



Asymptotic dependency structure of multiple signals

Asymptotic equipartition property for diagrams of probability spaces

Rostislav Matveev¹  · Jacobus W. Portegies²

Received: 7 February 2018 / Revised: 16 August 2018 / Published online: 12 October 2018
© The Author(s) 2018

Abstract

We formalize the notion of the dependency structure of a collection of *multiple* signals, relevant from the perspective of information theory, artificial intelligence, neuroscience, complex systems and other related fields. We model multiple signals by commutative diagrams of probability spaces with measure-preserving maps between some of them. We introduce the asymptotic entropy (pseudo-)distance between diagrams, expressing how much two diagrams differ from an information-processing perspective. If the distance vanishes, we say that two diagrams are asymptotically equivalent. In this context, we prove an asymptotic equipartition property: any sequence of tensor powers of a diagram is asymptotically equivalent to a sequence of homogeneous diagrams. This sequence of homogeneous diagrams expresses the relevant dependency structure.

Keywords Asymptotic equipartition property · Entropy distance · Diagrams of probability spaces · Multiple signals

1 Introduction

According to usual modeling assumptions in information theory, a discrete signal is cut into a collection of long words of length n , whose particular representation is irrelevant (each word is considered as an atomic object without inner structure), and small errors are allowed. If the signal is modeled as a sequence of independently, identically distributed random variables, there is only one relevant variable determining

✉ Rostislav Matveev
matveev@mis.mpg.de

Jacobus W. Portegies
j.w.portegies@tue.nl

¹ Max-Planck-Institut für Mathematik in der Naturwissenschaften, Inselstr. 22, 04103 Leipzig, Germany

² Eindhoven University of Technology, De Zaale, Eindhoven, The Netherlands

the signal, namely the entropy: the exponential growth rate of the number of typical words of length n . We elaborate on this point of view below in Sect. 1.1.

Similarly, if one probes a measure-preserving dynamical system at a discrete sequence of times with a finite-output measurement device and counts measurement trajectories of length n , while discarding rarely appearing, untypical ones, one arrives at the notion of entropy of a system-measurement pair. Entropy, in this case, is the exponential growth rate of the number of typical trajectories with respect to the length n . The supremum of such entropies over varying measurement devices is the Kolmogorov–Sinai entropy of a measure-preserving dynamical system. According to a theorem of Ornstein [17], the entropy is the only invariant of the isomorphism classes of certain types of dynamical systems (Bernoulli shifts).

In information theory, but also in artificial intelligence, neuroscience and the theory of complex systems, one usually studies *multiple* signals at once. Likewise, a dynamical system is often observed with multiple measurement devices simultaneously. In these cases, one assumes in addition that the relations between the signals are essential. In this article we characterize, under these modeling assumptions, the relevant invariants in multiple signals, that are obtained as i.i.d. samples from random variables. We explain this in more detail in Sects. 1.3 and 1.4.

We will now explain our point of view on entropy for a single signal, that is, for a single probability space.

1.1 Probability spaces and their entropy

First we consider a finite probability space $X = (S, p)$, where S is a finite set, and p is a probability measure on S . For simplicity, assume for now that the measure p has full support. Next, we consider the, so-called, Bernoulli sequence of probability spaces

$$X^n = (S^n, p^{\otimes n})$$

where S^n denotes the n -fold Cartesian product of S , and $p^{\otimes n}$ is the n -fold product measure.

The *entropy* of X is the exponential growth rate of the *observable cardinality* of tensor powers of X . The observable cardinality, loosely speaking, is the cardinality of the set X^n after the biggest possible subset of small measure has been removed. It turns out that the observable cardinality of X^n might be much smaller than $|S|^n$, the cardinality of the whole of X^n , in the following sense.

The *Asymptotic Equipartition Property* states that for every $\varepsilon > 0$ and $n \gg 0$ one can find a, so-called, *typical subset* $A_\varepsilon^{(n)} \subset S^n$, defined as a subset that takes up almost all of the mass of X^n and the probability distribution on $A_\varepsilon^{(n)}$ is almost uniform on the normalized logarithmic scale, as stated in the following theorem, see [8].

Theorem 1.1 (Asymptotic equipartition property) *Suppose $X = (S, p)$ is a finite probability space. Then, for every $\varepsilon > 0$ and every $n \gg 0$ there exists a subset $A_\varepsilon^{(n)} \subset X^n$ such that*

1. $p^{\otimes n}(A_\varepsilon^{(n)}) > 1 - \varepsilon$
2. For any $a, a' \in A_\varepsilon^{(n)}$ holds

$$\left| \frac{\ln p^{\otimes n}(a)}{n} - \frac{\ln p^{\otimes n}(a')}{n} \right| < \varepsilon.$$

Moreover, if $A_\varepsilon^{(n)}$ and $B_\varepsilon^{(n)}$ are two subsets of X^n satisfying the two conditions above, then their cardinalities satisfy

$$\left| \frac{\ln |A_\varepsilon^{(n)}|}{n} - \frac{\ln |B_\varepsilon^{(n)}|}{n} \right| < 2\varepsilon. \tag{1}$$

The cardinality $|A_\varepsilon^{(n)}|$ may be much smaller than $|S|^n$, but it will still grow exponentially with n . Even though there are generally many choices for such a set $A_\varepsilon^{(n)}$, in view of the property (1) in Theorem 1.1, the exponential growth rate with respect to n is well-defined up to 2ε .

The limit of the growth rate as $\varepsilon \rightarrow 0+$ is called the entropy of X

$$\text{Ent}(X) := \lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \ln |A_\varepsilon^{(n)}|. \tag{2}$$

This point of view on entropy goes back to the original idea of Boltzmann [3,4], according to which entropy is the logarithm of the number of equiprobable states, that a system, comprised of many identical weakly interacting subsystems, may take on. It was further developed and applied to Information Theory by Shannon [19], and in the context of dynamical systems by Kolmogorov and Sinai [12,13,21].

Entropy is especially easy to evaluate if the space is uniform, since for any finite probability space with the uniform distribution the observable cardinality is equal to the cardinality of the whole space and therefore

$$\text{Ent}(X) = \ln |X|. \tag{3}$$

For non-uniform spaces, the entropy can be evaluated by the well-known formula

$$\text{Ent}(X) = - \sum_{x \in S_X} p_X(x) \ln p_X(x).$$

1.2 Asymptotic equivalence

If X_1 and X_2 are probability spaces with the same entropy, there is a bijection between their typical sets of sequences of length n , for the plain reason that they can be chosen to have the same cardinality. It means that up to a change of code (of representation) and an error that becomes small as n gets large, the spaces X_1^n and X_2^n are equivalent. In the same sense, both X_1^n and X_2^n are equivalent to a uniform measure space with cardinality $e^{n\text{Ent}(X_i)}$.

In [10], Gromov formalized this concept of asymptotic equivalence. With his definition, two Bernoulli sequences of measure spaces X_1^n and X_2^n are asymptotically equivalent if there exists an “almost-measure-preserving” “almost-bijection”

$$X_1^n \xleftrightarrow{f} X_2^n$$

Even though we were greatly influenced by ideas in [10], we found that Gromov’s definition does not extend easily to situations in which *multiple signals* are processed at the same time, or when a dynamical system is probed with several measurement devices at once.

1.3 Diagrams of probability spaces

We model multiple signals by *diagrams* of probability spaces. By a diagram of probability spaces we mean a commutative diagram of probability spaces and measure-preserving maps between some of them. We will give a precise definition in Sect. 2.4, but will now consider particular examples of diagrams called *two-fans*.

A two-fan

$$X \xleftarrow{\pi_X} U \xrightarrow{\pi_Y} Y$$

is a triple of probability spaces $X = (\underline{X}, p_X), Y = (\underline{Y}, p_Y)$ and $U = (\underline{U}, p_U)$, and two measure-preserving maps π_X and π_Y . We will restrict ourselves for now to the case in which the underlying set of U is the Cartesian product of the underlying sets of X and Y , $\underline{U} = \underline{X} \times \underline{Y}$, and π_X and π_Y are the ordinary projections. Such a situation arises, for example, when a complex dynamical system, such as a living cell or a brain, is observed via two measuring devices.

Generalizing from the case of single signals, we might want to say that two-fans

$$\begin{aligned} X_1 &\longleftarrow U_1 \longrightarrow Y_1 \\ X_2 &\longleftarrow U_2 \longrightarrow Y_2 \end{aligned}$$

are asymptotically equivalent if for large n there exist almost measure-preserving, almost-bijections

$$\begin{array}{ccc} X_1^n & U_1^n & Y_1^n \\ \updownarrow f & \updownarrow h & \updownarrow g \\ X_2^n & U_2^n & Y_2^n \end{array}$$

Without additional assumptions, asymptotic equivalence classes for two-fans would be completely determined by the entropies of the constituent spaces.

However, such an asymptotic equivalence relation would be too coarse. Consider the three examples of two-fans are shown in Fig. 1, which is to be interpreted in the following way. Each of the spaces X_i and $Y_i, i = 1, 2, 3$, have cardinality six and a uniform distribution, where the weight of each atom is $\frac{1}{6}$. The spaces U_i have cardinality 12 and the distribution is also uniform with all weights being $\frac{1}{12}$. The

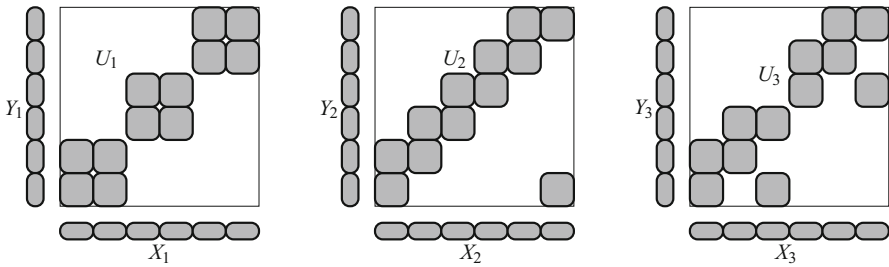


Fig. 1 Examples of pairs of probability spaces together with joint distributions

support of the measure on U_i 's is colored grey on the pictures. The maps from U_i to X_i and Y_i are coordinate projections.

In view of Eq. (3) we have for each $i = 1, 2, 3$,

$$\begin{aligned} \text{Ent}(X_i) &= \ln 6 \\ \text{Ent}(Y_i) &= \ln 6 \\ \text{Ent}(U_i) &= \ln 12. \end{aligned}$$

However, common information-processing techniques can still differentiate between the two-fans, by calculating solutions to information-optimization problems. This observation is sometimes expressed by saying that mutual information does not capture the full dependency structure that is relevant from an information-processing perspective. Information-optimization problems play an important role in information theory [25], causal inference [18], artificial intelligence [23], information decomposition [5], robotics [1], and neuroscience [9].

The additional assumption that the relations between the signals is relevant and should be preserved by the asymptotic equivalence results in the requirement that the following diagram commutes

$$\begin{array}{ccccc} X_1^n & \longleftarrow & U_1^n & \longrightarrow & Y_1^n \\ \uparrow f & & \uparrow h & & \uparrow g \\ X_2^n & \longleftarrow & U_2^n & \longrightarrow & Y_2^n \end{array}$$

However, with this generalization of an asymptotic equivalence to diagrams, we were not able to prove a corresponding Asymptotic Equipartition Property or even to prove the transitivity of the relation.

1.4 The entropy distances and asymptotic equivalence for diagrams

Instead of finding an almost measure-preserving bijection between large parts of the two spaces, we consider a stochastic coupling (transportation plan, joint distribution) between a pair of spaces and measure its deviation from being an isomorphism of probability spaces (a measure-preserving bijection). Such a measure of deviation from being an isomorphism then leads to the notion of intrinsic entropy distance, and its

stable version—the asymptotic entropy distance, as explained in Sect. 3. We say two sequences of diagrams are asymptotically equivalent if the asymptotic entropy distance between them vanishes.

The intrinsic entropy distance is an intrinsic version of a distance between random variables going by many different names, such as entropy distance, shared information distance and variation of information. It was reinvented many times by different people, among them Shannon [20], Kolmogorov, Sinai and Rokhlin. It appears in the proof of the theorem about generating partitions for ergodic systems by Kolmogorov and Sinai, see for example [22].

The intrinsic version of the entropy distance between probability spaces was introduced by Kovacevic et al. [14] and by Vidyasagar [24]. They showed that the involved minimization problem is NP-hard. Methods to find approximate solutions are discussed in [6,11].

1.5 Asymptotic equipartition property

With the notion of asymptotic equivalence induced by the asymptotic entropy distance, we *can* prove an asymptotic equipartition property for diagrams. Whereas the asymptotic equipartition property for single probability spaces states that high tensor powers of probability spaces can be approximated by *uniform* measure spaces, the Asymptotic Equipartition Property Theorem for diagrams, Theorem 6.1, states that sequences of successive tensor powers of a diagram can be approximated in the asymptotic entropy distance by a sequence of *homogeneous* diagrams.

Homogeneous diagrams have the property that the symmetry group acts transitively on the support of the measures of the constituent spaces. Two-fans shown on Fig. 1 are particular examples of homogeneous diagrams.

Homogeneous probability spaces are just uniform probability spaces, while homogeneous *diagrams* are, unlike homogeneous probability spaces, rather complex objects. Nonetheless, they seem to be simpler than arbitrary diagrams of probability spaces for the types of problems that we would like to address.

In a subsequent article we show that the optimal values in Information-Optimization problems only depend on the asymptotic class of a diagram and that they are continuous with respect to the asymptotic entropy distance; in many cases, the optimizers are continuous as well. The Asymptotic Equipartition Property implies that for the purposes of calculating optimal values and approximate optimizers, one only needs to consider homogeneous diagrams and this can greatly simplify computations.

Summarizing, the Asymptotic Equipartition Property and the continuity of Information-Optimization problems are important justifications for the choice of asymptotic equivalence relation and the introduction of the intrinsic and asymptotic Kolmogorov–Sinai distances.

1.6 Definitions and results in random variable context

In this article, we use the language of probability spaces and their commutative diagrams rather than the language of random variables, because we often encounter

situations in which their joint distributions are not defined, are variable, or even do not exist.

