, but maximal entropy. In other words, we have:

$$I_p(X : Y) = \inf_{\substack{q \in P(X) \\ r \in P(Y)}} D(p(X, Y)||q(X)r(Y)). \quad (18)$$

The distribution given by  $p(X)p(Y)$  is optimal in the sense that:

$$p(X)p(Y) = \arg \min_{\substack{q \in P(X) \\ r \in P(Y)}} D(p(X, Y) || q(X)r(Y)) . \quad (19)$$

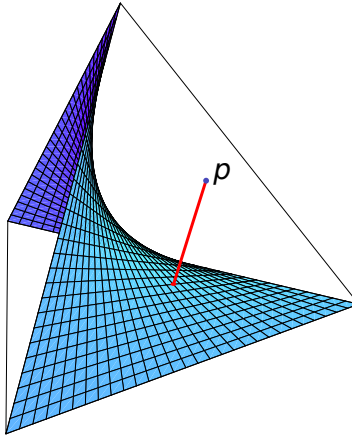


Figure 1: For two binary nodes, the family of product distributions is a surface in a 3-dimensional simplex.

The divergence between  $p$  and a submanifold is, as usual in geometry, the “distance” between  $p$  and the “closest point” on that submanifold, which in our case is the geodesic projection w.r.t. the mixture connection.

**Extension to channels.** We can use the same insight with channels. Instead of a joint distribution on  $N$  nodes, we consider a channel from an input  $X$  to an output  $Y$ . Suppose we have a family  $\mathcal{E}$  of channels, and a channel  $k$  that may not be in  $\mathcal{E}$ . Then, just as in geometry, we can define the “distance” between  $k$  and  $\mathcal{E}$ .

**Definition 2.** Let  $p$  be an input distribution. The divergence between a channel  $k$  and a family of channels  $\mathcal{E}$  is given by:

$$D_p(k || \mathcal{E}) := \inf_{m \in \mathcal{E}} D_p(k || m) . \quad (20)$$

If the minimum is uniquely realized, we call the channel

$$\pi_{\mathcal{E}} k := \arg \min_{m \in \mathcal{E}} D_p(k || m) \quad (21)$$

the *KL-projection* of  $k$  on  $\mathcal{E}$  (and simply “a” KL-projection if it is not unique).

We will always work with compact families, so the minima will always be realized, and for strictly positive  $p$  they will be unique (see Section 4 for the details).

We will consider families  $\mathcal{E}$  for which the KL-divergence satisfies a Pythagorean equality (see Figure 2 below for some intuition):

$$D_p(k||m) = D_p(k||\pi_{\mathcal{E}}k) + D_p(\pi_{\mathcal{E}}k||m) \quad (22)$$

for every  $m \in \mathcal{E}$ . These families (technically, closures of exponential families) are defined in Section 4.

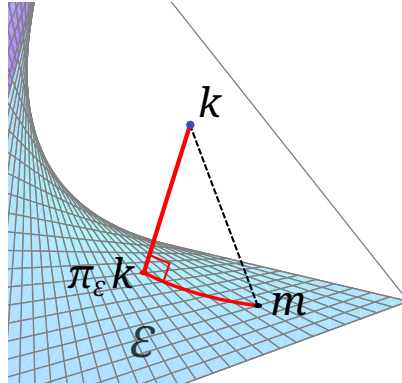


Figure 2: Illustration of the Pythagoras theorem for projections

**One input.** Consider first one input node  $X$ , with input distribution  $p(X)$ , and one output node  $Y$ . A *constant channel*  $k$  in  $K(X; Y)$  is a channel whose entries do not depend on  $X$  (more precisely:  $k(x; y) = k(x'; y)$  for any  $x, x', y$ ). This denomination is motivated by the following properties:

- They correspond to channels that do not use any information from the input to generate the output.
- The output distribution given by  $k$  is a probability distribution on  $Y$  which does not depend on  $X$ .
- Deterministic constant channels are precisely constant functions.

We call  $\mathcal{E}_0$  the family of constant channels. Take now any channel  $k \in K(X; Y)$ . If we want to quantify the dependence in  $k$  of  $Y$  on  $X$  we can then look at the divergence of  $k$  from the constant channels:

$$d_1(k) := D_p(k||\mathcal{E}_0) . \quad (23)$$

The minimum is realized in  $\pi_{\mathcal{E}_0}k$ . We have that:

$$d_1(k) = D_p(k||\pi_{\mathcal{E}_0}k) = \sum_{x,y} p(x) k(x; y) \log \frac{k(x; y)}{\pi_{\mathcal{E}_0}k(y)} \quad (24)$$

$$= H_{p\pi_{\mathcal{E}_0}k}(Y) - H_{pk}(Y|X) = I_{pk}(X : Y) , \quad (25)$$

so that consistently with our intuition, the dependence of  $Y$  on  $X$  is just the mutual information. From the channel point of view, it is simply the divergence from the constant channels. (A rigorous calculation is done in Section 4.)



