

Information and Complexity, Or: Where Is the Information?



Nihat Ay, Nils Bertschinger, Jürgen Jost, Ekehard Olbrich,
and Johannes Rauh

Abstract We demonstrate how a systematic theory of complexity emerges from information theoretical concepts. The complexity of a structure may refer to the difficulty of its description, the encoding of its regularities or the relations between its elements, components or parts. All such measures can be and usually are quantified with the help of information theoretical concepts. We first describe those concepts and then use them to analyze how complexity emerges from interactions between parts or conversely, can be decomposed into individual, joint, or complementary contributions of those parts. We can also use these theoretical concepts to understand the interaction between a system and its environment and the relations between different levels of aggregation in complex systems.

Keywords Information theory · Mutual information · Complexity measures · Information decomposition · Levels and scales

1 Introduction

The concepts of information and complexity seem to be intricately linked. Complexity notions are quantified in information theoretical terms, and a general principle might say that a structure is the more complex, the more information is needed to

N. Ay · J. Jost (✉) · E. Olbrich · J. Rauh
Max Planck Institute for Mathematics in the Sciences, Leipzig, Germany
e-mail: jost@mis.mpg.de

N. Ay · J. Jost
Santa Fe Institute for the Sciences of Complexity, Santa Fe, NM, USA

N. Bertschinger
Frankfurt Institute for Advanced Studies, Frankfurt am Main, Germany

Goethe University, Frankfurt am Main, Germany

describe or build it. That principle, however, needs some qualification. One should distinguish between—usually useful—information about regularities of a structure or a process and—often useless—information about random details. The question is not only *information about what?*, but also *where is that information?*, that is, whether and how it is or can be internally stored in a system with limited capacity, at which level of a process information is needed to predict the continuation of a process, and where it can be found in a distributed system. In the latter case, we should, however, not only look for information that is exclusively located somewhere or that is shared between entities, but should also consider complementary or synergistic information, that is, information that only emerges when several sources are combined.

These lecture notes describe what is currently known about these questions, and they develop the underlying theoretical concepts and elucidate them at simple examples. Also, when we can quantify complexity concepts, we can also try to optimize the corresponding complexity measures. This will also be systematically discussed.

These notes are the result of a series of lectures that one of us (JJ) delivered at the Summer School in Como in July, 2018. They present work that we have done jointly during the last few years. JJ thanks Elisa Mastrogiacomio and Sergio Albeverio for organizing a very stimulating school, and the participants and the other lecturers, in particular Luciano Boi, Ivar Ekeland and Frank Riedel, for stimulating discussions.

2 Background: Principles of Information Theory

2.1 Shannon Information

The basic concept is that of the *Shannon Information* [52] of a random variable X , or equivalently, of a probability distribution p , when the possible values x_i of X are realized with probabilities $p_i = p(x_i)$. These probabilities satisfy $0 \leq p_i \leq 1$ for all i , with the normalization $\sum_i p_i = 1$. The Shannon information or entropy then is

$$H(X) = H(p_1, \dots, p_n) = - \sum_i p_i \log_2 p_i \text{ (bits)}. \quad (1)$$

This is the expected *reduction of uncertainty*, i.e., the *information gain*, if we learn which concrete value x_i of the random variable X from a known distribution p with probabilities $p_i = p(x_i)$ is realized.

This is the basic example: When we have two possible events occurring with equal probability $1/2$ (an unbiased coin) we thus gain $\log_2 2 = 1$ bit of information when we observe the outcome.

A fair dice yields $\log_2 6$ bits of information.

2.2 Mutual and Conditional Information

We now consider the situation where we have an additional random variable Y .

In the example of the dice, we could let $Y = 0$ (resp. 1) for an odd (even) result, each with probability $1/2$. According to the basic example, we have $H(Y) = 1$ bit. When we know Y , there remain only 3 possibilities for the value of X , each with probability $1/3$.

This leads us to the concept of *conditional information*; in this example, the remaining uncertainty about X when knowing Y is

$$H(X|Y) = \log_2 3. \quad (2)$$

Thus, the uncertainty about the value of X is reduced from $\log_2 6$ to $\log_2 3$ bit when knowing Y .

The joint information is related to the conditional information by

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y). \quad (3)$$

Thus, $H(X, Y) \leq H(X) + H(Y)$, and $<$ if X and Y are not independent. In the example, we have $H(X, Y) = H(X)$, since the value of X determines that of Y .

The information gain about X from knowing Y is called the *mutual information* of X and Y ,

$$MI(X : Y) = H(X) - H(X|Y). \quad (4)$$

In our example $MI(X : Y) = \log_2 6 - \log_2 3 = \log_2 2 = 1$ bit. From Y , we gain 1 bit of information about X .

The mutual information is symmetric,

$$MI(X : Y) = MI(Y : X). \quad (5)$$

The difference structure is perhaps the most important aspect. In many respects, (4) is more important and fundamental than (1), because we always have some prior knowledge, expressed here through Y , when we observe some X . Thus, the mutual information $MI(X : Y)$ tells us how much we can already infer about X when we know Y . By then observing X , we only gain the additional information $H(X|Y)$.

Summary:

$$H(X) = MI(X : Y) + H(X|Y) \quad (6)$$

$H(X)$ = how much you learn from observing X

$MI(X : Y)$ = how much you learn about X by observing Y

$H(X|Y)$ = how much you learn from observing X when you already know Y .