Some relations between the probability spaces can be easily represented by commutative diagrams of probability spaces, such as by a diamond diagram, Sect. 2.5.5, while the description with random variables is complex and not easily interpretable. The diagrams also provide a geometric overview of various entropy identities and inequalities.

Since the language of random variables will be more familiar to many readers, we now present our main result in these terms.

For random variables X, Y, Z etc., we denote by $\underline{X}, \underline{Y}, \underline{Z}$ the target sets, and by X, Y, Z the probability spaces with the induced distributions.

In general, there is a relation between k -tuples of random variables and diagrams of a certain type, involving a space for every subset $I \subset \{1, \dots, k\}$.

For example, a pair of random variables X, Y (defined on the same probability space) gives rise to a two-fan,

$$X \longleftarrow X \times Y \longrightarrow Y$$

where X, Y and $X \times Y$ are the target spaces of the random variables X, Y and (X, Y) endowed with their respective laws (i.e. the pushforward of the probability measure).

However, not every type of diagram corresponds to a tuple of random variables.

The entropy distance between two k -tuples $X = (X_1, \dots, X_k)$ and $Y = (Y_1, \dots, Y_k)$ is defined by

$$\text{kd}(X, Y) := \sum_{I \subset \{1, \dots, k\}} (2\text{Ent}(X_I Y_I) - \text{Ent}(X_I) - \text{Ent}(Y_I)).$$

A random k -tuple (H_1, \dots, H_k) is called homogeneous if for every two elements

$$\begin{aligned} (h_1, \dots, h_k) &\in \underline{H}_1 \times \dots \times \underline{H}_k \\ (h'_1, \dots, h'_k) &\in \underline{H}_1 \times \dots \times \underline{H}_k \end{aligned}$$

there exists a k -tuple of invertible maps $f_1 : \underline{H}_1 \rightarrow \underline{H}_1, \dots, f_k : \underline{H}_k \rightarrow \underline{H}_k$ such that $(f_1 \circ H_1, \dots, f_k \circ H_k)$ is equal to (H_1, \dots, H_k) in distribution and $f_i(h_i) = h'_i$. This condition is strictly stronger than the requirement that all the distributions are uniform.

Given n random k -tuples

$$X(i) = (X_1(i), \dots, X_k(i)), \quad i = 1, \dots, n$$

we can naturally construct the k -tuple

$$(X(1), \dots, X(n))$$

defined by

$$(X(1), \dots, X(n))_j = (X_j(1), \dots, X_j(n)), \quad j = 1, \dots, k.$$

It is difficult to formulate our main result, Theorem 6.1, in full generality using the language of random variables. However, the following theorem is an immediate corollary.

Theorem 1 *Let $(X(i) : i \in \mathbb{N})$ be a sequence of i.i.d. random tuples defined on a standard probability space.*

Define random k -tuples $Y(n)$ by

$$Y(n) := (X(1), \dots, X(n)).$$

Then, there exists a sequence of homogeneous random k -tuples $H(n) = (H_1(n), \dots, H_k(n))$, where $H_i(n)$ takes values in \underline{X}_i^n , such that

$$\frac{1}{n}kd(Y(n), H(n)) = O\left(n^{-1/2} \ln^{3/2}(n)\right).$$

2 Category of probability spaces and diagrams

In this section we present the basic setup used throughout the article. We will start by explaining how probability spaces and (equivalence classes) of measure-preserving maps between them form a *category*. This point of view on probability theory was already advocated in [2, 10].

Category theory yields simple definitions of diagrams of probability spaces and morphisms between them and allows for precise and relatively short proofs. The setup is also convenient when couplings (joint distributions) between probability spaces are absent or variable.

2.1 Categories

Below we briefly review elementary category theory. We refer the reader to the first chapter of [15] for a more extensive introduction.

A *category* \mathbf{C} is an abstract mathematical structure that captures the idea of a collection of spaces and structure-preserving maps between them, such as groups and homomorphisms, vector spaces and linear maps, and topological spaces and continuous maps. Categories consist of a collection of objects (which need not to be sets), a collection of morphisms (which need not to be maps), and a rule for composing morphisms.

More formally a category consists of

- A class of objects $\text{Obj}_{\mathbf{C}}$;
- A class of morphisms $\text{Hom}_{\mathbf{C}}(A, B)$ for every pair of objects $A, B \in \text{Obj}_{\mathbf{C}}$. For a morphism $f \in \text{Hom}_{\mathbf{C}}(A, B)$ one usually writes $f: A \rightarrow B$. Object A will be

called the domain and B the target of f , and we say that f is a morphism from A to B ;

- For each triple of objects A, B and C , a binary, associative operation, called composition,

$$\circ : \text{Hom}_{\mathbf{C}}(B, C) \times \text{Hom}_{\mathbf{C}}(A, B) \rightarrow \text{Hom}_{\mathbf{C}}(A, C)$$

$$(g, f) \mapsto g \circ f$$

- For every object $A \in \text{Obj}_{\mathbf{C}}$ an identity morphism $\mathbf{1}_A : A \rightarrow A$, with the property that for every $f : A \rightarrow B$ and every $g : B \rightarrow A$,

$$f \circ \mathbf{1}_A = f, \quad \mathbf{1}_A \circ g = g.$$

A morphism $f : A \rightarrow B$ is an *isomorphism* if there exists a morphism $g : B \rightarrow A$ such that $f \circ g = \mathbf{1}_B$ and $g \circ f = \mathbf{1}_A$.

Category theory becomes a very powerful tool when *functors* and their *natural transformations* are considered. Functors can be seen as homomorphisms between categories. In turn, *natural transformations* are homomorphisms between functors.

A (covariant) *functor* $\mathcal{X} : \mathbf{C} \rightarrow \mathbf{D}$ between two categories \mathbf{C} and \mathbf{D} , maps objects and morphisms in \mathbf{C} to objects and morphisms in \mathbf{D} , respectively. It satisfies the following additional properties: For every morphism $f : A \rightarrow B$ in \mathbf{C} the image $\mathcal{X}(f)$ is a morphism from $\mathcal{X}(A)$ to $\mathcal{X}(B)$ in \mathbf{D} and composition is preserved,

$$\mathcal{X}(g \circ f) = \mathcal{X}(g) \circ \mathcal{X}(f)$$

for any pair of morphisms $f : A \rightarrow B$ and $g : B \rightarrow C$.

A *natural transformation* between functors $\mathcal{X}, \mathcal{Y} : \mathbf{C} \rightarrow \mathbf{D}$ is a family η of morphisms in category \mathbf{D} , indexed by objects in \mathbf{C} : For every $A \in \text{Obj}_{\mathbf{C}}$, there is a morphism $\eta_A : \mathcal{X}(A) \rightarrow \mathcal{Y}(A)$, such that for every morphism $f : A \rightarrow B$ the diagram

$$\begin{array}{ccc} \mathcal{X}(A) & \xrightarrow{\mathcal{X}(f)} & \mathcal{X}(B) \\ \downarrow \eta_A & & \downarrow \eta_B \\ \mathcal{Y}(A) & \xrightarrow{\mathcal{Y}(f)} & \mathcal{Y}(B) \end{array}$$

commutes, that is

$$\eta_B \circ \mathcal{X}(f) = \mathcal{Y}(f) \circ \eta_A.$$

2.2 Probability spaces and reductions

We will now describe the category **Prob**. The objects in **Prob** are *finite probability spaces*. A finite probability space X is a pair (S, p) , where S is a (not necessarily finite) set and $p : 2^S \rightarrow [0, 1]$ is a probability measure, such that there is a *finite* subset of S with full measure. We denote by $\underline{X} = \text{supp } p$ the support of the measure and by $|X| := |\text{supp } p_X|$ its cardinality. Slightly abusing the language, we call this quantity

the *cardinality* of X . We will no longer explicitly mention that the probability spaces we consider are finite. We will also write p_X where we truly mean its density with respect to the counting measure.

We say that a map $f: X \rightarrow Y$ between two probability spaces X and Y is *measure-preserving* if the push-forward f_*p_X equals p_Y . This means that for every $A \subset Y$,

$$(f_*p_X)(A) := p_X(f^{-1}A) = p_Y(A).$$

We say that two measure-preserving maps $f: X \rightarrow Y$ are equivalent if they agree on a set of full measure. We call an equivalence class of measure-preserving maps from X to Y a *reduction*.

The morphisms in the category **Prob** are exactly the *reductions* between finite probability spaces. At this stage one might want to check that **Prob** is indeed a category, and this is guaranteed as the composition of two reductions is again a reduction.

2.3 Isomorphisms, automorphisms and homogeneity

Now that we have organized probability spaces and reductions into a category, we get concepts such as isomorphism for free: Two probability spaces X and Y are isomorphic in the category **Prob** if and only if there exists a measure-preserving bijection between the supports of the measures on X and Y . If X and Y are isomorphic, they have the same cardinality. The automorphism group $\text{Aut}(X)$ is the group of all self-isomorphisms of X .

A probability space X is called *homogeneous* if the automorphism group $\text{Aut}(X)$ acts transitively on the support \underline{X} of the measure. For the category **Prob**, this turns out to be a complicated way of saying that the measure on X is uniform on its support, but when we consider *diagrams* later, there will be no such simple implication. Homogeneity is an isomorphism invariant and we will denote the subcategory of homogeneous spaces by **Prob_h**.

There is a product in **Prob** (which is not a product in the sense of category theory!) given by the Cartesian product of probability spaces, that we will denote by $X \otimes Y := (\underline{X} \times \underline{Y}, p_X \otimes p_Y)$, where $p_X \otimes p_Y$ is the (independent) product measure. There are canonical reductions $X \otimes Y \rightarrow X$ and $X \otimes Y \rightarrow Y$ given by projections to factors. For a pair of reductions $f_i: X_i \rightarrow Y_i$, $i = 1, 2$ their tensor product is the reduction $f_1 \otimes f_2: X_1 \otimes X_2 \rightarrow Y_1 \otimes Y_2$, which is equal to the class of the Cartesian product of maps representing f_i 's. The product leaves the subcategory of homogeneous spaces invariant. If one of the factors in the product is replaced by an isomorphic space, then the product stays in the same isomorphism class.

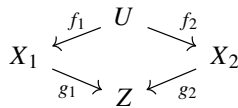
We close this section with a technical remark. The category **Prob** is not a *small category*. However it has a small full subcategory, that contains an object for every isomorphism class in **Prob** and for every pair of objects in it, it contains all the available morphisms between them and is closed under the product. From now on we imagine that such a subcategory was chosen and fixed and replaces **Prob** in all considerations below.

2.4 Diagrams of probability spaces

Essentially, a diagram $\mathcal{X} = \{X_i; f_{ij}\}$ is a commutative diagram in **Prob** consisting of a finite number of probability spaces and reductions between some of them. We have seen an example of the two-fan diagram in the introduction

$$X \xleftarrow{f} U \xrightarrow{g} Y$$

Another example is the diamond diagram



We require the diagram to be commutative, that is

$$g_1 \circ f_1 = g_2 \circ f_2.$$

A morphism $\rho: \mathcal{X} \rightarrow \mathcal{Y}$ between two diagrams $\mathcal{X} = \{X_i; f_{ij}\}$ and $\mathcal{Y} = \{Y_i; g_{ij}\}$ of the same shape is a collection of reductions between corresponding individual objects $\rho_i: X_i \rightarrow Y_i$, that commute with the reductions within each diagram, $\rho_j \circ f_{ij} = g_{ij} \circ \rho_i$.

We need to keep track of the combinatorial structure of the collection of reductions within a diagram. There are several possibilities for doing so:

- the reductions form a directed, acyclic graph which is transitively closed;
- the spaces in the diagram form a poset;
- the underlying combinatorial structure could be recorded as a finite category.

The last option seems to be most convenient since it has many operations, that are necessary for our analysis, already built-in. Besides, we need at times to iterate the construction of commutative diagrams, to create diagrams of diagrams, which is readily available in the category-theory framework but is cumbersome in the other contexts.

A (finite) poset category \mathbf{G} is a finite category such that for every two objects O_1 and O_2 there is at most one morphism between them in either direction:

$$|\text{Hom}_{\mathbf{G}}(O_1, O_2) \cup \text{Hom}_{\mathbf{G}}(O_2, O_1)| \leq 1.$$

For instance, the poset category \mathbf{A}_2 ,

$$O_1 \xleftarrow{\pi_1} O_{12} \xrightarrow{\pi_2} O_2$$

is a category with three objects $\{O_1, O_{12}, O_2\}$ and two non-identity morphisms $\pi_1: O_{12} \rightarrow O_1$ and $\pi_2: O_{12} \rightarrow O_2$.

A two-fan is then a diagram indexed by \mathbf{A}_2 : we assign to each object in \mathbf{A}_2 a probability space and to each morphism in \mathbf{A}_2 a reduction.

In general, then, a *diagram of probability spaces* indexed by a poset category \mathbf{G} is a functor $\mathcal{X} : \mathbf{G} \rightarrow \mathbf{Prob}$. The requirement that \mathcal{X} is a functor and not just a map between objects and morphisms (combined with the assumption that there is only one morphism between objects), is exactly the requirement that the diagrams should be commutative.

The collection of all diagrams of probability spaces indexed by a fixed poset category \mathbf{G} forms the so-called category of functors

$$\mathbf{Prob}\langle \mathbf{G} \rangle := [\mathbf{G}, \mathbf{Prob}].$$

The objects of $\mathbf{Prob}\langle \mathbf{G} \rangle$ are diagrams, that is functors from \mathbf{G} to \mathbf{Prob} , while morphisms in $\mathbf{Prob}\langle \mathbf{G} \rangle$ are natural transformations between them. We will refer to the morphisms in $\mathbf{Prob}\langle \mathbf{G} \rangle$ as reductions as well.