**Two inputs.** Consider now two input nodes with input probability  $p$  and one output node. We can again define the family  $\mathcal{E}_0$  of constant channels, and the same calculations give:

$$D_p(k|\mathcal{E}_0) = I_{pk}(X_1, X_2 : Y) . \quad (26)$$

This time, though, we can say a lot more: the quantity above can be decomposed. In analogy with the independence definition for probability distributions, we would like to define a split channel as a product channel of its parts:  $p(y|x_1, x_2) = p(y|x_1)p(y|x_2)$ . Unfortunately, the term on the right would be in general not normalized, so we replace the condition by a weaker one. We call the channel  $k(X_1, X_2; Y)$  *split* if it can be written as:

$$k(x_1, x_2; y) = \phi_0(x_1, x_2) \phi_1(x_1; y) \phi_2(x_2; y) \quad (27)$$

for some *functions*  $\phi_0, \phi_1, \phi_2$ , which in general are not themselves channels (in particular,  $\phi_i(x_i; y) \neq p(y|x_i)$ ). We call  $\mathcal{E}_1$  the family of split channels. This family corresponds to those channels that do not have any synergy. This is a special case of an exponential family, analogous to the family of product distributions of Figure 1. The examples “single node” and “split channel” in the next section belong exactly to this family. Take now any channel  $k(X_1, X_2; Y)$ . In analogy with mutual information, we call *synergy* the divergence:

$$d_2(k) := D_p(k|\mathcal{E}_1) . \quad (28)$$

Simply speaking, our synergy is quantified as the deviation of the channel from the set  $\mathcal{E}_1$  of channels without synergy.

We can now project  $k$  first to  $\mathcal{E}_1$ , and then to  $\mathcal{E}_0$ . Since  $\mathcal{E}_0$  is a subfamily of  $\mathcal{E}_1$ , the following Pythagoras relation holds from (22):

$$D_p(k|\pi_{\mathcal{E}_0}k) = D_p(k|\pi_{\mathcal{E}_1}k) + D_p(\pi_{\mathcal{E}_1}k|\pi_{\mathcal{E}_0}k) . \quad (29)$$

If in analogy with the one-input case we call the last quantity  $d_1$ , we get from (26) and (28):

$$I_{pk}(X_1, X_2 : Y) = d_2(k) + d_1(k) . \quad (30)$$

The term  $d_1$  measures how much information comes from single nodes (but it does not tell which nodes). The term  $d_2$  measures how much information comes from the synergy of  $X_1$  and  $X_2$  in the channel. The example “XOR” in the next section will show this.

If we call  $\mathcal{E}_2$  the whole  $K(X; Y)$ , we get  $\mathcal{E}_0 \subset \mathcal{E}_1 \subset \mathcal{E}_2$  and:

$$d_i(k) := D_p(\pi_{\mathcal{E}_i}k|\pi_{\mathcal{E}_{i-1}}k) . \quad (31)$$

**Three inputs.** Consider now three nodes  $X_1, X_2, X_3$  with input probability  $p$ , and a channel  $k$ . We have again:

$$D_p(k|\mathcal{E}_0) = I_{pk}(X_1, X_2, X_3 : Y) . \quad (32)$$

This time we can decompose the mutual information in different ways. We can for example look at split channels, i.e. in the form:

$$k(x_1, x_2, x_3; y) = \phi_0(x) \phi_1(x_1; y) \phi_2(x_2; y) \phi_3(x_3; y) \quad (33)$$

for some  $\phi_0, \phi_1, \phi_2, \phi_3$ . As in the previous case, we call this family  $\mathcal{E}_1$ . Or we can look at more interesting channels, the ones in the form:

$$k(x_1, x_2, x_3; y) = \phi_0(x) \phi_{12}(x_1, x_2; y) \phi_{13}(x_1, x_3; y) \phi_{23}(x_2, x_3; y) \quad (34)$$

for some  $\phi_0, \phi_{12}, \phi_{13}, \phi_{23}$ . We call this family  $\mathcal{E}_2$ , and it is easy to see that:

$$\mathcal{E}_0 \subset \mathcal{E}_1 \subset \mathcal{E}_2 \subset \mathcal{E}_3, \quad (35)$$

where  $\mathcal{E}_0$  denotes again the constant channels, and  $\mathcal{E}_3$  denotes the whole  $K(X; Y)$ . We define again:

$$d_i(k) := D_p(\pi_{\mathcal{E}_i} k | | \pi_{\mathcal{E}_{i-1}} k). \quad (36)$$

This time, the Pythagorean relation can be nested, and it gives us:

$$I_{pk}(X_1, X_2, X_3 : Y) = d_3(k) + d_2(k) + d_1(k), \quad (37)$$

The difference between pairwise synergy and threewise synergy is shown in the “XOR” example in the next section.

Now that we have introduced the measure for a small number of input, we can study the examples from the literature [5], and show that our measure is consistent with the intuition. The general case will be more in rigor introduced in Section 4.