We can iterate the conditioning process with another random variable Z , to get the *conditional mutual information*

$$MI(X : Y|Z) = H(X|Z) - H(X|Y, Z). \tag{7}$$

$MI(X : Y|Z)$ quantifies how much additional mutual information between X and Y can be gained when we already know Z .

Careful: While always $H(X|Z) \leq H(X)$, we do *not necessarily* have $MI(X : Y|Z) \leq MI(X : Y)$.

Example: The XOR function (exclusive or):

x	y	z
0	0	0
1	0	1
0	1	1
1	1	0

where X, Y assume their two values independently with probability $1/2$ each. Thus, $MI(X : Y) = MI(X : Z) = MI(Y : Z) = 0$, but $MI(X : Y|Z) = MI(X : Z|Y) = MI(Y : Z|X) = 1$, because knowing the values of two of the variables determines that of the third.

2.3 Maximum Entropy

E. Jaynes' maximum-entropy principle [24]: Take the least informative estimate possible on the given information, that is, don't put any information into your model that is not based on the observed data. Look for p with maximal entropy $H(p)$ under the constraint that the expectation values of certain observables f_α be reproduced,

$$E_p f_\alpha = \sum_i f_\alpha^i p_i \text{ for } \alpha = 1, \dots, A. \tag{8}$$

The solution is an exponential distribution

$$p_j = \frac{1}{Z} \exp\left(\sum_\alpha \lambda_\alpha f_\alpha^j\right) \text{ with } Z = \sum_i \exp\left(\sum_\alpha \lambda_\alpha f_\alpha^i\right). \tag{9}$$

In particular, when there are no observations,

$$p_j = \frac{1}{n} \text{ for } j = 1, \dots, n. \tag{10}$$

2.4 Kullback-Leibler Divergence

A reference for the information geometric concepts that will be introduced and used here and in the sequel is [5]. The *Kullback-Leibler divergence* (*KL-divergence* for short) or *relative entropy* for two probability distributions p, q

$$D(p\|q) = \begin{cases} \sum_i p_i \log_2 \frac{p_i}{q_i} & \text{if } \text{supp } p \subset \text{supp } q \\ \infty & \text{else} \end{cases} \tag{11}$$

is positive ($D(p\|q) > 0$ if $p \neq q$), but not symmetric, as in general, $D(p\|q) \neq D(q\|p)$.

Example: The mutual information is the KL-divergence between the joint distribution and the product of the marginals,

$$MI(X : Y) = D(p(x, y)\|p(x)p(y)). \tag{12}$$

Among all distributions $p(x, y)$ with the same marginals $p(x) = \sum_y p(x, y)$, $p(y) = \sum_x p(x, y)$, the product distribution $p(x)p(y)$ has the largest entropy. This is, of course, a special case of Jaynes’ principle. That is, when we only know the marginals, Jaynes’ principle would suggest to take the product distribution as our estimate.

Example: The space of all probability distributions on two binary variables is a 3-dimensional simplex. It contains the 2-dimensional subfamily of product distributions. The extreme points of the simplex are the Dirac measures $\delta^{(x,y)}$, $x, y = 0, 1$. Maximization of the distance from the family of product distributions leads to distributions with support cardinality two (perfect correlation or anticorrelation) [4].

The formal way of expressing Jaynes’ principle is to project a given distribution onto the product family \mathcal{E} to maximize entropy while preserving the marginals, with π denoting that projection,

$$\begin{aligned} D(p\|\mathcal{E}) &:= \inf_{q \in \mathcal{E}} D(p\|q) = D(p\|\pi(p)) \\ &= H_{\pi(p)}(X, Y) - H_p(X, Y). \end{aligned} \tag{13}$$

3 Complexity

In this section, we want to introduce and discuss complexity concepts. But what is *complexity*? Some possible answers (see [5, 6] for a systematic discussion): **Complexity** is

1. the minimal effort or the minimal resources needed to describe or generate an object. Examples of such complexity concepts include algorithmic complexity

- (Kolmogorov [31], Chaitin [14], Solomonoff [57]); computational complexities; entropy (Shannon [52]), or entropy rate (Kolmogorov [31], Sinai [56]).
2. the minimal effort or the minimal resources needed to describe or generate the *regularities* or the *structure* of an object. Examples of such complexity concepts include Kolmogorov minimal sufficient statistics and related notions, stochastic complexity (Rissanen [46]), effective complexity (Gell-Mann and Lloyd [18]), excess entropy [53], also known as effective measure complexity [21], forecasting complexity [64], also introduced as statistical complexity by Crutchfield, Young, Shalizi [15, 51].
 3. the extent to which an object, as a whole, is more than the sum of its parts (Aristotle [1]), that is, the extent to which the whole cannot be understood by the analysis of the parts of the system in isolation, but only by also considering their interactions.

In order to systematically explore these aspects, we start with the most basic concept, that of algorithmic complexity [14, 31, 57] (see [33] for a systematic exposition). This concept expresses 1) in its purest form.

3.1 Algorithmic Complexity

The *algorithmic complexity* of an object, such as a number or a piece of text, is the length of the shortest computer program that generates or produces the object as output.¹ Typically, one cannot compute this complexity, but only provide an upper bound by producing a computer program, but does not know whether this is the shortest possible one.

From a conceptual perspective, the basic premise is that irregular or random structures have the highest algorithmic complexity, because they do not admit a short description. In other words, we want to characterize the complexity of a structure by the difficulty of its description. That is, we ask the question: How much can the description of a structure be simplified by utilizing regularities?