Let us go through the simple example of two-fans: we look at a reduction $\eta : \mathcal{X} \rightarrow \mathcal{Y}$ between two-fan diagrams $\mathcal{X}, \mathcal{Y} : \mathbf{A}_2 \rightarrow \mathbf{Prob}$. The reduction η , being a natural transformation between \mathcal{X} and \mathcal{Y} , is illustrated by the commutative diagram

$$\begin{array}{ccccc} \mathcal{X}(O_1) & \xleftarrow{\mathcal{X}(\pi_1)} & \mathcal{X}(O_{12}) & \xrightarrow{\mathcal{X}(\pi_2)} & \mathcal{X}(O_2) \\ \downarrow \eta(O_1) & & \downarrow \eta(O_{12}) & & \downarrow \eta(O_2) \\ \mathcal{Y}(O_1) & \xleftarrow{\mathcal{Y}(\pi_1)} & \mathcal{Y}(O_{12}) & \xrightarrow{\mathcal{Y}(\pi_2)} & \mathcal{Y}(O_2) \end{array}$$

Thus, a reduction of a two-fan is a family of reductions of probability spaces indexed by the objects in the poset category \mathbf{A}_2 such that the diagram commutes.

For a diagram $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{G} \rangle$, the poset category \mathbf{G} will be called the *combinatorial type* of \mathcal{X} . For a poset category \mathbf{G} or a diagram $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{G} \rangle$ we denote by $\llbracket \mathbf{G} \rrbracket = \llbracket \mathcal{X} \rrbracket$ the number of objects in the category \mathbf{G} .

An object O in a poset category \mathbf{G} will be called a *source*, if it is not a target of any morphism except for the identity. Likewise a *sink* object is not a domain of any morphism, except for the identity morphism. If a category contains a *unique* source object, the object is called the initial object and such a category will be called complete.

The above terminology transfers to diagrams indexed by \mathbf{G} : A source space in $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{G} \rangle$ is one that is not a target space of any reduction within the diagram, a sink space is not the domain of any non-trivial reduction and \mathcal{X} is called complete if \mathbf{G} is, i.e. if it has a unique source space.

The tensor product of probability spaces extends to a tensor product of diagrams. For $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}\langle \mathbf{G} \rangle$, such that $\mathcal{X} = \{X_i; f_{ij}\}$ and $\mathcal{Y} = \{Y_i; g_{ij}\}$ define

$$\mathcal{X} \otimes \mathcal{Y} := \{X_i \otimes Y_i; f_{ij} \otimes g_{ij}\}.$$

The construction of the category of commutative diagrams could be applied to any category, not just \mathbf{Prob} . Two additional cases will be of interest to us.

Denote by \mathbf{Set} the category of finite sets and surjective maps. Then all of the above constructions could be repeated for sets instead of probability spaces. Thus we could talk about the category of diagrams of sets $\mathbf{Set}\langle \mathbf{G} \rangle$.

Given a reduction $f : X \rightarrow Y$ between two probability spaces, the restriction $f : \underline{X} \rightarrow \underline{Y}$ is a well-defined surjective map. Given a diagram $\mathcal{X} = \{X_i; f_{ij}\}$ of

probability spaces, there is an underlying diagram of sets, obtained by taking the supports of measures on each level and restricting reductions on these supports. We will denote it by $\underline{\mathcal{X}} = \{ \underline{X}_i; \underline{f}_{ij} \}$, where $\underline{X}_i := \text{supp } p_{X_i}$. Thus we have a forgetful functor

$$\underline{\cdot} : \mathbf{Prob}\langle \mathbf{G} \rangle \rightarrow \mathbf{Set}\langle \mathbf{G} \rangle.$$

We could also repeat the construction of commutative diagrams to form a category of diagrams of diagrams. Thus given two poset categories \mathbf{G} and \mathbf{H} we can form a category $\mathbf{Prob}\langle \mathbf{G}, \mathbf{H} \rangle := \mathbf{Prob}\langle \mathbf{G} \rangle\langle \mathbf{H} \rangle$. We will rarely need anything beyond a two-fan of diagrams. There is a natural isomorphism

$$\mathbf{Prob}\langle \mathbf{G} \rangle\langle \mathbf{H} \rangle \equiv \mathbf{Prob}\langle \mathbf{H} \rangle\langle \mathbf{G} \rangle.$$

Thus, for example, a two-fan of \mathbf{G} -diagrams could be equivalently considered as a \mathbf{G} -diagram of two-fans, see also Sect. 2.5.3.

2.5 Examples of diagrams

We now consider some examples of poset categories and corresponding diagrams, that will be important in what follows.

2.5.1 Singleton

We denote by \bullet the poset category with a single object. Clearly diagrams indexed by \bullet are just probability spaces and we have $\mathbf{Prob} \equiv \mathbf{Prob}\langle \bullet \rangle$.

2.5.2 Chains

The chain \mathbf{C}_n of length $n \in \mathbb{N}$ is a category with n objects $\{O_i\}_{i=1}^n$ and morphisms from O_i to O_j whenever $i \geq j$. A diagram $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{C}_n \rangle$ is a chain of reductions

$$\mathcal{X} = (X_n \rightarrow X_{n-1} \rightarrow \cdots \rightarrow X_1).$$

2.5.3 Two-fan

The two-fan Λ_2 is a category with three objects $\{O_1, O_{12}, O_2\}$ and two non-identity morphisms $O_{12} \rightarrow O_1$ and $O_{12} \rightarrow O_2$. A diagram indexed by a two-fan will also be called a two-fan.

Essentially, a two-fan $(X \leftarrow Z \rightarrow Y)$ is a triple of probability spaces and a pair of reductions between them.

A reduction of a two-fan $\mathcal{F} = (X \leftarrow Z \rightarrow Y)$ to another two-fan $\mathcal{F}' = (X' \leftarrow Z' \rightarrow Y')$ is a triple of reductions $Z \rightarrow Z', Y \rightarrow Y'$ and $X \rightarrow X'$ that commute with the reductions within each fan, so that the diagram on Fig. 2a is commutative.

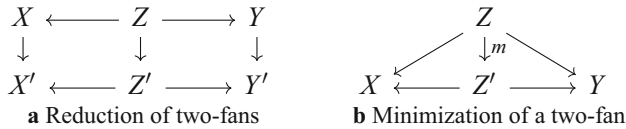


Fig. 2 Two-fans, their reductions and minimizations

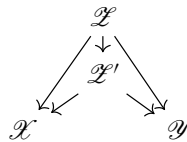
A two-fan $(X \leftarrow Z \rightarrow Y)$ is called *minimal* if for any super-diagram $\{X, Y, Z, Z'\}$ shown on Fig. 2b the reduction $m : Z \rightarrow Z'$ must be an isomorphism. Minimal two-fans are also called *couplings* in probability theory.

For any two-fan $(X \leftarrow Z \rightarrow Y)$ of probability spaces there always exist a unique (up to isomorphism), minimal two-fan $(X \leftarrow Z' \rightarrow Y)$, that can be included in the diagram shown on Fig. 2b. The minimization can be constructed by taking $Z' := \underline{X} \times \underline{Y}$ as a set and considering a probability distribution on Z' induced by a map $Z \rightarrow Z'$, that is the Cartesian product of the reductions $\underline{Z} \rightarrow \underline{X}$ and $\underline{Z} \rightarrow \underline{Y}$ in the original two-fan. Thus, the inclusion of a pair of probability spaces X and Y as sink vertices in a minimal two-fan is equivalent to specifying a joint distribution on $\underline{X} \times \underline{Y}$.

Note that minimality of a two-fan is defined in purely categorical terms. Even though the definition applies to two-fans of morphisms in any category, the minimization need not to exist. However as the next proposition asserts, if minimization of any two-fan exists in a category \mathbf{C} , then it also exists in a category of diagrams over \mathbf{C} .

Proposition 2.1 *Let \mathbf{G} be a poset category, and let $\mathcal{X} = \{X_i; a_{ij}\}$, $\mathcal{Y} = \{Y_i; b_{ij}\}$ and $\mathcal{Z} = \{Z_i; c_{ij}\}$ be three \mathbf{G} -diagrams. Then*

1. A two-fan $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}) \in \mathbf{Prob}(\mathbf{G}, \mathbf{A}_2)$ of \mathbf{G} -diagrams is minimal if and only if the constituent two-fans of probability spaces $\mathcal{F}_i = (X_i \leftarrow Z_i \rightarrow Y_i)$ are all minimal.
2. For any two-fan $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$ of \mathbf{G} -diagrams its minimal reduction exists, that is, there exists a minimal two-fan $\mathcal{F}' = (\mathcal{X} \leftarrow \mathcal{Z}' \rightarrow \mathcal{Y})$ included in the following diagram



The proof of Proposition 2.1 can be found on page 274.

2.5.4 Co-fan

A co-fan \mathbf{V} is a category with three objects and morphisms

$$\mathbf{V} = (O_1 \rightarrow O_\bullet \leftarrow O_2).$$

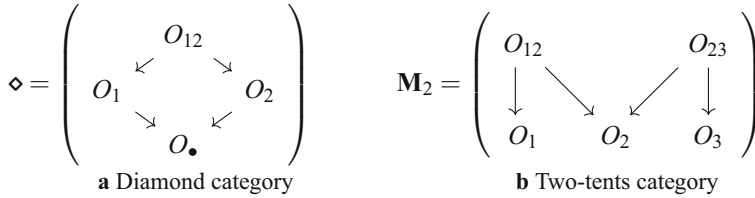


Fig. 3 Diamond and two-tents categories

2.5.5 A diamond diagram

A “diamond” diagram is indexed by a diamond category \diamond , that consists of a two-fan and a co-fan, as shown on Fig. 3.

Of course, there is also a morphism $O_{12} \rightarrow O_{\bullet}$, which lies in the transitive closure of the given four morphisms. We will often skip writing morphisms that are implied by the transitive closure.

A diamond diagram will be called *minimal* if the top two-fan in it is minimal.

2.5.6 “Two-tents” diagram

The “two-tents” category \mathbf{M}_2 consists of five objects, of which two are sources and three are sinks, and morphisms are as in Fig. 3b.

Thus, a typical two-tents diagram consists of five probability spaces and reductions as in

$$\mathcal{X} = (X \leftarrow U \rightarrow Y \leftarrow V \rightarrow Z).$$

The probability spaces U and V are sources and X, Y and Z are sinks.

2.5.7 Full diagram

The full category \mathbf{A}_n on n objects is a category with objects $\{O_I\}_{I \in 2^{\{1, \dots, n\}} \setminus \{\emptyset\}}$ indexed by all non-empty subsets $I \in 2^{\{1, \dots, n\}}$ and a morphism from O_I to O_J , whenever $J \subseteq I$.

A diagram \mathcal{X} indexed by a full category will be called *minimal*, if for every two-fan in it, it also contains a minimal two-fan with the same sink vertices. If $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{A}_n \rangle$ is minimal full diagram of probability spaces, then the set $\mathcal{X}(O_I)$ can be considered as a subset of the product $\prod_{i \in I} \mathcal{X}(O_i)$, while reductions are just coordinate projections.

For an n -tuple of random variables X_1, \dots, X_n one may construct a minimal full diagram $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{A}_n \rangle$ by considering all joint distributions and “marginalization” reductions. We denote such a diagram by $\langle X_1, \dots, X_n \rangle$. On the other hand, the reductions from the initial space to the sink vertices of a full diagram can be viewed as random variables on the domain of definition given by the (unique) initial space.

Suppose $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{A}_n \rangle$ is a minimal full diagram with sink vertices X_1, \dots, X_n . It is convenient to view \mathcal{X} as a distribution on the Cartesian product of the underlying sets of the sink vertices:

$$p_{\mathcal{X}} \in \Delta(\underline{X}_1 \times \cdots \times \underline{X}_n)$$

where ΔS stands for the space of all probability distribution on a finite set S .

Once the underlying sets of the sink spaces are fixed, there is a one-to-one correspondence between the full minimal diagrams and distributions as above.

As a corollary of Proposition 2.1 we also obtain the following characterization of minimal full diagrams of any \mathbf{G} -diagrams of probability spaces.

Corollary 2.2 *Let \mathbf{G} be an arbitrary poset category. Then*

1. *A full diagram \mathcal{F} of \mathbf{G} -diagrams is minimal, if and only if the constituent full diagrams of probability spaces \mathcal{F}_i are all minimal.*
2. *For any full diagram $\mathcal{F} \in \mathbf{Prob}\langle \mathbf{G}, \mathbf{A}_n \rangle$ of \mathbf{G} -diagrams there exists another minimal full diagram $\mathcal{F}' \in \mathbf{Prob}\langle \mathbf{G}, \mathbf{A}_n \rangle$ with the same sink entries and a reduction $\mu: \mathcal{F} \rightarrow \mathcal{F}'$, such that μ restricts to an isomorphism on sink entries of \mathcal{F} . Moreover, \mathcal{F}' is unique upto isomorphism.*

2.6 Constant diagrams

Suppose X is a probability space and \mathbf{G} is a poset category. One may form a *constant \mathbf{G} -diagram* by considering a functor that maps all objects in \mathbf{G} to X and all the morphisms to the identity morphism $X \xrightarrow{\text{Id}} X$. We denote such a constant diagram by $X^{\mathbf{G}}$ or simply by X , when \mathbf{G} is clear from the context. Any constant diagram is automatically minimal.

If $\mathcal{Y} = \{Y_i; f_{ij}\}$ is another \mathbf{G} -diagram, then a reduction $\rho: \mathcal{Y} \rightarrow X^{\mathbf{G}}$ (which we write sometimes simply as $\rho: \mathcal{Y} \rightarrow X$) is a collection of reductions $\rho_i: Y_i \rightarrow X$, such that

$$f_{ij} \circ \rho_i = \rho_j.$$

Let $\mathcal{X} = \{X_i; f_{ij}\}$ be a complete diagram with the initial space X_0 . Then there is a canonical reduction

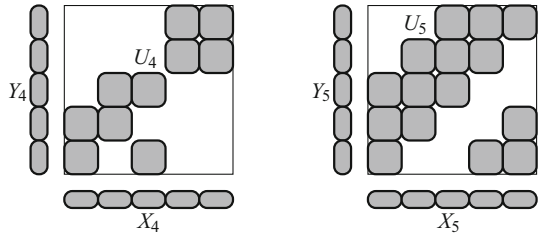
$$\rho: X_0^{\mathbf{G}} \rightarrow \mathcal{X}. \tag{4}$$

with components

$$\rho_i := f_{0i}.$$

By $\{\bullet\}$ we denote a one-point probability space. The constant \mathbf{G} -diagram $\{\bullet\}^{\mathbf{G}}$ is a unit with respect to the product in $\mathbf{Prob}\langle \mathbf{G} \rangle$.