### 3 Examples

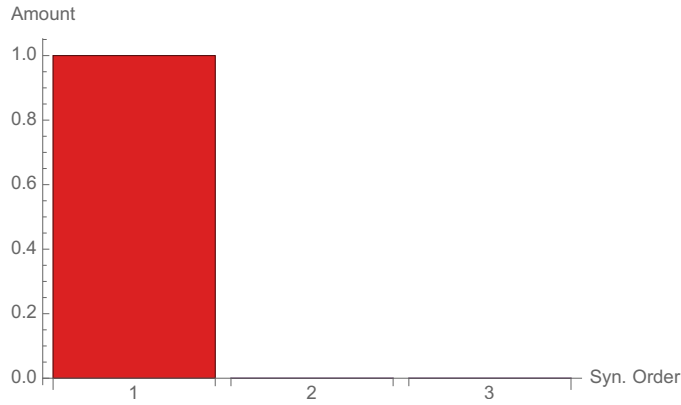
Here we present some examples of decomposition for well-known channels. All the quantities have been computed using an algorithm analogous to iterative scaling (as in [11]).

**Single Node Channel.** The easiest example is considering a channel which only depends on  $X_1$ , i.e.:

$$I(X : Y) = I(X_1 : Y). \quad (38)$$

For example, consider 3 binary input nodes  $X_1, X_2, X_3$  with constant input probability, and one binary output node  $Y$  which is an exact copy of  $X_1$ .

Then we have exactly one bit of single node information, and no higher order terms. Geometrically,  $k$  lies in  $\mathcal{E}_1$ , so the only nonzero divergence in equation (37) is  $d_1(k)$ . As one would expect,  $d_2(k)$  and  $d_3(k)$  vanish, as there is no synergy of order 2 and 3.



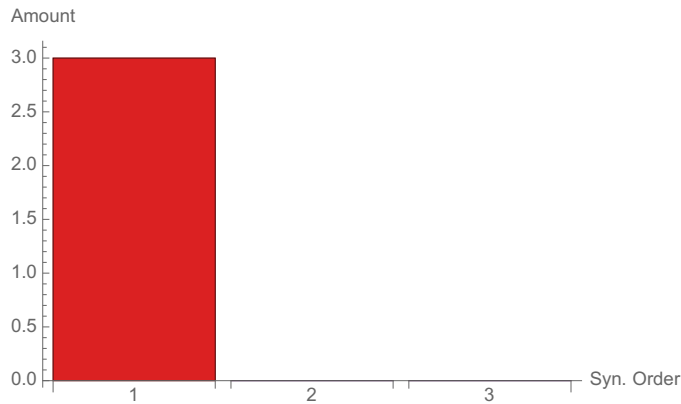
**Split Channel.** The second easiest example is a more general channel which obeys equation (33). In particular, consider 3 binary input nodes  $X_1, X_2, X_3$  with constant input probability (so, the  $x_i$  are independent), and output  $Y = X_1 \times X_2 \times X_3$ . As channel we simply take the identity map  $(x_1, x_2, x_3) \mapsto (x_1, x_2, x_3) \in Y$ . In this particular case:

$$I(X : Y) = \sum_i I(X_i; Y) . \quad (39)$$

We have 3 bits of mutual information, which are all single node (but from different nodes). Since:

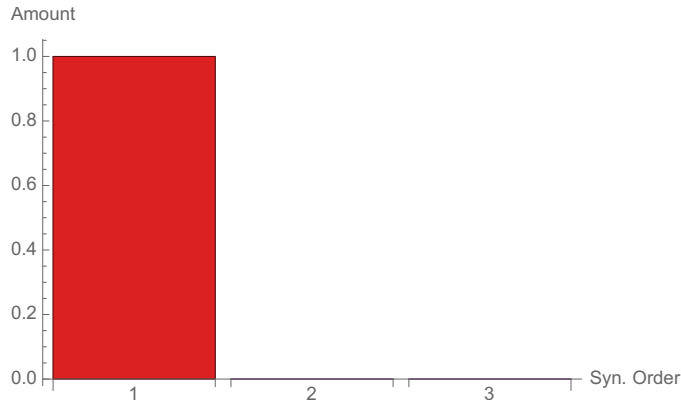
$$k(x_1, x_2, x_3; y) = \phi_1(x_1, y_1) \phi_2(x_2, y_2) \phi_3(x_3, y_3) , \quad (40)$$

which is a special case of (33),  $k \in \mathcal{E}_1$ , and so  $d_2(k)$  and  $d_3(k)$  in equation (37) are again zero.



**Correlated Inputs.** Consider 3 perfectly correlated binary nodes, each one with uniform marginal probability. As output take a perfect copy of one (hence,

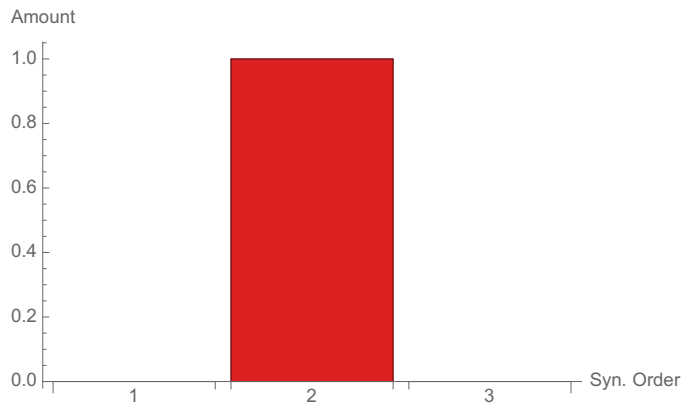
all) of the inputs. We have again one bit of mutual information, which could come from any of the nodes, but no synergy, as no two nodes are interacting in the channel. The input distribution has correlation, but this has no effect on the channel, since the channel is simply copying the value of  $X_1$  (or  $X_2$  or  $X_3$ , equivalently). Therefore again  $k \in \mathcal{E}_1$ . Of the terms in equation (37), again the only non-zero is  $d_1(k)$ .