- Very simple structures need not be simplified any further.
- Random structures cannot be simplified.
- *Computational complexity* (see for instance the expositions in [38, 39]): Running time of shortest computer program that can generate the structure: A simple structure is produced quickly, whereas for a random one, everything has to be explicit in the program, and so, it does not need to run for a long time either.
- Random structures are not of interest for themselves, but only as members of an ensemble; it therefore suffices to describe the latter (Gell-Mann and Lloyd [18]).

¹ To make the complexity of different objects comparable, one needs to agree on a predetermined programming language; usually, one assumes some universal Turing machine, and changing that Turing machine will introduce an additive constant in the upper bounds.

3.2 External and Internal Complexity

The question that arises from the above concept of algorithmic complexity is how to compute it, that is, how to find the shortest description of a given structure. Quite apart from the fact that this depends on the choice of the device we use to evaluate it (in theory: some universal Turing machine, and the choice of that Turing machine then introduces an additive constant), in practice, we have only bounded means to represent a structure. Thus: What do we want to know? We want to

1. know a rich and complex structure,
2. but represent it most efficiently.

More formally, we want to

1. maximize *external complexity*,
2. but minimize *internal complexity*.

This perspective was introduced in [25]. For an application in pattern classification, see for instance [3].

3.3 Optimization Principles

Organisms live in and interact with a complex environment, see for instance [61] (for a measure theoretical approach, see [7]), and need to maintain their own autopoiesis [37]. A modern society consists of several complex subsystems that follow their own rules, but need to interact with each other [35, 36]. With the concept of Shannon information, we can formulate some abstract principles that either maximize or minimize some kind of complexity (we follow [26] here). The basic versions, however, lead to trivial results, as we shall now see.

1. Gain as much information as possible: Look at random patterns
2. Avoid surprises: Look at blank screen
3. Try to predict future sensory inputs as accurately as possible on the basis of the current ones (and perhaps try to bring yourself into a state where this is possible [17])
4. Try to manipulate the environment such that the results of own actions are as accurately predictable as possible [32].
5. Maximize

$$\begin{aligned}
 & MI(S_{t+1} : E_t) - MI(S_{t+1} : E_t | S_t) \\
 = & H(S_{t+1}) - H(S_{t+1} | E_t) - H(S_{t+1} | S_t) + H(S_{t+1} | E_t, S_t) \quad (14)
 \end{aligned}$$

to establish the strongest possible correlation between the current state E_t of the environment and future sensory data S_{t+1} , but such that this correlation can already be predicted from the current input S_t [9]

To proceed further, let us discuss some questions.

1. Q: Why should a system model an external probability distribution?
A: To make predictions on the basis of regularities
2. Q: How can this be achieved in an environment that is vastly more complex than the system itself?
A: Detect regularities
3. Q: How to detect regularities?
A: Because of 2), the system is forced to compress.

These answers have some consequences in various fields:

- Psychology: Use heuristics [19, 54, 55]
- Cognition: External versus internal complexity [25]
- Statistics: Avoid overfitting
- Statistical learning theory: Start with models with few parameters and gradually increase as you learn (Vapnik-Chervonenkis) [59, 60].

3.4 Correlations in Time Series

We can also use the information theoretical notions to evaluate the complexity of a time series in terms of the correlations that it exhibits. A time series X_t , $t \in \mathbb{N}$ could possibly have

- No regularities: $H(X_t|X_{t-1}) = H(X_t)$
- the Markov property: $H(X_t|X_{t-1}, X_{t-2}, \dots) = H(X_t|X_{t-1})$, or
- Long term correlations, as in texts, genetic sequences, ...

To evaluate this, we quantify how much new information is gained when one already knows n consecutive symbols and then sees the $(n + 1)$ st. (Grassberger [21]).

For which n is this largest? When n is small, one perhaps cannot predict much, and if n is large, one may be able to guess the rest anyway.

The larger this n , the more complex the sequence.

For genetic sequences, $n \sim 14$ [47], for amino acid sequences (proteins) $n \sim 5$.

In literature analysis, such a principle can be used to evaluate the complexity of language [16].

A more sophisticated concept is the *genon* concept of molecular biology [30, 48, 49].

3.5 Complementarity

Instead of trying to predict the environment, one can also let the environment do the computation itself (see [26]).

If you want to catch a ball, you do not use Newtonian mechanics to compute the trajectory, but simply run so that the ball appears under a constant angle. The environment computes the trajectory, and you only need to sample. This outsourcing of computation represents one mechanism for the compression mentioned in Sect. 3.3.

More generally, embodied cognition has emerged as a new paradigm in robotics [43].

3.6 Hierarchical Models and Complexity Measures

In this section, we follow [5, 6]. Returning to Jaynes' approach, we could maximize entropy while preserving marginals among subsets of variables. For instance, for a distribution on 3 variables, we could prescribe all single and pairwise marginals.

Assume that we have a state set V that consists of the possible values of N variables. We then consider the hierarchy

$$\mathfrak{S}_1 \subseteq \mathfrak{S}_2 \subseteq \dots \subseteq \mathfrak{S}_{N-1} \subseteq \mathfrak{S}_N := 2^V, \quad (15)$$

where \mathfrak{S}_k is the family of subsets of V with $\leq k$ elements, from which we get the set of probability distributions $\mathcal{E}_{\mathfrak{S}_k}$ with dependencies of order $\leq k$. For instance, $\mathcal{E}_{\mathfrak{S}_1}$ is the family of distributions that are simply the products of their marginals. In particular, for a probability distribution in this family, there are no correlations between the probabilities of two or more of the variables. In $\mathcal{E}_{\mathfrak{S}_2}$, we then allow for pairwise correlations, but no triple or higher order ones.