Fig. 4 Non-homogeneous two-fans consisting of uniform spaces



2.7 Homogeneous diagrams

A diagram $\mathcal{X} \in \mathbf{Prob}(\mathbf{G})$ indexed by some poset category \mathbf{G} is called *homogeneous* if its automorphism group $\text{Aut}(\mathcal{X})$ acts transitively on every probability space in \mathcal{X} . Three examples of homogeneous diagrams were given in the introduction. The subcategory of all homogeneous diagrams indexed by \mathbf{G} will be denoted $\mathbf{Prob}(\mathbf{G})_h$.

In fact, for \mathcal{X} to be homogeneous it is sufficient that the $\text{Aut}(\mathcal{X})$ acts transitively on every source space in \mathcal{X} . Thus, if \mathcal{X} is complete with initial space X_0 , to check homogeneity it is sufficient to check the transitivity of the action of the symmetries of \mathcal{X} on X_0 .

Any subdiagram of a homogeneous diagram is also homogeneous. In particular, all the individual spaces of a homogeneous diagram are homogeneous

$$\mathbf{Prob}(\mathbf{G})_h \subset \mathbf{Prob}_h(\mathbf{G}).$$

However homogeneity of the whole of the diagram is a stronger property than homogeneity of the individual spaces in the diagram, thus in general

$$\mathbf{Prob}(\mathbf{G})_h \neq \mathbf{Prob}_h(\mathbf{G}).$$

Two examples of non-homogeneous two-fans are shown in Fig. 4. The pictures are to be interpreted in the same way as the pictures in Fig. 1.

A single probability space is homogeneous if and only if there is a representative in its isomorphism class with uniform measure and the same holds true for chain diagrams, for the co-fan or any other diagram that does not contain a two-fan. However, for more complex diagrams, for example for two-fans, no such simple description is available.

2.7.1 Universal construction of homogeneous diagrams

Examples of homogeneous diagrams could be constructed in the following manner. Suppose Γ is a finite group and $\{H_i\}$ is a collection of subgroups. Consider a collection of sets $\underline{X}_i := \Gamma/H_i$ and consider a natural surjection $f_{ij} : \underline{X}_i \rightarrow \underline{X}_j$ whenever H_i is a subgroup of H_j . Equipping each \underline{X}_i with the uniform distribution one can turn the diagram of sets $\{\underline{X}_i; f_{ij}\}$ into a homogeneous diagram of probability spaces. It will be complete if there is a smallest subgroup (under inclusion) among H_i 's.

Such a diagram will be complete and minimal, if together with any pair of groups H_i and H_j in the collection, their intersection $H_i \cap H_j$ also belongs to the collection $\{H_i\}$.

In fact, any homogeneous diagram arises this way. Suppose diagram $\mathcal{X} = \{X_i; f_{ij}\}$ is homogeneous, then we set $\Gamma = \text{Aut}(\mathcal{X})$ and choose a collection of points $x_i \in X_i$ such that $f_{ij}(x_i) = x_j$ and denote by $H_i := \text{Stab}(x_i) \subset \Gamma$. Then, if one applies the construction of the previous paragraph to Γ , with the collection of subgroups $\{H_i\}$, one recovers the original diagram \mathcal{X} upto isomorphism.

2.8 Conditioning

Suppose a diagram \mathcal{X} contains a fan

$$\mathcal{F} = \left(X \xleftarrow{f} Z \xrightarrow{g} Y \right).$$

Given a point $x \in X$ with a non-zero weight one may consider *conditional probability distributions* $p_Z(\cdot | x)$ on \underline{Z} , and $p_Y(\cdot | x)$ on \underline{Y} . The distribution $p_Z(\cdot | x)$ is supported on $f^{-1}(x)$ and is defined by the property that for any function $f : Z \rightarrow \mathbb{R}$ holds

$$\int_Z f d p_Z = \int_X \left[\int_Z f(z) d p_Z(z|x) \right] d p_X(x)$$

and is given by

$$p_Z(z|x) = \frac{p_Z(z)}{p_X(x)}.$$

The distribution $p_Y(\cdot | x)$ is the pushforward of $p_Z(\cdot | x)$ under g

$$p_Y(\cdot | x) = g_* p_Z(\cdot | x).$$

Recall that if \mathcal{F} is minimal, the underlying set of Z can be assumed to be the product $\underline{X} \times \underline{Y}$. In that case

$$p_Y(y|x) = \frac{p_Z(x, y)}{p_X(x)}.$$

We denote the corresponding space $Y|x := (\underline{Y}, p_Y(\cdot | x))$.

Under some assumptions it is possible to condition a whole sub-diagram of \mathcal{X} . More specifically, if a diagram \mathcal{X} contains a sub-diagram \mathcal{Y} and a probability space X satisfying the condition that *there exists a space Z in \mathcal{X} that reduces to all the spaces in \mathcal{Y} and to X* , then we may condition the whole of \mathcal{Y} on $x \in X$ given that $p_X(x) > 0$.

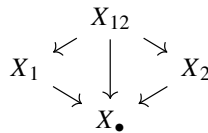
For $x \in X$ with positive weight we denote by $\mathcal{Y} \lfloor x$ the diagram of spaces in \mathcal{Y} conditioned on $x \in X$. The diagram $\mathcal{Y} \lfloor x$ has the same combinatorial type as \mathcal{Y} and will be called the *slice* of \mathcal{Y} over $x \in X$. Note that the space X itself may or may not belong to \mathcal{Y} . The conditioning $\mathcal{Y} \lfloor x$ may depend on the choice of a fan between \mathcal{Y} and X , however when \mathcal{X} is complete the conditioning $\mathcal{Y} \lfloor x$ is well-defined and is independent of the choice of fans.

Suppose now that there are two subdiagram \mathcal{Y} and \mathcal{Z} in \mathcal{X} and in addition \mathcal{Z} is a constant diagram, $\mathcal{Z} = \mathbf{Z}^{\mathbf{G}'}$ for some poset category \mathbf{G}' . Let $z \in \mathbf{Z}$, then $\mathcal{Y} \lfloor z$ is well defined and is independent of the choice of the space in \mathcal{Z} , the element of which z is to be considered.

If \mathcal{X} is homogeneous, then $\mathcal{Y} \lfloor x$ is also homogeneous and its isomorphism class does not depend on the choice of $x \in \underline{X}$.

2.9 Entropy

We define entropy by the limit in Eq. 2. Entropy satisfies the so-called Shannon inequality, see for example [8]. Namely for any minimal diamond diagram



the following inequality holds,

$$\text{Ent}(X_1) + \text{Ent}(X_2) \geq \text{Ent}(X_{12}) + \text{Ent}(X_{\bullet}). \tag{5}$$

Furthermore, entropy is additive with respect to the tensor product, that is, for a pair of probability spaces $X, Y \in \mathbf{Prob}$ holds

$$\text{Ent}(X \otimes Y) = \text{Ent}(X) + \text{Ent}(Y). \tag{6}$$

Conditional entropy $\text{Ent}(X \lfloor Y)$ is defined for a pair X, Y of probability spaces included in a minimal two-fan ($X \leftarrow Z \rightarrow Y$) as

$$\text{Ent}(X \lfloor Y) := \text{Ent}(Z) - \text{Ent}(X).$$

The above quantity is always non-negative in view of Shannon inequality (5). Moreover, the following identity holds, see [8]

$$\text{Ent}(X \lfloor Y) = \int_Y \text{Ent}(X \lfloor y) d p_Y(y). \tag{7}$$

For a \mathbf{G} -diagram $\mathcal{X} = \{X_i; f_{ij}\}$ define the entropy homomorphism

$$\text{Ent}_*: \mathbf{Prob}(\mathbf{G}) \rightarrow \mathbb{R}^{\llbracket \mathbf{G} \rrbracket}, \quad \{X_i; f_{ij}\} \mapsto (\text{Ent}(X_i)) \in \mathbb{R}^{\llbracket \mathbf{G} \rrbracket}.$$

It will be convenient for us to equip the target $\mathbb{R}^{[\mathbf{G}]}$ with the ℓ^1 -norm. Thus

$$|\text{Ent}_*(\mathcal{X})|_1 = \sum_{i=1}^{[\mathbf{G}]} \text{Ent}(X_i).$$

If \mathcal{X} is a complete \mathbf{G} -diagram with initial space X_0 , then by Shannon inequality (5) there is an obvious estimate

$$\text{Ent}(X_0) \leq |\text{Ent}_*(\mathcal{X})|_1 \leq \llbracket \mathcal{X} \rrbracket \cdot \text{Ent}(X_0).$$

3 The entropy distance

We turn the space of diagrams into a pseudo-metric space by introducing the intrinsic entropy distance and asymptotic entropy distance. The intrinsic entropy distance is obtained by taking an infimum of the entropy distance over all possible joint distributions on two probability spaces.

3.1 Entropy distance and asymptotic entropy distance

3.1.1 Entropy distance in the case of single probability spaces

For a two-fan $\mathcal{F} = (X \leftarrow Z \rightarrow Y)$ define a “distance” $\text{kd}(\mathcal{F})$ between probability spaces X and Y with respect to \mathcal{F} by

$$\begin{aligned} \text{kd}(\mathcal{F}) &:= \text{Ent}(Z|Y) + \text{Ent}(Z|X) \\ &= 2\text{Ent}(Z) - \text{Ent}(X) - \text{Ent}(Y). \end{aligned}$$

If a two-fan \mathcal{F} satisfies $\text{kd}(\mathcal{F}) = 0$, then both reductions in \mathcal{F} are isomorphisms. Thus, essentially $\text{kd}(\mathcal{F})$ is some measure of the deviation of the statistical map defined by \mathcal{F} from being a deterministic bijection between X and Y .

The minimal reduction \mathcal{F}' of \mathcal{F} satisfies

$$\text{kd}(\mathcal{F}') \leq \text{kd}(\mathcal{F}). \tag{8}$$

For a pair of probability spaces X, Y define the *intrinsic entropy distance* as

$$\mathbf{k}(X, Y) := \inf \{ \text{kd}(\mathcal{F}) : \mathcal{F} = (X \leftarrow Z \rightarrow Y) \text{ is a two-fan} \}. \tag{9}$$

The optimization takes place over all two-fans with sink spaces X and Y . In view of inequality (8) one could as well optimize over the space of *minimal* two-fans, which we will also refer to as *couplings* between X and Y . The tensor product of X and Y trivially provides a coupling and the set of couplings is compact, therefore an optimum is always achieved and it is finite.

The bivariate function $\mathbf{k}: \mathbf{Prob} \times \mathbf{Prob} \rightarrow \mathbb{R}_{\geq 0}$ defines a notion of pseudo-distance and it vanishes exactly on pairs of isomorphic probability spaces. This follows directly from the Shannon inequality (5), and a more general statement will be proven in Proposition 3.1 below.

3.1.2 Entropy distance for complete diagrams

The definition of entropy distance for complete diagrams repeats almost literally the definition for single spaces. We fix a complete poset category \mathbf{G} and will be considering diagrams from $\mathbf{Prob}(\mathbf{G})$.

Consider three such diagrams $\mathcal{X} = \{X_i, f_{ij}\}$, $\mathcal{Y} = \{Y_i, g_{ij}\}$ and $\mathcal{Z} = \{Z_i, h_{ij}\}$ from $\mathbf{Prob}(\mathbf{G})$. Recall that a two-fan $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$ can also be viewed as a \mathbf{G} -diagram of two-fans

$$\mathcal{F}_i = (X_i \leftarrow Z_i \rightarrow Y_i).$$

Define

$$\begin{aligned} \text{kd}(\mathcal{F}) &:= \sum_i \text{kd}(\mathcal{F}_i) \\ &= \sum_i (2\text{Ent}(Z_i) - \text{Ent}(X_i) - \text{Ent}(Y_i)). \end{aligned}$$

The quantity $\text{kd}(\mathcal{F})$ vanishes if and only if the fan \mathcal{F} provides isomorphisms between all individual spaces in \mathcal{X} and \mathcal{Y} that commute with the inner structure of the diagrams, that is, it provides an isomorphism between \mathcal{X} and \mathcal{Y} in $\mathbf{Prob}(\mathbf{G})$.

The *intrinsic entropy distance between diagrams* is defined in analogy with the case of single probability spaces

$$\mathbf{k}(\mathcal{X}, \mathcal{Y}) := \inf \{ \text{kd}(\mathcal{F}) : \mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}) \},$$

where the infimum is over all two-fans of \mathbf{G} -diagrams with sink vertices \mathcal{X} and \mathcal{Y} .

The following proposition records that the intrinsic entropy distance is in fact a pseudo-distance on $\mathbf{Prob}(\mathbf{G})$, provided \mathbf{G} is a complete poset category (that is when \mathbf{G} has a unique initial space).

Proposition 3.1 *Let \mathbf{G} be a complete poset category. Then the bivariate function*

$$\mathbf{k}: \mathbf{Prob}(\mathbf{G}) \times \mathbf{Prob}(\mathbf{G}) \rightarrow \mathbb{R}$$

is a pseudo-distance on $\mathbf{Prob}(\mathbf{G})$.

Moreover, two diagrams $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}(\mathbf{G})$ satisfy $\mathbf{k}(\mathcal{X}, \mathcal{Y}) = 0$ if and only if \mathcal{X} is isomorphic to \mathcal{Y} in $\mathbf{Prob}(\mathbf{G})$.

The idea of the proof is very simple. In the case of single probability spaces X, Y, Z a coupling between X and Z can be constructed from a coupling between X and Y and

a coupling between Y and Z by *adhesion* on Y , see [16]. The triangle inequality then follows from Shannon inequality. However, since we are dealing with diagrams the combinatorial structure requires careful treatment. Therefore, we provide a detailed proof on page 276.