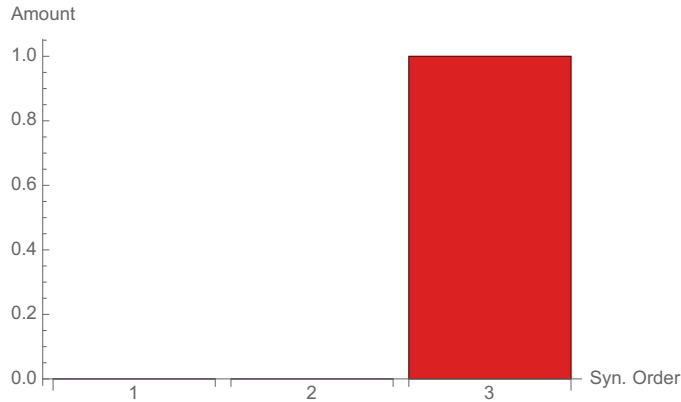
This example in the literature is used to motivate the notion of redundancy. A “redundant channel” is in our decomposition exactly equivalent to a single node channel, since it contains exactly the same amount of information.

**Parity (XOR).** The standard example of synergy is given by the XOR function, and more generally by the parity function between two or more nodes.

For example, consider 3 binary input nodes  $X_1, X_2, X_3$  with constant input probability, and one binary output node  $Y$  which is given by  $X_1 \vee X_2$ . We have 1 bit of mutual information, which is purely arising from a pairwise synergy (of  $X_1$  and  $X_2$ ), so this time  $k \in \mathcal{E}_2$ . The function XOR is *pure* synergy, so  $d_2(k)$  is the only non-zero term in (37).

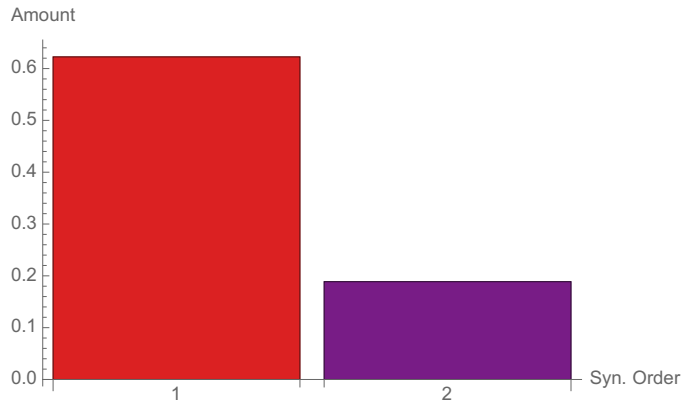


If instead  $Y$  is given by the threewise parity function, or  $X_1 \vee X_2 \vee X_3$ , we have again 1 bit of mutual information, which now is purely arising from a threewise synergy, so here  $k \in \mathcal{E}_3$ , and the only non-zero term in (37) is  $d_3(k)$ .



In these examples there are no terms of lower order synergy, but the generic elements of  $\mathcal{E}_2$  and  $\mathcal{E}_3$  usually do contain a nonzero lower part. Consider for instance the next example.

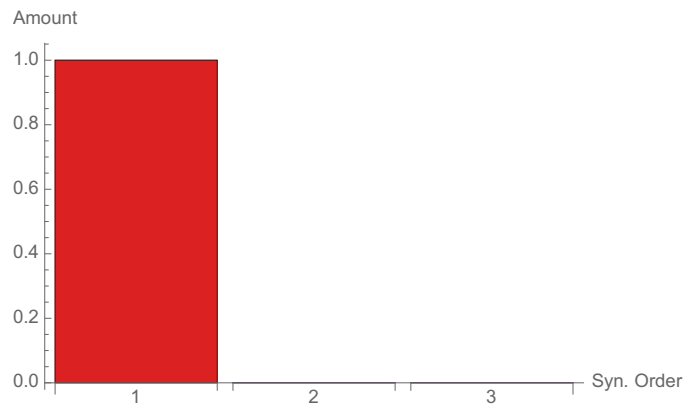
**AND and OR.** The other two standard logic gates, AND and OR, share the same decomposition. Consider two binary nodes  $X_1, X_2$  with uniform probability, and let  $Y$  be  $X_1 \vee X_2$  (or  $X_1 \wedge X_2$ ). There is again one bit of mutual information, which comes mostly from single nodes, but also from synergy.