We point out that one can also consider other families of subsets of V and the corresponding probability distributions. For instance, when V is the ordered set of integers $\{1, \dots, N\}$, one could consider the family of those subsets that consist of uninterrupted strings of length $\leq k$. This will be our choice when we discuss the excess entropy below.

We let $\pi_{\mathfrak{S}_k}$ be the projection on $\mathcal{E}_{\mathfrak{S}_k}$, $p^{(k)} := \pi_{\mathfrak{S}_k}(p)$. For instance, $p^{(1)}$ is the product distribution with the same marginals as p .

We have the important Pythagorean relation

$$D(p^{(l)} \parallel p^{(m)}) = \sum_{k=m}^{l-1} D(p^{(k+1)} \parallel p^{(k)}), \quad (16)$$

for $l, m = 1, \dots, N-1, m < l$. In particular,

$$D(p \parallel p^{(1)}) = \sum_{k=1}^{N-1} D(p^{(k+1)} \parallel p^{(k)}). \quad (17)$$

If we take configurations with dependencies of order $\leq k$, we get the **Complexity measure** [6] with weight vector $\alpha = (\alpha_1, \dots, \alpha_{N-1}) \in \mathbb{R}^{N-1}$

$$C_\alpha(p) := \sum_{k=1}^{N-1} \alpha_k D(p \parallel p^{(k)}) \tag{18}$$

$$= \sum_{k=1}^{N-1} \beta_k D(p^{(k+1)} \parallel p^{(k)}), \tag{19}$$

with $\beta_k := \sum_{l=1}^k \alpha_l$.

$p^{(k)}$ is the distribution of highest entropy among all those with the same correlations of order $\leq k$ as p .

Thus, we consider a weighted sum of the higher order correlation structure.

Examples:

- Tononi-Sporns-Edelman complexity [58]: $\alpha_k = \frac{k}{N}$ addresses the issue of the interplay between differentiation and integration in complex systems (for an analysis of system differentiation from an information theoretical perspective, see also [28])
- Stationary stochastic process X_n : Conditional entropy

$$h_p(X_n) := H_p(X_n \mid X_1, \dots, X_{n-1}).$$

Entropy rate or Kolmogorov–Sinai entropy [31, 56]

$$h_p(X) := \lim_{n \rightarrow \infty} h_p(X_n) = \lim_{n \rightarrow \infty} \frac{1}{n} H_p(X_1, \dots, X_n), \tag{20}$$

Excess entropy (Grassberger [21])

$$\begin{aligned} E_p(X) &:= \lim_{n \rightarrow \infty} \sum_{k=1}^n (h_p(X_k) - h_p(X)) \\ &= \lim_{n \rightarrow \infty} (H_p(X_1, \dots, X_n) - nh_p(X)) \end{aligned} \tag{21}$$

$$= \lim_{n \rightarrow \infty} \underbrace{\sum_{k=1}^{n-1} \frac{k}{n-k} D(p_n^{(k+1)} \parallel p_n^{(k)})}_{=: E_p(X_n)}, \tag{22}$$

where we choose \mathfrak{S}_k as the sequences of integers $j + 1, j + 2, \dots, j + \ell$ with $\ell \leq k$. The excess entropy measures the non-extensive part of the entropy, i.e. the amount of entropy of each element that *exceeds* the entropy rate.

3.7 Interactions Between Levels

The question of emergence, that is, how a higher level that is (at least partially) autonomous from lower levels, arises in many disciplines. For example, classical mechanics arises from an underlying quantum structure, but the laws of classical mechanics are causally closed, in the sense that for computing trajectories of Newtonian particles, we do not need information from the quantum level. Likewise, human genetics rests on the laws of Mendel and does not need to consider an underlying biochemical level. In other fields it is often not so clear, however, to what extent laws operate autonomously at a certain level without needing permanent or at least regular access to some lower level. For instance, does it suffice for understanding macroeconomic processes to consider relations between macroeconomic variables, or is an input from the microeconomic level essentially needed? Or can one understand social dynamics without access to the psychic and mental states of the participating individuals? For a general discussion of the issue of emergence from the perspective developed in the present contribution, see for instance [29].

Here, we describe the approach of [41, 42] (and refer to [41] for references to earlier work). We consider a structure

$$\begin{array}{ccc}
 \widehat{X} & \xrightarrow{\psi} & \widehat{X}' \\
 \pi \uparrow & & \uparrow \pi \\
 X & \xrightarrow{\phi} & X'
 \end{array}$$

with basic level X, X' and higher level $\widehat{X}, \widehat{X}'$; an arrow $Y \rightarrow Y'$ represents a discrete time step where X, X' form a Markov process, with transition kernel ϕ , which can be observed at the higher level $\widehat{X}, \widehat{X}'$ in a lossy fashion.

The higher level could result from averaging or aggregating the lower level. Think of \widehat{X} as a coarse-graining of X given by an observation map π .

We can propose several criteria for the upper process being closed in the sense that it depends on the lower process only through some initialization.

I Informational closure: The higher process is informationally closed, i.e. there is no information flow from the lower to the higher level. Knowledge of the microstate will not improve predictions of the macrostate.

$$MI(\widehat{X}' : X | \widehat{X}) = 0 \tag{23}$$

where the conditional mutual information

$$MI(\widehat{X}' : X | \widehat{X}) = H(\widehat{X}' | \widehat{X}) - H(\widehat{X}' | X) \tag{24}$$

measures the reduction in uncertainty about \widehat{X}' when knowing X instead of only \widehat{X} .