It is important to note, that the proof uses the fact that \mathbf{G} is complete. In fact, even though the definition of \mathbf{k} could be easily extended to some bivariate function on the space of diagrams of any fixed combinatorial type, it fails to satisfy the triangle inequality in general, because the composition of couplings requires completeness of \mathbf{G} .

3.1.3 The asymptotic entropy distance

Let \mathbf{G} be a complete poset category. We will show in Corollary 3.5 below, that the sequence

$$n \mapsto \mathbf{k}(\mathcal{X}^n, \mathcal{Y}^n)$$

is sublinear and therefore the following limit exists.

$$\kappa(\mathcal{X}, \mathcal{Y}) := \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{k}(\mathcal{X}^n, \mathcal{Y}^n). \tag{10}$$

We call its value, $\kappa(\mathcal{X}, \mathcal{Y})$, the *asymptotic entropy distance* between two diagrams $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}(\mathbf{G})$.

As a corollary of Proposition 3.1 and definition (10) we immediately obtain that the asymptotic entropy distance is a homogeneous pseudo-distance on $\mathbf{Prob}(\mathbf{G})$.

Corollary 3.2 *Let \mathbf{G} be a complete poset category. Then the bivariate function*

$$\kappa: \mathbf{Prob}(\mathbf{G}) \times \mathbf{Prob}(\mathbf{G}) \rightarrow \mathbb{R}$$

is a pseudo-distance on $\mathbf{Prob}(\mathbf{G})$ satisfying the following properties: for any pair of diagrams $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}(\mathbf{G})$

1. $\kappa(\mathcal{X}, \mathcal{Y}) \leq \mathbf{k}(\mathcal{X}, \mathcal{Y})$
2. *for any $n \in \mathbb{N}_0$ holds $\kappa(\mathcal{X}^n, \mathcal{Y}^n) = n \cdot \kappa(\mathcal{X}, \mathcal{Y})$.*

We will see later that there are instances when $\kappa < \mathbf{k}$, moreover there are pairs of non-isomorphic diagrams with vanishing asymptotic entropy distance between them.

In the next subsection we derive some elementary properties of the intrinsic entropy distance and the asymptotic entropy distance.

3.2 Properties of (asymptotic) entropy distance

3.2.1 Tensor product

We show that the tensor product on the space of diagrams is 1-Lipschitz. Later this will allow us to give a simple description of tropical diagrams, that is, of points in the

asymptotic cone of $\mathbf{Prob}\langle \mathbf{G} \rangle$, as limits of certain sequences of “classical” diagrams, as will be discussed in a subsequent article.

Proposition 3.3 *Let \mathbf{G} be a complete poset category. Then with respect to the Kolmogorov distance on $\mathbf{Prob}\langle \mathbf{G} \rangle$ the tensor product*

$$\otimes : (\mathbf{Prob}\langle \mathbf{G} \rangle, k)^2 \rightarrow (\mathbf{Prob}\langle \mathbf{G} \rangle, k)$$

is 1-Lipschitz in each variable, that is, for every triple $\mathcal{X}, \mathcal{Y}, \mathcal{Y}' \in \mathbf{Prob}\langle \mathbf{G} \rangle$ the following bound holds

$$k(\mathcal{X} \otimes \mathcal{Y}, \mathcal{X} \otimes \mathcal{Y}') \leq k(\mathcal{Y}, \mathcal{Y}').$$

This statement is a direct consequence of additivity of entropy with respect to the tensor product. Details can be found on page 279.

It follows directly from definition (10) and Proposition 3.3, that the asymptotic entropy distance enjoys a similar property.

Corollary 3.4 *Let \mathbf{G} be a complete poset category. Then with respect to the asymptotic entropy distance on $\mathbf{Prob}\langle \mathbf{G} \rangle$ the tensor product*

$$\otimes : (\mathbf{Prob}\langle \mathbf{G} \rangle, \kappa)^2 \rightarrow (\mathbf{Prob}\langle \mathbf{G} \rangle, \kappa)$$

is 1-Lipschitz in each variable.

As another corollary we obtain the subadditivity properties of the intrinsic entropy distance and asymptotic entropy distance.

Corollary 3.5 *Let \mathbf{G} be a complete poset category and let $\mathcal{X}, \mathcal{Y}, \mathcal{U}, \mathcal{V} \in \mathbf{Prob}\langle \mathbf{G} \rangle$, then*

$$k(\mathcal{X} \otimes \mathcal{U}, \mathcal{Y} \otimes \mathcal{V}) \leq k(\mathcal{X}, \mathcal{Y}) + k(\mathcal{U}, \mathcal{V})$$

and

$$\kappa(\mathcal{X} \otimes \mathcal{U}, \mathcal{Y} \otimes \mathcal{V}) \leq \kappa(\mathcal{X}, \mathcal{Y}) + \kappa(\mathcal{U}, \mathcal{V}).$$

It implies in particular that shifts are non-expanding maps in $(\mathbf{Prob}\langle \mathbf{G} \rangle, k)$ or $(\mathbf{Prob}\langle \mathbf{G} \rangle, \kappa)$.

Corollary 3.6 *Let \mathbf{G} be a complete poset category and $\delta = k, \kappa$ be either intrinsic entropy distance or asymptotic entropy distance on $\mathbf{Prob}\langle \mathbf{G} \rangle$. Let $\mathcal{U} \in \mathbf{Prob}\langle \mathbf{G} \rangle$. Then the shift map*

$$\mathcal{U} \otimes \cdot : (\mathbf{Prob}\langle \mathbf{G} \rangle, \delta) \rightarrow (\mathbf{Prob}\langle \mathbf{G} \rangle, \delta), \quad \mathcal{X} \mapsto \mathcal{U} \otimes \mathcal{X}$$

is a non-expanding map with respect to either intrinsic entropy distance or asymptotic entropy distance.

Less obvious is the fact that κ is, in fact, translation invariant and in particular, $(\mathbf{Prob}(\mathbf{G}), \kappa)$ satisfies the cancellation property. This is the subject of Proposition 3.7 below, which was communicated to us by Tobias Fritz.

Proposition 3.7 *For any triple of diagrams $\mathcal{X}, \mathcal{Y}, \mathcal{U}$ holds*

$$\kappa(\mathcal{X} \otimes \mathcal{U}, \mathcal{Y} \otimes \mathcal{U}) = \kappa(\mathcal{X}, \mathcal{Y}).$$

The proof of the lemma can be found on page 280.

3.2.2 Entropy

Recall that we defined the entropy function

$$\mathbf{Ent}_* : \mathbf{Prob}(\mathbf{G}) \rightarrow \mathbb{R}^{[\mathbf{G}]}$$

by evaluating the entropy of all individual spaces in a \mathbf{G} -diagram. The target space $\mathbb{R}^{[\mathbf{G}]}$ will be endowed with the ℓ^1 -norm with respect to the natural coordinate system. With such a choice, the entropy function is 1-Lipschitz with respect to the Kolmogorov distance on $\mathbf{Prob}(\mathbf{G})$.

Proposition 3.8 *Suppose \mathbf{G} is a complete poset category and $\delta = k, \kappa$ is either intrinsic entropy distance or asymptotic entropy distance on $\mathbf{Prob}(\mathbf{G})$. Then the entropy function*

$$\mathbf{Ent}_* : (\mathbf{Prob}(\mathbf{G}), \delta) \rightarrow (\mathbb{R}^{[\mathbf{G}]}, |\cdot|_1), \quad \mathcal{X} = \{X_i, f_{ij}\} \mapsto (\mathbf{Ent} X_i)_i \in \mathbb{R}^{[\mathbf{G}]}$$

is 1-Lipschitz.

Again, the proof of the proposition above is an application of Shannon’s inequality, see page 281 for details.

3.3 The Slicing Lemma

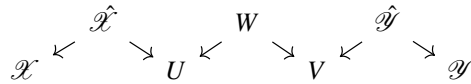
The Slicing Lemma, Proposition 3.9 below, allows to estimate the intrinsic entropy distance between two diagrams with the integrated intrinsic entropy distance between “slices”, which are diagrams obtained by conditioning on another probability space. It turned out to be a very powerful tool for estimation of the intrinsic entropy distance and will be used below on several occasions.

As described in Sect. 2.6, by a reduction of a diagram $\mathcal{X} = \{X_i, f_{ij}\}$ to a single space U we mean a collection of reductions $\{\rho_i : X_i \rightarrow U\}$ from the individual spaces in \mathcal{X} to U , that commute with the reductions within \mathcal{X}

$$\rho_j \circ f_{ij} = \rho_i.$$

Alternatively, whenever a single probability space appears as a domain or a target of a morphism to or from a \mathbf{G} -diagram, it should be replaced by a constant \mathbf{G} -diagram.

Proposition 3.9 (Slicing Lemma) *Suppose \mathbf{G} is a complete poset category and we are given $\mathcal{X}, \hat{\mathcal{X}}, \mathcal{Y}, \hat{\mathcal{Y}} \in \mathbf{Prob}(\mathbf{G})$ —four \mathbf{G} -diagrams and $U, V, W \in \mathbf{Prob}$ —three probability spaces, that are included into the following three-tents diagram*



such that the two-fan $(U \leftarrow W \rightarrow V)$ is minimal. Then the following estimate holds

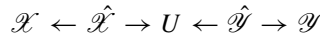
$$\begin{aligned} \mathbf{k}(\mathcal{X}, \mathcal{Y}) &\leq \int_W \mathbf{k}(\mathcal{X}[u], \mathcal{Y}[v]) d p_W(u, v) \\ &\quad + \llbracket \mathbf{G} \rrbracket \cdot kd(U \leftarrow W \rightarrow V) \\ &\quad + \sum_i [\text{Ent}(U[X_i]) + \text{Ent}(V[Y_i])]. \end{aligned}$$

The idea of the proof of the Slicing Lemma (page 281) is as follows. For every pair $(u, v) \in \underline{W}$ we consider an optimal two-fan \mathcal{G}_{uv} coupling $\mathcal{X}[u]$ and $\mathcal{Y}[v]$. These fans have the same underlying diagram of sets. Then we construct a coupling between \mathcal{X} and \mathcal{Y} as a convex combination of distributions of \mathcal{G}_{uv} 's weighted by $p_W(u, v)$. The estimates on the resulting two-fan then imply the proposition.

Various implications of the Slicing Lemma are summarized in the next corollary.

Corollary 3.10 *Let \mathbf{G} be a complete poset category, $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}(\mathbf{G})$ and $U \in \mathbf{Prob}$.*

1. Given a “two-tents” diagram



the following inequality holds

$$\mathbf{k}(\mathcal{X}, \mathcal{Y}) \leq \int_U \mathbf{k}(\mathcal{X}[u], \mathcal{Y}[u]) d p_U(u) + 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot \text{Ent}(U).$$

2. Given a fan



the following inequality holds

$$\mathbf{k}(\mathcal{X}, \mathcal{Y}) \leq \int_{UV} \mathbf{k}(\mathcal{X}[u], \mathcal{Y}) d p_U(u) + 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot \text{Ent}(U).$$

3. Let $\mathcal{X} \rightarrow U$ be a reduction, then

$$\mathbf{k}(\mathcal{X}, \mathcal{Y}) \leq \int_U \mathbf{k}(\mathcal{X}[u], \mathcal{Y}) d p_U(u) + \llbracket \mathbf{G} \rrbracket \cdot \text{Ent}(U).$$

4. For a co-fan $\mathcal{X} \rightarrow U \leftarrow \mathcal{Y}$ holds

$$k(\mathcal{X}, \mathcal{Y}) \leq \int_U k(\mathcal{X}|u, \mathcal{Y}|u) d p_U(u).$$

4 Distributions and types

In this section we recall some elementary inequalities for (relative) entropies and the total variation distance for distributions on finite sets. Furthermore, we generalize the notion of a probability distribution on a set to a distribution on a diagram of sets. Finally, we give a perspective on the theory of types, and also introduce types in the context of complete diagrams.

4.1 Distributions

4.1.1 Single probability spaces

For a finite set S we denote by ΔS the collection of all probability distributions on S . It is a unit simplex in the real vector space \mathbb{R}^S . We often use the fact that it is a compact, convex set, whose interior points correspond to fully supported probability measures on S .

For $\pi_1, \pi_2 \in \Delta S$ denote by $|\pi_1 - \pi_2|_1$ the total variation of the signed measure $(\pi_1 - \pi_2)$ and define the entropy of the distribution π_1 by

$$h(\pi_1) := - \sum_{x \in X} \pi_1(x) \ln \pi_1(x). \quad (11)$$

If, in addition, π_2 lies in the interior of ΔS define the relative entropy by

$$D(\pi_1 || \pi_2) := \sum_{x \in X} \pi_1(x) \ln \frac{\pi_1(x)}{\pi_2(x)}.$$

The entropy of a probability space is often defined through formula (11). It is a standard fact, and can be verified with the help of Lemma 4.2 below, that for $\pi \in \Delta S$ holds

$$h(\pi) = \text{Ent}(S, \pi) \quad (12)$$

which justifies the name “entropy” for the function $h: \Delta S \rightarrow \mathbb{R}$.

Define a *divergence ball* of radius $\varepsilon > 0$ centered at $\pi \in \text{Interior} \Delta S$ as

$$B_\varepsilon(\pi) := \{ \pi' \in \Delta S : D(\pi' || \pi) \leq \varepsilon \}. \quad (13)$$

For a fixed π and $\varepsilon \ll 1$ the ball $B_\varepsilon(\pi)$ also lies in the interior of ΔS .

The total variation norm and relative entropy are related by the following inequality.

Lemma 4.1 *Let S be a finite set, then for any $\pi_1, \pi_2 \in \Delta S$, Pinsker’s inequality holds*

$$|\pi_1 - \pi_2|_1 \leq \sqrt{2D(\pi_1 || \pi_2)}.$$

The claim of the Lemma, Pinsker’s inequality, is a well-known inequality in for instance information theory, and a proof can be found in [8].