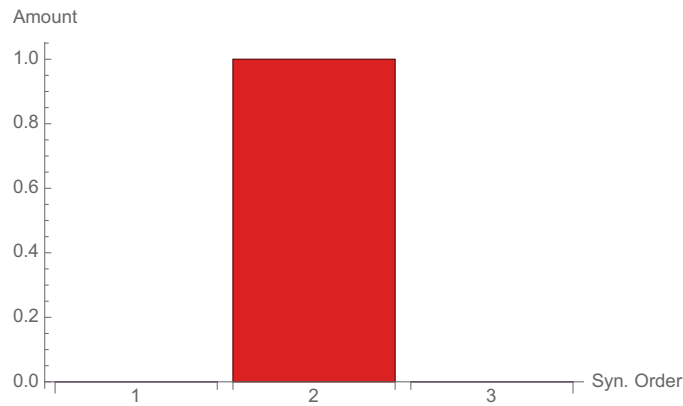
Geometrically, this means that AND and OR are channels in  $\mathcal{E}_2$  which lie close to the submanifold  $\mathcal{E}_1$ .

**XorLoses.** Here we present a slightly more complicated example, coming from [5]. We have three binary nodes  $X_1, X_2, X_3$ , where  $X_1, X_2$  have uniform probabilities, and an output node  $Y = X_1 \vee X_2$ , just like in the “XOR” example.

Now we take  $X_3$  to be perfectly correlated with  $Y = X_1 \vee X_2$ , so that  $Y$  could get the information either from  $X_3$  or from the synergy between  $X_1$  and  $X_2$ . We have one bit of mutual information, which can be seen as entirely coming from  $X_3$ , and so the synergy between  $X_1$  and  $X_2$  is not adding anything.



**XorDuplicate.** Again from [5]. We have 3 binary nodes  $X_1, X_2, X_3$ , where  $X_1, X_2$  have uniform probabilities, while  $X_3 = X_1$ . The output is  $X_1 \vee X_2 = X_3 \vee X_2$ , so it could get the information either from the synergy between  $X_1$  and  $X_2$ , or  $X_2$  and  $X_3$ . There is one bit of mutual information, which is coming from a pairwise interaction. Again, it does not matter between whom.



It should be clear from the examples here that decomposing only *by order*, and not by the specific subsets, is crucial. For example, in the “input correlation” example, there is no natural way to decide from *which single node* the information comes, even if it is clear that the interaction is of order 1.

## 4 General Case

Here we try to give a general formulation, for  $N$  inputs, of the quantities defined in Section 2. As in the introduction we call the set of input nodes  $V$  of cardinality  $N$ , and we consider a subset of the nodes  $I$ . We denote the joint random variable  $(X_i, i \in I)$  by  $X_I$ , and we denote the complement of  $I$  in  $V$  by  $I^c$ . The case  $N = 3$  in Section 2 should motivate the following definition.

**Definition 3.** Let  $I \subseteq V$ . We call  $F_I$  the space of functions who only depend on  $X_I$  and  $Y$ :

$$F_I := \{f \in F(X, Y) \mid f(x_I, x_{I^c}; y) = f(x_I, x'_{I^c}; y) \forall x_{I^c}, x'_{I^c}\}. \quad (41)$$

Let  $0 \leq i \leq N$ . We call  $F_i$  the space of channels which can be written as a product of functions of  $F_I$  with the order of  $I$  at most  $k$ :

$$\mathcal{E}_i := cl \left\{ k \in K(X; Y) \mid \exists \phi_I \in F_I, \phi_0 \in F(X) \mid k = \phi_0 \prod_I \phi_I; |I| \leq i \right\}, \quad (42)$$

where  $cl$  denotes the closure in  $K(X; Y)$ . Intuitively, this means that  $\mathcal{E}_i$  does not only contain terms in the form given in the curly brackets, but also limits of such terms. Stated differently, the closure of a set includes not only the set itself, but also its boundary. This is important, because when we project to a family, the projection may lie on the boundary. In order for the result to exist, the boundary must then be included.

This way:

- $\mathcal{E}_0$  is the space of constant channels;
- $\mathcal{E}_N$  is the whole  $K(X; Y)$ ;
- $\mathcal{E}_i \subseteq \mathcal{E}_j$  if and only if  $i \leq j$ ;
- For  $N \leq 3$  we recover exactly the quantities of Section 2.

The family  $\mathcal{E}_i$  is also the closure of the family in the form:

$$\left\{ \frac{1}{Z(X)} \exp \left( \sum_I f_I(X; Y) \right) \mid f_I \in F_I; |I| \leq i \right\}, \quad (43)$$

where:

$$Z(x) := \sum_y \exp \left( \sum_I f_I(x; y) \right). \quad (44)$$

Such families are known in the literature as *exponential families* (see for example [13]). In particular, it is compact (for finite  $N$ ), so that the infimum of any function on  $\mathcal{E}_i$  is always a minimum. This means that for a channel  $k$  and an input distribution  $p$ :

$$D_p(k \parallel \mathcal{E}_i) := \inf_{m \in \mathcal{E}_i} D_p(k \parallel m) = \min_{m \in \mathcal{E}_i} D_p(k \parallel m) \quad (45)$$

always exists. If it is unique, for example if  $p$  is strictly positive, we define the unique KL-projection as:

$$\pi_{\mathcal{E}_i} k := \arg \min_{m \in \mathcal{E}_i} D_p(k||m). \quad (46)$$

$\pi_{\mathcal{E}_i} k$  has the property that it defines the same output probability on  $Y$ .