II Observational commutativity: It makes no difference whether we perform the aggregation first and then observe the upper process, or we observe the process on the microstate level, and then lump together the states.

Kullback-Leibler divergence between the lower and the upper transition kernel from X to \widehat{X}' is 0, for some initial distribution on X .

$$I \Rightarrow II, \text{ and in deterministic case also } II \Rightarrow I. \quad (25)$$

(In I , probabilities at \widehat{X} , in II at X)

III Commutativity: There exists a transition kernel ψ such that the diagram commutes (Görnerup-Jacobi, 2010)

$$II \Rightarrow III, \text{ and in deterministic case also } III \Rightarrow II. \quad (26)$$

II : Transition kernels satisfy $\Psi = \Pi\Phi\Pi^T$

III : Transition kernels satisfy $\Psi\Pi = \Pi\Phi$

IV Markovianity: $\widehat{X}, \widehat{X}'$ forms again a Markov process (Shalizi-Moore, 2003).

$$I \Rightarrow IV, \text{ but } IV \not\Rightarrow III. \quad (27)$$

V Predictive efficiency: A more abstract formulation is that an emergent level corresponds to an efficiently predictable process, that is, one that can be predicted in its own terms, without permanent recourse to a lower level.

3.7.1 A Test Case: The Tent Map

We now evaluate the preceding concepts at the example of the tent map, following [40] (see also [2] for background).

$$T(x) = \begin{cases} 2x & \text{if } 0 \leq x \leq 1/2 \\ 2 - 2x & \text{else} \end{cases}$$

The tent map is a basic example of a chaotic dynamical iteration, because at every step differences between values can get doubled, and therefore, after several steps, even very tiny differences between initial values can become macroscopically large. The folding at $x = 1/2$ ensures that nevertheless the unit interval is mapped to itself. Thus, some differences also get reduced. Understanding this interplay between amplification and reduction of differences is surprisingly subtle, as one may also see in the following.

For a threshold value $\alpha \in [0, 1]$ we define the symbolic dynamics

$$\begin{aligned} \phi_\alpha : X &\rightarrow \hat{X} = \{0, 1\} \\ \phi_\alpha(x) &:= \begin{cases} 0 & \text{if } 0 \leq x < \alpha \\ 1 & \text{else} \end{cases} \end{aligned}$$

The sequence $x_n = T^n(x)$, for an initial value $x \in X$, yields the derived symbol dynamics $s_n = \phi_\alpha(x_n) \in \{0, 1\}$.

The probability of finding s_n in the state 0 is the probability that x_n lies in the interval $[0, \alpha]$ (which is α for the tent map).

We consider the symbolic dynamics derived from consecutive time steps

$$(s_{n+m}, s_{n+m-1}, \dots, s_n),$$

with $k \in \mathbb{N}$

$$s_k(x) = \begin{cases} 0 & \text{if } T^k(x) < \alpha \\ 1 & \text{if } T^k(x) \geq \alpha \end{cases}.$$

For comparison, we take a random sequence $\xi_n \in [0, 1]$ (uniformly, i.i.d.), and consider the corresponding symbolic dynamics

$$\sigma_n := \begin{cases} 0 & \text{if } \xi_n < \alpha \\ 1 & \text{if } \xi_n \geq \alpha. \end{cases}$$

The question now is: Are there systematic differences between the symbolic sequence s_n derived from iterations of the tent map and σ_n ?

For $\alpha = 1/2$, they look the same (in fact, we simply have a Bernoulli sequence: the values 0 and 1 occur with equal probability $1/2$; $p(0) = p(1) = 1/2$). If we don't know x , s_n looks as random as σ_n . The transition probabilities are

$$p(0|0) = p(1|0) = p(0|1) = p(1|1) = 1/2.$$

We next consider $\alpha = 2/3$. Put $x_n := T^n(x)$.

$\sigma_n = 0$ and $\sigma_n = 1$ occur independently with probabilities $2/3$ and $1/3$.

When $s_n = 1$, that is, $2/3 < x_n \leq 1$, then $0 \leq x_{n+1} < 2/3$, that is $s_{n+1} = 0$. Thus, there is no transition from 1 to 1. For the state $s_n = 0$, both transitions are equally likely: when $0 \leq x_n \leq 1/3$, we have $0 \leq x_{n+1} \leq 2/3$, that is, $s_{n+1} = 0$, while for $1/3 < x_n \leq 2/3$, we get $s_{n+1} = 1$. Thus, for s_n ,

$$p(0|0) = p(1|0) = 1/2, \quad p(0|1) = 1, \quad p(1|1) = 0$$

while for σ_n

$$p(0|0) = p(0|1) = 2/3, \quad p(1|0) = p(1|1) = 1/3.$$

This leads us to the concept of *forbidden sequences*. While for the threshold $\alpha = 1/2$, the symbolic dynamics of the tent map cannot be distinguished from that of a random

sequence, and is Markovian, in contrast, for the threshold $\alpha = 2/3$, the sequence 11 does not occur, and the symbolic dynamics is different from a random one, but still Markovian.

For other thresholds, we can also get longer forbidden sequences and non-Markovian symbolics.

Even from a random sequence ξ_n , we can derive non-Markovian symbolic dynamics.

Let $x^1, x^2 \in [0, 1]$; we consider the symbolic rule

$$s(x^1, x^2) = \begin{cases} 0 & \text{if } x^1 \leq x^2 \\ 1 & \text{if } x^2 < x^1. \end{cases}$$

For our random sequence, take $x^1 = \xi_n, x^2 = \xi_{n+1}$. Thus, we draw the points x^1, x^2 randomly and independently.