4.1.2 Distributions on diagrams

A map $f : S \rightarrow S'$ between two finite sets induces an affine map $f_* : \Delta S \rightarrow \Delta S'$.

For a diagram of sets $\mathcal{S} = \{S_i; f_{ij}\}$ we define the *space of distributions on the diagram \mathcal{S}* by

$$\Delta \mathcal{S} := \left\{ (\pi_i) \in \prod_i \Delta S_i : (f_{ij})_* \pi_i = \pi_j \right\}.$$

Essentially, an element of $\Delta \mathcal{S}$ is a collection of distributions on the sets S_i in \mathcal{S} that is consistent with respect to the maps f_{ij} . The consistency conditions $(f_{ij})_* \pi_i = \pi_j$ form a collection of linear equations with integer coefficients with respect to the standard convex coordinates in $\prod \Delta S_i$. Thus, $\Delta \mathcal{S}$ is a rational affine subspace in the product of simplices. In particular, $\Delta \mathcal{S}$ has a convex structure.

If \mathcal{S} is complete with initial set S_0 , then specifying a distribution $\pi_0 \in \Delta S_0$ uniquely determines distributions on all of the S_i ’s by setting $\pi_i := (f_{0i})_* \pi_0$. In such a situation we have

$$\Delta \mathcal{S} \cong \Delta S_0.$$

If \mathcal{S} is not complete and S_0, \dots, S_k is a collection of its source sets, then $\Delta \mathcal{S}$ is isomorphic to an affine subspace of the product $\Delta S_0 \times \dots \times \Delta S_k$ cut out by linear equations with integer coefficients corresponding to co-fans in \mathcal{S} with source sets among S_0, \dots, S_k .

To simplify notation, for a probability space X or a diagram \mathcal{X} we will write

$$\begin{aligned} \Delta X &:= \underline{\Delta X} \\ \Delta \mathcal{X} &:= \underline{\Delta \mathcal{X}}. \end{aligned}$$

4.2 Types

We now discuss briefly the theory of types. Types are special subspaces of tensor powers that consist of sequences with the same “empirical distribution” as explained in details below. For a more detailed discussion the reader is referred to [7,8]. We generalize the theory of types to complete diagrams of sets and complete diagrams of probability spaces.

The theory of types for diagrams, that are not complete, is more complex and will be addressed in a subsequent article.

4.2.1 Types for single probability spaces

Let S be a finite set. For $n \in \mathbb{N}$ denote also

$$\Delta^{(n)} S := \Delta S \cap \frac{1}{n} \mathbb{Z}^S$$

a collection of rational points in ΔS with denominator n . (We say that a rational number $r \in \mathbb{Q}$ has denominator $n \in \mathbb{N}$ if $r \cdot n \in \mathbb{Z}$)

Define the *empirical distribution map* $\mathbf{q} : S^n \rightarrow \Delta S$, that sends $(s_i)_{i=1}^n = \mathbf{s} \in S^n$ to the empirical distribution $\mathbf{q}(\mathbf{s}) \in \Delta S$ given by

$$\mathbf{q}(\mathbf{s})(a) = \frac{1}{n} \cdot |\{i : s_i = a\}| \quad \text{for any } a \in S.$$

Clearly the image of \mathbf{q} lies in $\Delta^{(n)} S$.

For $\pi \in \Delta^{(n)} S$, the space $T_\pi^n S := \mathbf{q}^{-1}(\pi)$ equipped with the uniform measure is called a *type* over π . The symmetric group \mathbb{S}_n acts on S^n by permuting the coordinates. This action leaves the empirical distribution invariant and therefore could be restricted to each type, where it acts transitively. Thus, for $\pi \in \Delta^{(n)} S$ the probability space $(T_\pi^n S, u)$ with u being a uniform (\mathbb{S}_n -invariant) distribution, is a homogeneous space.

Suppose $X = (\underline{X}, p)$ is a probability space. Let τ_n be the pushforward of $p^{\otimes n}$ under the empirical distribution map $\mathbf{q} : X^n \rightarrow \Delta X$. Clearly $\text{supp } \tau_n \subset \Delta^{(n)} X$, thus $(\Delta X, \tau_n)$ is a finite probability space. Therefore we have a reduction

$$\mathbf{q} : X^n \rightarrow (\Delta X, \tau_n) \tag{14}$$

which we call the *empirical reduction*. If $\pi \in \Delta^{(n)} X$ is such that $\tau_n(\pi) > 0$, then

$$T_\pi^n \underline{X} = X^n \lfloor \pi.$$

In particular, it follows that the right-hand side does not depend on the probability p on X as long as π is “compatible” to it.

The following lemma records some standard facts about types, which can be checked by elementary combinatorics and found in [8].

Lemma 4.2 *Let X be a probability space and $\mathbf{x} \in X^n$, then*

1. $|\Delta^{(n)} X| = \binom{n + |X|}{|X|} \leq e^{|X| \cdot \ln(n+1)} = e^{O(|X| \cdot \ln n)}$
2. $p^{\otimes n}(\mathbf{x}) = e^{-n[h(\mathbf{q}(\mathbf{x})) + D(\mathbf{q}(\mathbf{x}) || p)]}$
3. $e^{-n \cdot h(\pi) - |X| \cdot \ln(n+1)} \leq |T_\pi^n \underline{X}| \leq e^{n \cdot h(\pi)}$ or $|T_\pi^n \underline{X}| = e^{-n \cdot h(\pi) + O(|X| \cdot \ln n)}$
4. $e^{-n \cdot D(\pi || p) - |X| \cdot \ln(n+1)} \leq \tau_n(\pi) = p^{\otimes n}(T_\pi^n \underline{X}) \leq e^{-n \cdot D(\pi || p)}$ or $\tau_n(\pi) = e^{-n \cdot D(\pi || p) + O(|X| \cdot \ln n)}$

If $X = (\underline{X}, p_X)$ is a probability space with rational probability distribution with denominator n , then the type over p_X will be called the *true type* of X

$$T^n X := T_{p_X}^n \underline{X}.$$

As a corollary to Lemma 4.2 and equation (12) we obtain the following.

Corollary 4.3 *For a finite set S and $\pi \in \Delta^{(n)} S$ holds*

$$n \cdot h(\pi) - |S| \cdot \ln(n + 1) \leq \text{Ent}(T_\pi^n S) \leq n \cdot h(\pi).$$

Also, for a finite probability space $X = (S, p)$ with a rational distribution p with denominator n holds

$$n \cdot \text{Ent}(X) - |X| \cdot \ln(n + 1) \leq \text{Ent}(T^n X) \leq n \cdot \text{Ent}(X).$$

In particular,

$$\text{Ent}(T^n X) = n \cdot \text{Ent}(X) + o(|X| \cdot n).$$

The following important theorem is known as Sanov’s theorem. It can be easily derived from Lemma 4.2 or a proof can be found in [8].

Theorem 4.4 (Sanov’s Theorem) *Let $X = (S, p)$ be a finite probability space and let $\mathbf{q}: X^n \rightarrow (\Delta X, \tau_n)$ be the empirical reduction. Then for every $r > 0$,*

$$\tau_n(\Delta X \setminus B_r(p)) \leq e^{-n \cdot r + |X| \cdot \ln(n+1)}$$

where $B_r(p)$ is the divergence ball (relative entropy ball) defined in (13).

Combining the estimate in Theorem 4.4 with the Pinsker’s inequality in 4.1 we obtain the following corollary.

Corollary 4.5 *For a finite probability space $X = (S, p)$ holds*

$$\tau_n(\{\pi : |\pi - p|_1 \geq r\}) \leq e^{-\frac{1}{2}n \cdot r^2 + O(|X| \cdot \ln n)}.$$

4.3 Types for complete diagrams

In this subsection we generalize the theory of types for diagrams indexed by a complete poset category. The theory for a non-complete diagrams is more complex and will be addressed in our future work. Before we describe our approach we need some preparatory material.

Suppose we have a reduction $f : X \rightarrow Y$ between a pair of probability spaces. Then for any $n \in \mathbb{N}$ there is an induced reduction $f_* : (\Delta X, \tau_n) \rightarrow (\Delta Y, \tau_n)$ that can be included in the following diamond diagram

$$\begin{array}{ccc} X^n & \xrightarrow{f^n} & Y^n \\ \downarrow \mathbf{q} & & \downarrow \mathbf{q} \\ (\Delta X, \tau_n) & \xrightarrow{f_*} & (\Delta Y, \tau_n) \end{array}$$

that satisfies certain special condition, namely, the sides of the diamond are independent conditioned on the bottom space

$$\Delta X \perp\!\!\!\perp Y^n \downarrow \Delta Y.$$

In particular, for any $\pi \in \Delta X$ with $\tau_n(\pi) > 0$ and $\pi' = f_*\pi \in \Delta Y$ holds

$$Y^n \downarrow \pi = Y^n \downarrow \pi' \tag{15}$$

and there is a well-defined reduction

$$Tf : T_\pi^n X \rightarrow T_{\pi'}^n Y$$

for any $\pi \in \Delta^{(n)} X$ and $\pi' = f_*\pi \in \Delta^{(n)} Y$.

Now we are ready to give the definitions of types. Let $\mathcal{X} \in \mathbf{Prob}\langle \mathbf{G} \rangle$ be a complete diagram, $\mathcal{X} = \{X_i; f_{ij}\}$ with initial space X_0 and let $\pi \in \Delta^{(n)} \mathcal{X}$.

Define the type $T_\pi^n \mathcal{X}$ as the \mathbf{G} -diagram, whose individual spaces are types of the individual spaces of \mathcal{X} over the corresponding push-forwards of π

$$T_\pi^n \mathcal{X} := \{T_{\pi_i}^n X_i; Tf_{ij}\}.$$

Consider a symmetric group \mathbb{S}_n acting on \mathcal{X}^n by automorphisms permuting the coordinates. The action leaves the types $T_\pi^n \mathcal{X}$ invariant and it is transitive on the initial space $T_\pi^n X_0$. Thus, each type $T_\pi^n \mathcal{X}$ is a homogeneous diagram.

4.3.1 The empirical two-fan

Unlike in the cases of single probability spaces there is no empirical reduction from the power of \mathcal{X} to $\Delta \mathcal{X}$. It will be convenient for us to see the types as the power of the diagram conditioned on a distribution. This is achieved by including the power of diagram into a *empirical two-fan*.

Given a \mathbf{G} -diagram \mathcal{X} with initial space X_0 we construct the associated empirical two-fan with sink vertices \mathcal{X}^n and $(\Delta \mathcal{X}, \tau_n)^{\mathbf{G}}$ as the “composition” of the canonical reduction $(X_0)^{\mathbf{G}} \rightarrow \mathcal{X}$, Eq. (4) in Sect. 2.6, and the empirical reduction $X_0^n \rightarrow \Delta X_0 \cong \Delta \mathcal{X}$ in Eq. (14).

$$\mathcal{Q}_n(\mathcal{X}) = \left(\begin{array}{ccc} & (X_0^{\mathbf{G}})^{(n)} & \\ f_{0*}^n \swarrow & & \searrow \mathbf{q}^{\mathbf{G}} \\ \mathcal{X}^n & & (\Delta^{(n)} \mathcal{X}, \tau_n)^{\mathbf{G}} \end{array} \right) \tag{16}$$

The two-fan \mathcal{Q}_n is not necessarily minimal, but its minimal reduction can be constructed using Lemma 2.2 on page 252.

Let $\pi_0 \in (\Delta \mathcal{X}, \tau_n)$ with $\tau_n(\pi) > 0$ and $\pi_i = f_{0i} \pi_0$. Then within \mathcal{Q}_n holds

$$T_{\pi}^n X_i = X_i^n \lfloor \pi_i = X_i^n \lfloor \pi_0$$

by Eq. (15) and therefore

$$\mathcal{X}^n \lfloor \pi = T_{\pi}^n \mathcal{X}.$$

For every $n \in \mathbb{N}$ and $\pi \in \Delta^{(n)} X_0$ the type $T_{\pi}^n \mathcal{X}$ is a homogeneous diagram. Suppose that a complete diagram \mathcal{X} is such that the probability distribution p_0 on the initial set is rational with the denominator n , then we call $T_{p_0}^n \mathcal{X}$ the *true type* of \mathcal{X} and denote

$$T^n \mathcal{X} := T_{p_0}^n \mathcal{X}.$$

5 Distance between types

Our goal in this section is to estimate the intrinsic entropy distance between two types over two different distributions $\pi_1, \pi_2 \in \Delta^{(n)} \mathcal{S}$ in terms of the total variation distance $|\pi_1 - \pi_2|_1$.

For this purpose we use a “lagging” technique which is explained below. Practically, we couple different types by randomly removing and inserting the appropriate amount of symbols to pass from a trajectory of the one type to a trajectory of the other.

5.1 The lagging trick

Let Λ_{α} be a binary probability space,

$$\Lambda_{\alpha} := \left(\{ \blacksquare, \blacksquare \}; p_{\Lambda_{\alpha}}(\blacksquare) = \alpha \right)$$

and let $\mathcal{X} = \{(\underline{X}_i, p_i); f_{ij}\}$, $\mathcal{Z} = \{(\underline{Z}_i, q_i); g_{ij}\}$ be two diagrams indexed by a complete poset category \mathbf{G} and included in a minimal two-fan, i.e a coupling,

$$(\Lambda_\alpha)\mathbf{G} \xleftarrow{\lambda} \mathcal{Z} \xrightarrow{\rho} \mathcal{X}.$$

Assume further that the distribution q on \mathcal{Z} is rational with denominator $n \in \mathbb{N}$, that is $q \in \Delta^{(n)} \underline{\mathcal{Z}}$. It follows that p and p_{Λ_α} are also rational with the same denominator n .

We construct a *lagging two-fan*

$$\mathcal{L} := (T^{(1-\alpha)n}(\mathcal{X}[\square]) \xleftarrow{l} T^n \mathcal{Z} \xrightarrow{T\rho} T^n \mathcal{X}) \tag{17}$$

as follows. The right leg $T\rho$ of \mathcal{L} is induced by the right leg ρ of the original two-fan. The left leg

$$l: T^n \mathcal{Z} \rightarrow T^{(1-\alpha)n}(\mathcal{X}[\square])$$

is obtained by erasing symbols that reduce to \blacksquare and applying ρ to the remaining symbols. The target space for the reduction l is the true type of $\mathcal{X}[\square]$ which is “lagging” behind $T^n \mathcal{Z}$ by a factor of $(1 - \alpha)$. More specifically, the reduction l is constructed as follows.