**Definition 4.** Let  $k \in K(X; Y)$ , let  $1 \leq i \leq N$ . Then the  $i$ -wise synergy of  $k$  is (if the KL-projections are unique):

$$d_i(k) := D_p(\pi_{\mathcal{E}_i} k || \pi_{\mathcal{E}_{i-1}} k). \quad (47)$$

For more clarity, we call the 1-wise synergy “single node information” or “single-node dependence”.

For  $k \in K(X; Y) = \mathcal{E}_N$ , we can look at its divergence from  $\mathcal{E}_0$ . If we denote  $\pi_{\mathcal{E}_0} k$  by  $k_0$ :

$$D_p(k || \mathcal{E}_0) = D_p(k || k_0) = \sum_{x,y} p(x) k(x; y) \log \frac{k(x; y)}{k_0(y)}. \quad (48)$$

If  $k$  is not strictly positive, we take the convention  $0 \log(0/0) = 0$ , and we discard the zero terms from the sum. Since  $k_{0*} p = k_* p$  but  $k_0$  is constant in  $x$ , it can *not* happen that for some  $(x; y)$ ,  $k_0(x; y) = 0$  but  $k(x; y) \neq 0$ . (The very same is true for all KL-projections  $\pi_{\mathcal{E}_i} k$ , since  $D_p(\pi_{\mathcal{E}_i} k || 0) \leq D_p(k || k_0)$ .) For all other terms, (48) becomes:

$$D_p(k || \mathcal{E}_0) = \sum_{x,y} p(x) k(x; y) \log k(x; y) - \sum_{x,y} p(x) k(x; y) \log k_0(y) \quad (49)$$

$$= -H_{pk}(Y|X) - \sum_y k_* p(y) \log k_0(y) \quad (50)$$

$$= -H_{pk}(Y|X) - \sum_y k_{0*} p(y) \log k_0(y) \quad (51)$$

$$= -H_{pk}(Y|X) + H_{k_{0*} p}(Y) = -H_{pk}(Y|X) + H_{k_* p}(Y) \quad (52)$$

$$= I_{pk}(X : Y). \quad (53)$$

On the other hand, the Pythagorean relation (22) implies:

$$D_p(k || k_0) = D_p(k || \pi_{\mathcal{E}_{N-1}} k) + D_p(\pi_{\mathcal{E}_{N-1}} k || k_0), \quad (54)$$

and iterating:

$$D_p(k || k_0) = D_p(k || \pi_{\mathcal{E}_{N-1}} k) + D_p(\pi_{\mathcal{E}_{N-1}} k || \pi_{\mathcal{E}_{N-2}} k) + \dots + D_p(\pi_{\mathcal{E}_1} k || k_0). \quad (55)$$

In the end, we get:

$$I(X : Y) = \sum_{i=1}^N D_p(\pi_{\mathcal{E}_i} k || \pi_{\mathcal{E}_{i-1}} k) = \sum_{i=1}^N d_i(k). \quad (56)$$

This decomposition is always non-negative, and it depends on the input distribution. The terms in (56) can be in general difficult to compute exactly. Nevertheless, they can be approximated with iterative procedures.



## 5 Comparison with the Measure of [5] and [6]

The measure of synergy, or respectively complementary information, defined in [5] and [6], is:

$$CI_p(Y : X_1, X_2) := I_p(Y : X_1, X_2) - \min_{p^* \in \wedge} I_{p^*}(Y : X_1, X_2), \quad (57)$$

where  $\wedge$  is the space of prescribed marginals:

$$\wedge = \{q \in P(X_1, X_2, Y) \mid q(X_1, Y) = p(X_1, Y), q(X_2, Y) = p(X_2, Y)\}. \quad (58)$$

Our measure of synergy can be written, for two inputs, in a similar form:

$$d_2(k) = D_p(k \mid \pi_{\mathcal{E}_1} k) = I_p(Y : X_1, X_2) - \min_{p^* \in \Delta} I_{p^*}(Y : X_1, X_2), \quad (59)$$

where  $\Delta$ , in addition to the constraints of  $\wedge$ , prescribes also the input:

$$\Delta = \{q \in P(X_1, X_2, Y) \mid q(X_1, Y) = p(X_1, Y), q(X_2, Y) = p(X_2, Y), q(X_1, X_2) = p(X_1, X_2)\}. \quad (60)$$

Clearly  $\Delta \subseteq \wedge$ , so:

$$\min_{p^* \in \Delta} I_{p^*}(Y : X_1, X_2) \geq \min_{p^* \in \wedge} I_{p^*}(Y : X_1, X_2), \quad (61)$$

which implies that:

$$d_2(k) \leq CI_p(Y : X_1, X_2). \quad (62)$$

We argue that not prescribing the input leads to overestimating synergy, because the subtraction in (57) includes a possible difference in the correlation of the input distributions.