The state probabilities are again $p(0) = p(1) = 1/2$, but the transition probabilities now depend on the history. The more 1s we have seen, the less likely it is to see another 1, because then ξ_n is expected to be very small, hence most likely, $\xi_{n+1} > \xi_n$.

We now analyze the *information flow* of this example. The information flow between the micro-level corresponding to state x_n and the coarse-grained level s_n is the conditional mutual information

$$MI(s_{n+1} : x_n | s_n) = H(s_{n+1} | s_n) - H(s_{n+1} | s_n, x_n).$$

Since s_{n+1} is fully determined by x_n , the second term vanishes,

$$MI(s_{n+1} : x_n | s_n) = H(s_{n+1} | s_n),$$

i.e., the information flow = conditional entropy on the coarse grained level, which has a local minimum at $\alpha = 2/3$.

Instead of drawing information from below, the upper level system relies on its memory.

4 Information Decomposition

We finally turn to the concept of information decomposition. To motivate it, we start with the *transfer entropy* [50]²

² Such a principle had already been introduced by the econometrician Granger [20] who wrote “We say that Y_t is causing X_t if we are better able to predict X_t using all available information than if the information apart from Y_t had been used.” In the econometric literature, this principle was applied only in linear settings. As [8], explained, the transfer entropy can be seen as an operationalization of this principle in a general context.

$$TE(Z \rightarrow X) := MI(X_+ : Z_- | X_-) \tag{28}$$

where the subscript—refers to the past and + to the future. $TE(Z \rightarrow X)$ quantifies the amount of information contained in Z about the future of X that cannot be obtained from its own past.

Problem: $X_+ = \text{XOR}(X_-, Z_-)$:

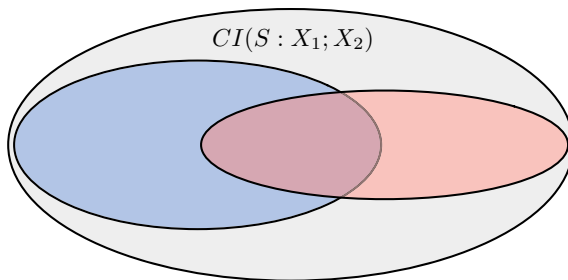
Here, the information in Z_- is only useful together with that of X_- . The transfer entropy cannot distinguish this situation from one where X_- does not contribute and Z_- determines X_+ by itself.

This problem is addressed by information decomposition. It was started by Williams and Beer [63] (but their measure I_{\min} of shared information does not distinguish whether different random variables carry the same information or just the same amount of information), and continued by Harder, Salge, Polani [23], Griffith and Koch [22], Bertschinger, Rauh, Olbrich, Ay, Banerjee, Jost [12, 13, 44, 45], and taken up by many other people (see for instance the references in [34]), with applications in different fields, like neuroscience [62]. There is no optimal solution, but that of Bertschinger, Rauh, Olbrich, Jost, Ay [13] (called the *BROJA* decomposition in the community) is currently the most widely accepted.

To describe our approach, we consider three random variables X_1, X_2 and S . The (total) mutual information $MI(S : X_1, X_2)$ quantifies the total information that is gained about S if the outcomes of X_1 and X_2 are known. How do X_1 and X_2 contribute to this information? For two explanatory variables, we expect four contributions to $MI(S : X_1, X_2)$:

$$\begin{aligned} MI(S : X_1, X_2) = & SI(S : X_1; X_2) \quad \text{shared information} \\ & + UI(S : X_1 \setminus X_2) \quad \text{unique information of 1} \\ & + UI(S : X_2 \setminus X_1) \quad \text{unique information of 2} \\ & + CI(S : X_1; X_2) \quad \text{complementary or synergistic information.} \end{aligned}$$

Here, $UI(S : X_1 \setminus X_2)$ is the information that X_1 has, but X_2 does not have, $SI(S : X_1; X_2)$ is the information that both of them have individually. Perhaps the most interesting term is the last, $CI(S : X_1; X_2)$, the information that only emerges if X_1 and X_2 pool their knowledge. This term is best illustrated in the **XOR** example discussed below.



We consider some examples. **AND**

x_1	x_2	s	$p(x_1, x_2, s)$
0	0	0	1/4
1	0	0	1/4
0	1	0	1/4
1	1	1	1/4

Here, x_1 and x_2 jointly determine s , but cannot be fully recovered from s .

When 1 has the value $x_1 = 0$, she can exclude $s = 1$, and analogously for 2.

Thus, when they both see 0, they share the information that $s = 0$.

The mechanism loses some information. When X_1, X_2 are i.i.d.,

$$H(X_1, X_2) = 2 \text{ bits,}$$

but

$$H(S) = MI(S : X_1, X_2) = -\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} \approx .811 \text{ bits.}$$

In general, we may have both correlations between the input variables and relations created by the mechanism that computes S .

We next recall **XOR** from Sect. 2.2:

x_1	x_2	s
0	0	0
1	0	1
0	1	1
1	1	0

Neither 1 nor 2 can determine the value of S by herself, but the value of the other is needed for that. This is a clear case of synergistic information only.

Our approach: Unique and shared information should only depend on the marginal distribution of the pairs (S, X_1) and (S, X_2) . This idea can be explained from an operational interpretation of unique information: Namely, if X_1 has unique information about S (with respect to X_2), then there must be some way to exploit this information. More precisely, there must be a situation in which X_1 can use this information to perform better at predicting the outcome of S .