Let $\lambda_j: Z_j \rightarrow \Lambda_\alpha$ be the components of the reduction $\lambda: \mathcal{Z} \rightarrow \Lambda_\alpha$. Given $\bar{z} = (z_i)_{i=1}^n \in T^n Z_j$ define the subset of indices

$$I_{\bar{z}} := \{i : \lambda_j(z_i) = \square\}$$

and define the j th component of l by

$$l_j((z_i)_{i=1}^n) := (\rho(z_i))_{i \in I_{\bar{z}}}.$$

By equivariance each l_j is a reduction of homogeneous spaces, since the inverse image of any point has the same cardinality. Moreover the reductions l_j commute with the reductions in $T^n \mathcal{Z}$ as explained in Sect. 4.3 and therefore l is a reduction of diagrams.

The next lemma uses the lagging two-fan to estimate the intrinsic entropy distance between its sink diagrams.

Lemma 5.1 *Let $\mathcal{X}, \mathcal{Z} \in \mathbf{Prob}(\mathbf{G})$ be two diagrams indexed by a complete poset category \mathbf{G} and included in a minimal two-fan*

$$\Lambda_\alpha \xleftarrow{\lambda} \mathcal{Z} \xrightarrow{\rho} \mathcal{X}$$

where distribution on \mathcal{Z} is rational with denominator $n \in \mathbb{N}$. Then

$$k(T^{(1-\alpha)n}(\mathcal{X}[\square]), T^n \mathcal{X})$$

$$\begin{aligned} &\leq n \cdot \llbracket G \rrbracket \cdot [2\text{Ent}(\Lambda_\alpha) + \alpha \cdot \ln |X_0|] + 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot |X_0| \cdot \ln(n + 1) \\ &= n \cdot \llbracket G \rrbracket \cdot [2\text{Ent}(\Lambda_\alpha) + \alpha \cdot \ln |X_0|] + O(|X_0| \cdot \ln n). \end{aligned}$$

It is an immediate consequence of the Slicing Lemma, in particular Corollary 3.10 part (2) that

$$\mathbf{k}(\mathcal{X}[\square], \mathcal{X}) \leq \llbracket G \rrbracket \cdot [2\text{Ent}(\Lambda_\alpha) + \alpha \cdot \ln |X_0|].$$

By the subadditivity of the intrinsic entropy distance,

$$\mathbf{k}((\mathcal{X}[\square]^{\otimes n}, \mathcal{X}^{\otimes n})) \leq n \cdot \llbracket G \rrbracket \cdot [2\text{Ent}(\Lambda_\alpha) + \alpha \cdot \ln |X_0|].$$

This bound is almost the estimate in Lemma 5.1, except Lemma 5.1 estimates the distance between types rather than tensor powers. We will soon see that tensor powers and types are very close in the intrinsic entropy distance. However, for the purpose of the proof of Lemma 5.1, it suffices to know that their entropies are close, an estimate that is provided by Corollary 4.3.

Proof of Lemma 5.1 We will use the lagging two-fan constructed in Eq. (17), namely

$$\mathcal{L} := (T^{(1-\alpha)n}(\mathcal{X}[\square]) \xleftarrow{l} T^n \mathcal{X} \xrightarrow{T\rho} T^n \mathcal{X})$$

as a coupling to estimate the intrinsic entropy distance

$$\mathbf{k}(T^{(1-\alpha)n}(\mathcal{X}[\square]), T^n \mathcal{X}) \leq \text{kd}(\mathcal{L}).$$

Recall that by Corollary 4.3 for a probability space X with a rational distribution we have

$$n \cdot \text{Ent}(X) - |X| \cdot \ln(n + 1) \leq \text{Ent}(T^n X) \leq n \cdot \text{Ent}(X).$$

Thus we can estimate $\text{kd}(\mathcal{L})$ as follows

$$\begin{aligned} \text{kd}(\mathcal{L}) &= \sum_i \left[(\text{Ent}(T^n Z_i) - \text{Ent}(T^n X_i)) \right. \\ &\quad \left. + (\text{Ent}(T^n Z_i) - \text{Ent}(T^{(1-\alpha)n}(X_i[\square]))) \right] \\ &\leq n \cdot \sum_i \left[(\text{Ent}(Z_i) - \text{Ent}(X_i)) + (\text{Ent}(Z_i) - (1 - \alpha)\text{Ent}(X_i[\square])) \right] \\ &\quad + 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot |X_0| \cdot \ln(n + 1). \end{aligned}$$

By minimality of the original two-fan and Shannon inequality (5) we have a bound

$$\text{Ent}(Z_i) - \text{Ent}(X_i) \leq \text{Ent}(\Lambda_\alpha).$$

The second part in the sum can be estimated using relation (7) as follows

$$\begin{aligned} \text{Ent}(Z_i) - (1 - \alpha)\text{Ent}(X_i|\square) &= \text{Ent}(\Lambda_\alpha) + \text{Ent}(X_i|\Lambda_\alpha) - (1 - \alpha)\text{Ent}(X_i|\square) \\ &= \text{Ent}(\Lambda_\alpha) + (1 - \alpha)\text{Ent}(X_i|\square) + \alpha\text{Ent}(X_i|\blacksquare) - (1 - \alpha)\text{Ent}(X_i|\square) \\ &\leq \text{Ent}(\Lambda_\alpha) + \alpha \cdot \ln |X_i|. \end{aligned}$$

Combining all of the above we obtain the estimate in the conclusion of the lemma. \square

5.2 Distance between types

In this section we use the lagging trick as described above to estimate the distance between types over two different distributions in $\Delta\mathcal{S}$ where \mathcal{S} is a complete diagram of sets.

Proposition 5.2 *Suppose \mathcal{S} is a complete \mathbf{G} -diagram of sets with initial set S_0 . Suppose $p, q \in \Delta^{(n)}\mathcal{S}$ and let $\alpha = \frac{1}{2}|p_0 - q_0|_1$. Then*

$$\begin{aligned} k(T_p^n\mathcal{S}, T_q^n\mathcal{S}) &\leq 2n \cdot \|\mathbf{G}\| \cdot [\alpha \cdot \ln |S_0| + 2\text{Ent}(\Lambda_\alpha)] + 4\|\mathbf{G}\| \cdot |X_0| \cdot \ln(n + 1) \\ &= 2n \cdot \|\mathbf{G}\| \cdot [\alpha \cdot \ln |S_0| + 2\text{Ent}(\Lambda_\alpha)] + O(|X_0| \cdot \ln n). \end{aligned}$$

The idea of the proof is to write p and q as a convex combination of a common distribution \hat{p} and “small amounts” of p^+ and q^+ , respectively. Then we use the lagging trick to estimate distances between types over p and \hat{p} , as well as between types over q and \hat{p} . We now present details of the proof.

Proof of Proposition 5.2 Recall that for a complete diagram \mathcal{S} with initial set S_0 we have

$$\Delta\mathcal{S} \cong \Delta S_0. \tag{18}$$

Our goal now is to write p and q as the convex combination of three other distributions \hat{p} , p^+ and q^+ as in

$$\begin{aligned} p &= (1 - \alpha) \cdot \hat{p} + \alpha \cdot p^+ \\ q &= (1 - \alpha) \cdot \hat{p} + \alpha \cdot q^+. \end{aligned}$$

We could do it the following way. Let $\alpha := \frac{1}{2}|p_0 - q_0|_1$. If $\alpha = 1$ then the proposition follows trivially by constructing a tensor-product fan, so from now on we assume that $\alpha < 1$. Define three probability distributions \hat{p}_0 , p_0^+ and q_0^+ on S_0 by setting for every $x \in S_0$

$$\begin{aligned} \hat{p}_0(x) &:= \frac{1}{1 - \alpha} \min \{p_0(x), q_0(x)\} \\ p_0^+ &:= \frac{1}{\alpha}(p_0 - (1 - \alpha)\hat{p}_0) \end{aligned}$$

$$q_0^+ := \frac{1}{\alpha}(q_0 - (1 - \alpha)\hat{p}_0).$$

Denote by $\hat{p}, p^+, q^+ \in \Delta\mathcal{S}$ the distributions corresponding to $\hat{p}_0, p_0^+, q_0^+ \in \Delta S_0$ under the affine isomorphism (18). Thus we have

$$\begin{aligned} p &= (1 - \alpha) \cdot \hat{p} + \alpha \cdot p^+ \\ q &= (1 - \alpha) \cdot \hat{p} + \alpha \cdot q^+. \end{aligned}$$

Now we construct a pair of two-fans of **G**-diagrams

$$\Lambda_\alpha \longleftarrow \tilde{\mathcal{X}} \longrightarrow \mathcal{X} \quad \text{and} \quad \Lambda_\alpha \longleftarrow \tilde{\mathcal{Y}} \longrightarrow \mathcal{Y} \tag{19}$$

by setting

$$\begin{aligned} \mathcal{X} &:= (\mathcal{S}, p) \\ \mathcal{Y} &:= (\mathcal{S}, q) \\ \tilde{X}_i &:= \left(S_i \times \underline{\Delta}_\alpha; \tilde{p}_i(s, \square) = (1 - \alpha)\hat{p}_i(s), \tilde{p}_i(s, \blacksquare) = \alpha \cdot p_i^+(s) \right) \\ \tilde{Y}_i &:= \left(S_i \times \underline{\Delta}_\alpha; \tilde{q}_i(s, \square) = (1 - \alpha)\hat{p}_i(s), \tilde{q}_i(s, \blacksquare) = \alpha \cdot q_i^+(s) \right) \\ \tilde{\mathcal{X}} &:= \left\{ \tilde{X}_i; f_{ij} \times \text{Id} \right\} \\ \tilde{\mathcal{Y}} &:= \left\{ \tilde{Y}_i; f_{ij} \times \text{Id} \right\}. \end{aligned}$$

The reductions in (19) are given by coordinate projections. We have the following isomorphisms

$$\mathcal{X} \llcorner \square \cong \mathcal{Y} \llcorner \square \cong (\mathcal{S}, \hat{p}).$$

To estimate the distance between types we now apply Lemma 5.1 to the fans in (19)

$$\begin{aligned} \mathbf{k}(T_p^n \mathcal{S}, T_q^n \mathcal{S}) &= \mathbf{k}(T^n \mathcal{X}, T^n \mathcal{Y}) \\ &\leq \mathbf{k}(T^n \mathcal{X}, T^{(1-\alpha)n}(\mathcal{X} \llcorner \square)) + \mathbf{k}(T^{(1-\alpha)n}(\mathcal{Y} \llcorner \square), T^n \mathcal{Y}) \\ &\leq 2n \cdot \llbracket \mathbf{G} \rrbracket \cdot [\alpha \cdot \ln |S_0| + 2\text{Ent}(\Lambda_\alpha)] + 4\llbracket \mathbf{G} \rrbracket \cdot |X_0| \cdot \ln(n + 1). \end{aligned}$$

□

6 Asymptotic equipartition property for diagrams

Below we prove that any Bernoulli sequence of complete diagrams can be approximated by a sequence of homogeneous diagrams. This is essentially the *Asymptotic Equipartition Theorem for diagrams*.

Theorem 6.1 *Suppose $\mathcal{X} \in \mathbf{Prob}(\mathbf{G})$ is a complete diagram of probability spaces. Then there exists a sequence $\tilde{\mathcal{H}} = (\mathcal{H}_n)_{n=0}^\infty$ of homogeneous diagrams of the same combinatorial type as \mathcal{X} such that for all $n \geq |X_0|$*

$$\frac{1}{n} \mathbf{k}(\mathcal{X}^{\otimes n}, \mathcal{H}_n) \leq C(|X_0|, \llbracket \mathbf{G} \rrbracket) \cdot \sqrt{\frac{\ln^3 n}{n}} \tag{20}$$

where X_0 is the initial space of \mathcal{X} and $C(|X_0|, \llbracket \mathbf{G} \rrbracket)$ is a constant only depending on $|X_0|$ and $\llbracket \mathbf{G} \rrbracket$.

Proof Denote by $\mathcal{S} = \underline{\mathcal{X}}$ the underlying diagram of sets and by $p_{\mathcal{X}}$ the true distribution on \mathcal{S} , such that

$$\mathcal{X} = (\mathcal{S}, p_{\mathcal{X}}).$$

We will construct the approximating homogeneous sequence by taking types over rational approximations of $p_{\mathcal{X}}$ in $\Delta \mathcal{S}$, that converge sufficiently fast to the true distribution $p_{\mathcal{X}}$.

More specifically, we select rational distributions $p_n \in \Delta^{(n)} \mathcal{S}$ such that

$$|p_n - p_{\mathcal{X}}|_1 \leq \frac{|S_0|}{n}.$$

As homogeneous spaces \mathcal{H}_n we set $\mathcal{H}_n = T_{p_n}^n \mathcal{S}$. We will show that the intrinsic entropy distance between \mathcal{H}_n and \mathcal{X}^n satisfies the required estimate (20).

First we apply slicing along the empirical two-fan

$$\mathcal{Q}_n(\mathcal{X}) = \left(\mathcal{X}^n \leftarrow \tilde{\mathcal{X}}^{(n)} \rightarrow (\Delta \mathcal{S}, \tau_n)^{\mathbf{G}} \right)$$

defined in Sect. 4.3, Eq. (16) on page 267.