For example, consider  $X_1, X_2, Y$  binary and correlated, but not *perfectly* correlated. (For perfectly correlated nodes, as in Section 3,  $\Delta = \wedge$ , so there is no difference between the two measures.) In detail, consider the channel:

$$k_\beta(x_1, x_2; y) := \frac{\exp(\beta y (x_1 + x_2))}{\sum_{y'} \exp(\beta y' (x_1 + x_2))}, \quad (63)$$

and the input distribution:

$$p_\alpha(x_1, x_2) := \frac{\exp(\alpha x_1 x_2)}{\sum_{x'_1, x'_2} \exp(\alpha x'_1 x'_2)}. \quad (64)$$

For  $\alpha, \beta \rightarrow \infty$ , the correlation becomes perfect, and the two measures of synergy are both zero. For  $0 < \alpha, \beta < \infty$ , our measure  $d_2(k_\beta)$  is zero, as clearly  $k_\beta \in \mathcal{E}_1$ .  $CI$  is more difficult to compute, but we can give a (non-zero) lower bound in the following way. First we fix two values  $\beta = \beta_0, \alpha = \alpha_0$ . We consider the joint distribution  $p_{\alpha_0} k_{\beta_0}$ , and look at the marginals:

$$p_{\alpha_0} k_{\beta_0}(X_1, Y), \quad p_{\alpha_0} k_{\beta_0}(X_2, Y). \quad (65)$$

We define the family  $\wedge$  as the set of joint probabilities which have exactly these marginals. If we increase  $\beta$ , we can always find an  $\alpha$  such that the marginals do not change:

$$p_{\alpha}k_{\beta}(X_1, Y) = p_{\alpha_0}k_{\beta_0}(X_1, Y), \quad p_{\alpha}k_{\beta}(X_2, Y) = p_{\alpha_0}k_{\beta_0}(X_2, Y), \quad (66)$$

i.e. such that  $p_{\alpha}k_{\beta} \in \wedge$ . Now we can look at the mutual information of  $p_{\alpha}k_{\beta}$  and of  $p_{\alpha_0}k_{\beta_0}$ . If they differ, and (for example) the former is larger, then:

$$\begin{aligned} I_{p_{\alpha}k_{\beta}}(Y : X_1, X_2) - I_{p_{\alpha_0}k_{\beta_0}}(Y : X_1, X_2) \\ \leq I_{p_{\alpha}k_{\beta}}(Y : X_1, X_2) - \min_{p^* \in \wedge} I_{p^*}(Y : X_1, X_2) = CI_{p_{\alpha}k_{\beta}} \end{aligned} \quad (67)$$

is a well-defined lower bound for  $CI_{p_{\alpha}k_{\beta}}$ . With a numerical simulation we can show graphically that the mutual information is indeed not constant within the families  $\wedge$ .

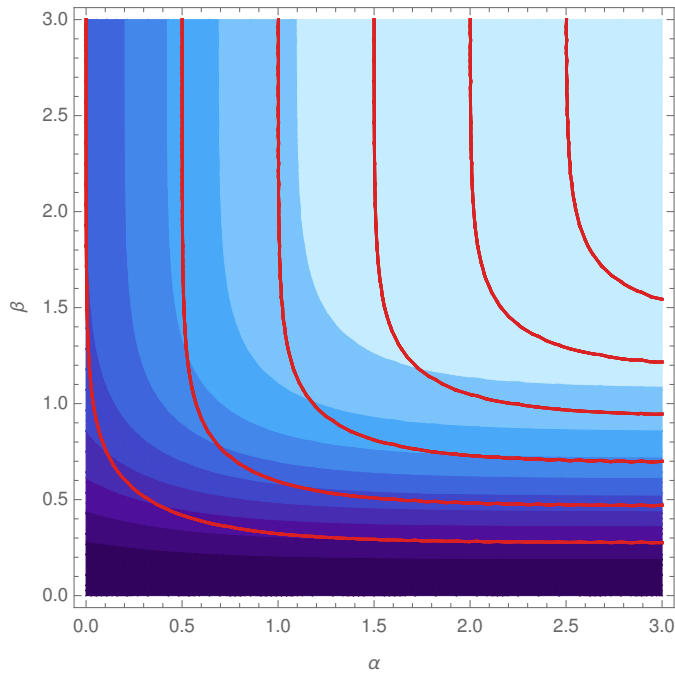


Figure 3: Mutual information and fixed marginals. The shades of blue represent the amount of  $I_p(Y : X_1, X_2)$  as a function of  $\alpha, \beta$  (brighter is higher). Each red line represents a family  $\wedge$  of fixed marginals. While the lines of fixed mutual information and the families of fixed marginals look qualitatively similar, they do not coincide exactly, which means that  $I_p$  varies within the  $\wedge$ .

From the picture we can see that the red lines (families  $\wedge$  for different initial values) approximate well the lines of constant mutual information, at least

qualitatively, but they are not exactly equal. This means that for most points  $p$  of  $\Lambda$ , the quantity:

$$CI_p(Y : X_1, X_2) := I_p(Y : X_1, X_2) - \min_{p^* \in \Lambda} I_{p^*}(Y : X_1, X_2) \quad (68)$$

will be non-zero. More explicitly, we can plot the increase in mutual information as  $p$  varies in  $\Lambda$ , for example as a function of  $\beta$ . This is always larger or equal than the difference between the mutual information and its minimum in  $\Lambda$  (i.e.  $CI$ ). We can see that it is positive, which implies that  $CI_p$  is also positive.