In this interpretation, 1 possesses unique information about S compared with 2, if there exists a reward function for which 1 can achieve a higher expected reward based on her value x_1 and her knowledge of the conditional distribution $p(s|x_1)$ than if she knew and utilized instead the conditional distribution of 2.

Thus, unique and shared information depend only on pairwise marginals. Only the synergistic information includes higher order dependencies. In that sense, synergy becomes a measure of higher order interactions, in the sense of information geometry.

From a conceptual perspective, and independently of the way the different terms in the decomposition are quantified, it is important to understand synergy, in order to clarify discussions that have become quite sterile, like the relative importance of genes and environment in biology. For a perspective in this direction, see [27].

References

1. Aristoteles: Philosophische Schriften 5. Metaphysik. Nach der Übersetzung von Hermann Bonitz bearbeitet von Horst Seidl, Felix Meiner, Hamburg (1995)
2. Atay, F., Jalan, S., Jost, J.: Randomness, chaos, and structure. *Complexity* **15**, 29–35 (2009)
3. Avdiyenko, L., Bertschinger, N., Jost, J.: Adaptive information-theoretical feature selection for pattern classification. In: *Computational Intelligence*, pp. 279–294. Springer International Publishing (2015)
4. Ay, N.: An Information-geometric approach to a theory of pragmatic structuring. *Ann. Prob.* **30**(1), 416–436 (2002)
5. Ay, N., Jost, J., Lê, H.V., Schwachhöfer, L.: *Information Geometry*. Springer, *Ergebnisse der Mathematik* (2017)
6. Ay, N., Olbrich, E., Bertschinger, N., Jost, J.: A geometric approach to complexity. *Chaos* **21**, 037103 (2011). <https://doi.org/10.1063/1.3638446>
7. Ay, N., Loehr, W.: The Umwelt of an embodied agent a measure-theoretic definition. *Theory Biosci.* **134**(3), 105–116 (2015)
8. Barnett, L., Barnett, A., Seth, A.: Granger causality and transfer entropy are equivalent for Gaussian variables. *Phys. Rev. Lett.* **103**(23), 38701 (2009)
9. Bertschinger, N., Olbrich, E., Ay, N., Jost, J.: Autonomy: an information theoretic perspective. *Biosystems* **91**, 331–345 (2008)
10. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J.: Shared information—new insights and problems in decomposing information in complex systems. *Proc. ECCS* 251–269 (2012)
11. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., Ay, N.: Quantifying unique information. *Entropy* **16**, 2161–2183 (2014). <https://doi.org/10.3390/e16042161>
12. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J.: Shared information—new insights and problems in decomposing information in complex systems. In: *Proceedings of the European Conference on Complex Systems, Brussels, Belgium, 2–7 September 2012*, pp. 251–269
13. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., Ay, N.: Quantifying unique information. *Entropy* **16**, 2161–2183 (2014)
14. Chaitin, G.: On the lengths of programs for computing finite binary sequences. *JACM* **13**, 547–569 (1966)
15. Crutchfield, J.P., Young, K.: Inferring statistical complexity. *Phys. Rev. Lett.* **63**, 105–108 (1989)
16. Efer, T., Heyer, G., Jost, J.: Shakespeare unter den Deutschen (C. Jansohn, Hrsg.)
17. Friston, K.: The free-energy principle: a unified brain theory? *Nat. Rev. Neurosc.* **11**, 127–138 (2010)
18. Gell-Mann, M., Lloyd, S.: Information measures, effective complexity, and total information. *Complexity* **2**, 44–52 (1996)
19. Gigerenzer, G., Todd, P.: *Simple Heuristics that Make us Smart*. Oxford University Press (1999)
20. Granger, C.: Investigating causal relations by econometric models and cross-spectral methods. *Econometrica* **37**, 424–438 (1969)
21. Grassberger, P.: Toward a quantitative theory of self-generated complexity. *Int. J. Theor. Phys.* **25**, 907–938 (1986)
22. Griffith, V., Koch, C.: Quantifying synergistic mutual information. In: Prokopenko, M. (ed.) *Guided Self-Organization: Inception*, vol. 9, pp. 159–190. Springer, Berlin/Heidelberg, Germany (2014)