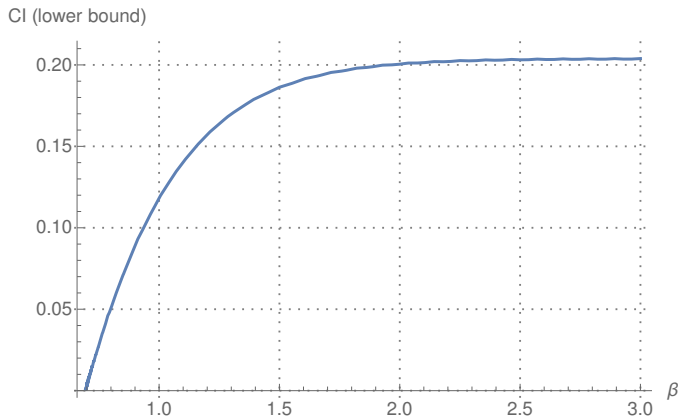


Figure 4: Lower bound for  $CI$  versus  $\beta$ . For each  $\beta \in [0.7, 3]$  we can find an  $\alpha$  such that the joint  $p_{\alpha k_{\beta}}$  lies in  $\Lambda$ . The increase in mutual information as  $\beta$  varies is a lower bound for  $CI$ , which is then in general non-zero.

We can see in Figure 3 that, especially for very large or very small values of  $\alpha$  and  $\beta$  (i.e. very strong or very weak correlation),  $CI$  captures the behaviour of mutual information quite well. These limits are precisely deterministic and constant kernels, for which most approaches in quantifying synergy coincide. This is the reason why the examples studied in [5] give quite a satisfying result for  $CI$  (in their notation,  $S_{VK}$ ). For the less studied (and computationally more complex) intermediate values, like  $1 < \alpha, \beta < 2$ , the approximation is instead far from accurate, and in that interval (see Figure 4) there is a sharp increase in  $I$ , which leads to overestimating synergy.

## 6 Conclusion

Using information geometry, we have defined a non-negative decomposition of the mutual information between inputs and outputs of a channel.

The decomposition divides the mutual information into contributions of the different orders of synergy in a unique way. It does *not*, however, divide the

mutual information into contributions of the different subsets of input nodes as Williams and Beer’s PID [4] would require.

For two inputs, our measure of synergy is closely related to the well-received quantification of synergy in [5] and [6]. Our measure though works in the desired way for an arbitrary (finite) number of inputs. Differently from [5] and [6], anyway, we do not define a measure for redundant or “shared” information, nor unique information of the single inputs or subsets.

The decomposition depends on the choice of an input distribution. In case of input correlation, redundant information is counted automatically only once. This way there is no need to quantify redundancy separately.

In general there is no way to compute our quantities in closed form, but they can be approximated by an iterative scaling algorithm (see for example [11]). The results are consistent with the intuitive properties of synergy, outlined in [4] and [5].

## References

- [1] McGill, W. L. *Multivariate information transmission*. Psychometrika, 19(2):97–116, 1954.
- [2] Schneidmann, E., Bialek, W., and Berry II, M. J. *Synergy, redundancy, and independence in population codes*. The Journal of Neuroscience, 23(37):11539–11553, 2003.
- [3] Schneidmann, E., Still, S., Berry II, M. J., and Bialek, W., *Network information and connected correlations*. Physical Review Letters, 91(23):238701–238704, 2003.
- [4] Williams, P. L. and Beer, R. D. *Nonnegative decomposition of multivariate information*. Preprint available on arXiv:1004.2151, 2010.
- [5] Griffith, V. and Koch, C. *Quantifying synergistic mutual information*. Preprint available on arXiv:1205.4265, 2014.
- [6] Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., and Ay, N. *Quantifying unique information*. Entropy, 16(4):2161, 2014.
- [7] Rauh, J., Bertschinger, N., Olbrich, E., and Jost, J. *Reconsidering unique information: towards a multivariate information decomposition*. Preprint available on arXiv:1404.3146, 2015.
- [8] Amari, S. *Information geometry on a hierarchy of probability distributions*. IEEE Transactions on information Theory, 47(5):1701–1709, 2001.
- [9] Ay, N. *Information geometry on complexity and stochastic interaction*. Entropy, 17(4):2432, 2015.
- [10] Harder, M., Salge, C., and Polani, D. *Bivariate measure of redundant information*. Phys. Rev. E, 87: 012130, 2013.

- [11] Csiszár, I. and Shields, P. C. *Information Theory and Statistics: A Tutorial*. Foundations and Trends in Communications and Information Theory, 1(4):417–528, 2004
- [12] Kakihara, Y. *Abstract Methods in Information Theory*. World Scientific, 1999.
- [13] Amari, S. and Nagaoka, H. *Methods of Information Geometry*. Oxford, 1993.