23. Harder, M., Salge, C., Polani, D.: A bivariate measure of redundant information. *Phys. Rev. E* **87**, 012130 (2013)
24. Jaynes, E.T.: *Probability Theory: The Logic of Science*. Cambridge University Press (2003)
25. Jost, J.: External and internal complexity of complex adaptive systems. *Theory Biosc.* **123**, 69–88 (2004)
26. Jost, J.: Sensorimotor contingencies and the dynamical creation of structural relations underlying percepts. In: *Proceedings of the Strüngmann Forum: Where's the Action? The Pragmatic Turn in Cognitive Science*. MIT Press, im Druck (2016)
27. Jost, J.: *Biological information*. *Theory Biosci.* (2020)
28. Jost, J., Bertschinger, N., Olbrich, E., Ay, N., Frankel, S.: An information theoretic approach to system differentiation on the basis of statistical dependencies between subsystems. *Phys. A* **378**, 1–10 (2007)
29. Jost, J., Bertschinger, N., Olbrich, E.: Emergence. A dynamical systems approach. *New Ideas Psych.* **28**, 265–273 (2010)
30. Jost, J., Scherrer, K.: Information theory, gene expression, and combinatorial regulation: a quantitative analysis. *Theory Biosci.* **133**, 1–21 (2014)
31. Kolmogorov, A.N.: Three approaches to the quantitative definition of information. *Probl. Inf. Transm.* **1**, 4–7 (1965)
32. Klyubin, A., Polani, D., Nehaniv, C.: Empowerment: a universal agent-centric measure of control. *Proc. IEEE CEC* **2005**, 128–135 (2005)
33. Li, M., Vitányi, P.: *An Introduction to Kolmogorov Complexity and Its Applications*. Springer (1997)
34. Lizier, J., Bertschinger, N., Jost, J., Wibral, M.: Information decomposition of target effects from multi-source interactions: perspectives on previous, current and future work. *Entropy* **20**, 307 (2018). <https://doi.org/10.3390/e20040307>
35. Luhmann, N.: *Soziale Systeme*. Suhrkamp, Frankfurt (1984)
36. Luhmann, N.: *Die Gesellschaft der Gesellschaft*, 2 Bde. Suhrkamp, Frankfurt (1997)
37. Maturana, H., Varela, F.: *Autopoiesis and Cognition*. Reidel, Boston (1979)
38. Moore, C., Mertens, S.: *The Nature of Computation*. Oxford University Press (2011)
39. Papadimitriou, C.: *Computational Complexity*. Addison Wesley (1994)
40. Pfante, O., Olbrich, E., Bertschinger, N., Ay, N., Jost, J.: Closure measures for coarse-graining of the tent map. *Chaos* **24**, 013136 (2014)
41. Pfante, O., Bertschinger, N., Olbrich, E., Ay, N., Jost, J.: Comparison between different methods of level identification. *Adv. Complex Syst.* **17**, 1450007 (21 Seiten) (2014)
42. Pfante, O., Bertschinger, N., Olbrich, E., Ay, N., Jost, J.: Wie findet man eine geeignete Beschreibungsebene für ein komplexes System? In: *Jahrbuch der Max-Planck-Gesellschaft*, im Druck (2016)
43. Pfeifer, R., Bongard, J.: *How the Body Shapes the Way we Think*. MIT Press (2007)
44. Rauh, J., Banerjee, P., Olbrich, E., Jost, J., Bertschinger, N.: On extractable shared information. *Entropy* **19**(7), 328 (2017)
45. Rauh, J., Bertschinger, N., Olbrich, E., Jost, J.: Reconsidering unique information: towards a multivariate information decomposition. In: *Proceedings of 2014 IEEE International Symposium on Information Theory (ISIT)*, Honolulu, HI, USA, 29 June–4 July 2014, pp. 2232–2236
46. Rissanen, J.: *Stochastic Complexity in Statistical Inquiry*. World Scientific (1989)
47. Sadosky, M.: Genes, information and sense: complexity and knowledge retrieval. *Theory Biosci.* **127**(2), 69–78 (2008)
48. Scherrer, K., Jost, J.: The gene and the Genon concept: a functional and information-theoretic analysis. *Mol. Syst. Biol.* **3**, 87, Epub 2007 Mar 13: EMBO and Nature Publishing Group
49. Scherrer, K., Jost, J.: Gene and Genon concept: coding versus regulation. *Theory Biosc.* **126**, 65–113 (2007)
50. Schreiber, T.: Measuring information transfer. *PRL* **85**, 461–464 (2000)
51. Shalizi, C.R., Crutchfield, J.P.: Computational mechanics: pattern and prediction, structure and simplicity. *J. Stat. Phys.* **104**, 817–879 (2001)

52. Shannon, C.: The Mathematical Theory of Communication (1948), reprinted in: Shannon, C., Weaver, W. The Mathematical Theory of Communication. Blahut, R., Hajek, B. (eds.), pp. 29–125. Univ. Illinois Press (1998)
53. Shaw, R.: The Dripping Faucet as a Model Chaotic System. Aerial Press, Santa Cruz (1984)
54. Simon, H.A.: A behavioral model of rational choice. *Quart. J. Econom.* **69**(1), 99–118 (1955)
55. Simon, H.A.: Rational choice and the structure of environments. *Psych. Rev.* **63**, 129–138 (1956)
56. Sinai, J.: On the concept of entropy for a dynamical system, (Russian) *Dokl. Akad. Nauk SSSR* **124**, 768–771 (1959)
57. Solomomoff, R.: A formal theory of inductive inference, Part I. *Inf. Control* **7**, 1–22. Part II. *Inf. Control* **7**, 224–254 (1964)
58. Tononi, G., Sporns, O., Edelman, G.M.: Measures of degeneracy and redundancy in biological networks. *PNAS* **96**, 3257–3267 (1999)
59. Vapnik, V.N.: The Nature of Statistical Learning Theory. Springer (1995)
60. Vapnik, V.N.: Statistical Learning Theory. Wiley-Interscience (1998)
61. von Uexküll, J.J.: Umwelt und Innenwelt der Tiere. In: Mildenerger, F., Herrmann, B. (eds.) *Klassische Texte der Wissenschaft*. Springer, Spektrum (2014)
62. Wibral, M., Priesemann, V., Kay, J.W., Lizier, J.T., Phillips, W.A.: Partial information decomposition as a unified approach to the specification of neural goal functions. *Brain Cogn.* **112**, 25–38 (2017)
63. Williams, P., Beer, R.: Nonnegative decomposition of multivariate information. [arXiv:1004.2515v1](https://arxiv.org/abs/1004.2515v1) (2010)
64. Zambella, D., Grassberger, P.: Complexity of forecasting in a class of simple models. *Complex Syst.* **2**, 269 (1988)