For the estimate below we use the fact that

$$\text{Ent}(\Delta \mathcal{S}, \tau_n) \leq \ln |\Delta^{(n)} \mathcal{S}| \leq |S_0| \cdot \ln(n + 1).$$

By slicing (see Corollary 3.10(2)) along the empirical two-fan we have

$$\begin{aligned} \mathbf{k}(T_{p_n}^n \mathcal{S}, \mathcal{X}^{\otimes n}) &\leq 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot \text{Ent}(\Delta \mathcal{S}, \tau_n) + \int_{\Delta \mathcal{S}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d \tau_n(\pi) \\ &\leq 2 \cdot \llbracket \mathbf{G} \rrbracket \cdot |S_0| \cdot \ln(n + 1) + \int_{\Delta \mathcal{S}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d \tau_n(\pi). \end{aligned}$$

To estimate the integral we split the domain into a small divergence ball $B_{\varepsilon_n} = B_{\varepsilon_n}(p_{\mathcal{X}})$ around the “true” distribution and its complement

$$\int_{\Delta \mathcal{S}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d \tau_n(\pi) = \int_{\Delta \mathcal{S} \setminus B_{\varepsilon_n}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d \tau_n(\pi) \tag{21}$$

$$+ \int_{B_{\varepsilon_n}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d\tau_n(\pi) \tag{22}$$

and we set the radius ε_n equal to

$$\varepsilon_n := (|S_0| + 1) \frac{\ln(n + 1)}{n}.$$

To estimate the first integral on the right-hand side of equality (22) note that the distance between two types over the same diagram of sets can always be crudely estimated by

$$2 \cdot \ln |S_0| \cdot \llbracket \mathbf{G} \rrbracket \cdot n.$$

Moreover, by Sanov’s theorem, Theorem 4.1, we can estimate the empirical measure of the complement of the divergence ball

$$\tau_n(\Delta \mathcal{S} \setminus B_{\varepsilon_n}) \leq e^{-n \cdot \varepsilon_n + |S_0| \cdot \ln(n+1)} \leq \frac{1}{n}$$

where we used the definition of ε_n to conclude the last inequality. Therefore we obtain

$$\begin{aligned} \int_{\Delta \mathcal{S} \setminus B_{\varepsilon_n}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d\tau_n(\pi) &\leq 2 \cdot \ln |S_0| \cdot \llbracket \mathbf{G} \rrbracket \cdot n \cdot \tau_n(\Delta \mathcal{S} \setminus B_{\varepsilon_n}) \\ &\leq 2 \cdot \ln |S_0| \cdot \llbracket \mathbf{G} \rrbracket. \end{aligned}$$

Define

$$\alpha_n = \frac{|S_0|}{n} + \sqrt{2\varepsilon_n}$$

if the right-hand side is smaller than 1 and set $\alpha_n = 1$ otherwise. Then every $\pi \in B_{\varepsilon_n}(p_{\mathcal{X}})$ satisfies $|p_n - \pi| \leq 2\alpha_n$ by Pinsker’s inequality (Lemma 4.1), and the triangle inequality. Consequently, by the estimate on the distance between types in Proposition 5.2

$$\begin{aligned} \int_{B_{\varepsilon_n}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d\tau_n(\pi) \\ \leq 2n \cdot \llbracket \mathbf{G} \rrbracket \cdot (\alpha_n \ln |S_0| + 2\text{Ent}(\Lambda_{\alpha_n})) + 4 \cdot \llbracket \mathbf{G} \rrbracket \cdot |S_0| \cdot \ln(n + 1). \end{aligned}$$

Using the definition of α_n and ε_n we find that

$$\int_{B_{\varepsilon_n}} \mathbf{k}(T_{p_n}^n \mathcal{S}, T_{\pi}^n \mathcal{S}) d\tau_n(\pi) = O\left(\sqrt{n \cdot \ln^3 n}\right)$$

and hence combining the above estimates

$$\frac{1}{n} \mathbf{k}(T_{p_n}^n \mathcal{S}, \mathcal{X}^{\otimes n}) = \mathcal{O} \left(\sqrt{\frac{\ln^3 n}{n}} \right).$$

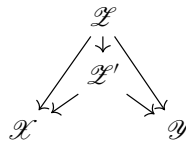
A more precise check shows that for $n \geq |S_0|$, the constants appearing in \mathcal{O} only depend on $|S_0|$ and $\mathbb{[G]}$. □

7 Technical proofs

This section contains some proofs that did not make it into the main text. The numbering of the claims in this section coincides with the numbering in the main text. Lemma that first appear in this section are numbered within section.

Proposition 2.1 *Let \mathbf{G} be a poset category, and let $\mathcal{X} = \{X_i; a_{ij}\}$, $\mathcal{Y} = \{Y_i; b_{ij}\}$ and $\mathcal{Z} = \{Z_i; c_{ij}\}$ be three \mathbf{G} -diagrams. Then*

1. *A two-fan $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}) \in \mathbf{Prob}(\mathbf{G}, \mathbf{A}_2)$ of \mathbf{G} -diagrams is minimal if and only if the constituent two-fans of probability spaces $\mathcal{F}_i = (X_i \leftarrow Z_i \rightarrow Y_i)$ are all minimal.*
2. *For any two-fan $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$ of \mathbf{G} -diagrams its minimal reduction exists, that is, there exists a minimal two-fan $\mathcal{F}' = (\mathcal{X} \leftarrow \mathcal{Z}' \rightarrow \mathcal{Y})$ included in the following diagram*



Before we go to the proof of Proposition 2.1, we will need the following lemma.

Lemma 7.1 *Suppose we are given a pair of two-fans of probability spaces*

$$\begin{aligned} \mathcal{F} &= (X \xleftarrow{\alpha} Z \xrightarrow{\beta} Y) \\ \mathcal{F}'' &= (X'' \xleftarrow{\alpha''} Z'' \xrightarrow{\beta''} Y'') \end{aligned}$$

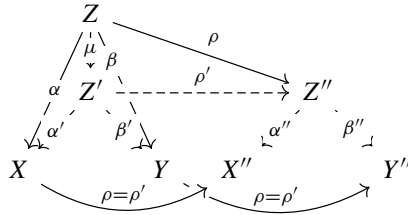
such that \mathcal{F}'' is minimal. Let

$$\mathcal{F} \xrightarrow{\mu} \mathcal{F}' = (X \xleftarrow{\alpha'} Z' \xrightarrow{\beta'} Y)$$

be a minimal reduction of \mathcal{F} . Then for any reduction $\rho : \mathcal{F} \rightarrow \mathcal{F}''$, there exists a reduction $\rho' : \mathcal{F}' \rightarrow \mathcal{F}''$ such that $\rho = \rho' \circ \mu$.

Proof of Lemma 7.1 We define ρ' on the sink spaces of \mathcal{F}' to coincide with ρ .

To prove the lemma we just need to provide a dashed arrow that makes the following diagram commutative



The reduction ρ' is constructed by simple diagram chasing and by using the minimality of \mathcal{F}'' . Suppose $z' \in Z'$ and $z_1, z_2 \in Z$ are such that $z' = \mu(z_1) = \mu(z_2)$. By commutativity of the solid arrows in the diagram above, we have

$$\alpha'' \circ \rho(z_1) = \rho \circ \alpha' \circ \mu(z_1) = \rho \circ \alpha' \circ \mu(z_2) = \alpha'' \circ \rho(z_2).$$

Similarly

$$\beta'' \circ \rho(z_1) = \beta'' \circ \rho(z_2).$$

Thus by minimality of \mathcal{F}'' it follows that $\rho(z_1) = \rho(z_2)$. Hence, ρ' can be constructed by setting $\rho'(z') = \rho(z_1)$. This finishes the proof of Lemma 7.1. □

Proof of Proposition 2.1 First we address claim (1) of the Proposition. Let $\mathbf{G} = \{O_i; m_{ij}\}$ be a poset category, $\mathcal{X}, \mathcal{Y}, \mathcal{Z} \in \mathbf{Prob}(\mathbf{G})$ be three \mathbf{G} -diagrams and $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$ be a two-fan. Recall that it can also be considered as a \mathbf{G} -diagram of two-fans

$$\mathcal{F} = \{\mathcal{F}_i; f_{ij}\}.$$

Any minimizing reduction

$$\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}) \longrightarrow \mathcal{F}' = (\mathcal{X} \leftarrow \mathcal{Z}' \rightarrow \mathcal{Y})$$

induces reductions

$$\mathcal{F}_i = (X_i \leftarrow Z_i \rightarrow Y_i) \longrightarrow \mathcal{F}'_i = (X_i \leftarrow Z'_i \rightarrow Y_i)$$

for all i in the index set I . It follows that if all \mathcal{F}_i 's are minimal, then so is \mathcal{F} .

Now we prove the implication in the other direction. Suppose \mathcal{F} is minimal. We have to show that all \mathcal{F}_i are minimal as well. Suppose there exist a non-minimal fan among \mathcal{F}_i 's. For an index $i \in I$ let

$$\check{J}(i) := \{j \in I : \text{Hom}_{\mathbf{G}}(O_j, O_i) \neq \emptyset\}$$

$$\hat{J}(i) := \{j \in I : \text{Hom}_{\mathbf{G}}(O_i, O_j) \neq \emptyset\}.$$

Choose an index i_0 such that

1. \mathcal{F}_{i_0} is not minimal
2. for any $j \in \hat{J}(i_0) \setminus \{i_0\}$ the two-fan \mathcal{F}_j is minimal.

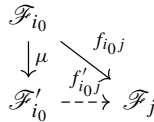
Consider now the minimal reduction $\mu : \mathcal{F}_{i_0} \rightarrow \mathcal{F}'_{i_0}$ and construct a two-fan $\mathcal{G} = \{\mathcal{G}_i; g_{ij}\}$ of \mathbf{G} -diagrams by setting

$$\mathcal{G}_i := \begin{cases} \mathcal{F}'_i & \text{if } i = i_0 \\ \mathcal{F}_i & \text{otherwise} \end{cases}$$

and

$$g_{ij} := \begin{cases} \mu \circ f_{ij} & \text{if } j = i_0 \text{ and } i \in \check{J}(i_0) \\ f'_{ij} & \text{if } i = i_0 \text{ and } j \in \hat{J}(i_0) \\ f_{ij} & \text{otherwise} \end{cases}$$

where f'_{i_0j} is the reduction provided by the Lemma 7.1 applied to the diagram



We thus constructed a non-trivial reduction $\mathcal{F} \rightarrow \mathcal{G}$ which is identity on the sink \mathbf{G} -diagrams \mathcal{X} and \mathcal{Y} . This contradicts the minimality of \mathcal{F} .

To address the second assertion of the Lemma 2.1 observe that the argument above gives an algorithm for the construction of a minimal reduction of any two-fan of \mathbf{G} -diagrams. □

Proposition 3.1 *Let \mathbf{G} be a complete poset category. Then the bivariate function*

$$\mathbf{k} : \mathbf{Prob}(\mathbf{G}) \times \mathbf{Prob}(\mathbf{G}) \rightarrow \mathbb{R}$$

is a pseudo-distance on $\mathbf{Prob}(\mathbf{G})$.

Moreover, two diagrams $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}(\mathbf{G})$ satisfy $\mathbf{k}(\mathcal{X}, \mathcal{Y}) = 0$ if and only if \mathcal{X} is isomorphic to \mathcal{Y} in $\mathbf{Prob}(\mathbf{G})$.

Proof The symmetry of \mathbf{k} is immediate. The non-negativity of \mathbf{k} follows from the fact that entropy of the target space of a reduction is not greater than the entropy of the domain, which is a particular instance of the Shannon inequality (5).

We proceed to prove the triangle inequality. We will make use of the following lemma

Lemma 7.2 *For a minimal full diagram of probability spaces*

$$\mathcal{Q} = \left(\begin{array}{ccccc} & & Q & & \\ & \swarrow & \downarrow & \searrow & \\ U & & W & & V \\ & \downarrow & \swarrow & \nwarrow & \downarrow \\ X & & Y & & Z \end{array} \right)$$

holds

$$kd(X \leftarrow W \rightarrow Z) \leq kd(X \leftarrow U \rightarrow Y) + kd(Y \leftarrow V \rightarrow Z).$$

The Lemma 7.2 follows immediately from Shannon inequality.

Suppose for now that $\mathbf{G} = \bullet$ and we are given three probability spaces X, Y, Z together with the optimal couplings $\mathcal{U} = (X \leftarrow U \rightarrow Y)$ and $\mathcal{V} = (Y \leftarrow V \rightarrow Z)$ in the sense of optimization problem (9). Together they form a two-tents diagram $\mathcal{T} = (X \leftarrow U \rightarrow Y \leftarrow V \rightarrow Z)$. If we can extend \mathcal{T} to a minimal full diagram \mathcal{Q} as in the assumption of Lemma 7.2, the triangle inequality would follow. The diagram $\mathcal{Q} = \mathbf{ad}(\mathcal{T})$ can be constructed by the so called adhesion, as explained below.

As explained in Sect. 2.5.7, to construct a minimal full diagram with sink vertices X, Y and Z it is sufficient to provide a distribution on $\underline{Q} := \underline{X} \times \underline{Y} \times \underline{Z}$ with the correct push-forwards. We do this by setting

$$p_Q(x, y, z) := \frac{p_U(x, y) \cdot p_V(y, z)}{p_Y(y)}.$$

It is straightforward to check that the appropriate restriction of the full diagram defined in the above manner is indeed the original two-tents diagram. Essentially, to extend we need to provide a relationship (coupling) between spaces X and Z and we do it by declaring X and Z independent conditioned on Y . This is an instance of operation called *adhesion*, see [16]. Thus we have shown that $\mathbf{k}: \mathbf{Prob} \times \mathbf{Prob} \rightarrow \mathbb{R}$ is a pseudo-distance.

Assume now that \mathbf{G} is an arbitrary complete poset category. Suppose $\mathcal{X} = \{X_i; f_{ij}\}$, $\mathcal{Y} = \{Y_i; g_{ij}\}$ and $\mathcal{Z} = \{Z_i; h_{ij}\}$ are \mathbf{G} -diagrams, with initial spaces being X_0, Y_0 and Z_0 , respectively. Let

$$\hat{\mathcal{U}} = (\mathcal{X} \leftarrow \mathcal{U} \rightarrow \mathcal{Y}) \quad \text{and} \quad \hat{\mathcal{V}} = (\mathcal{Y} \leftarrow \mathcal{V} \rightarrow \mathcal{Z})$$

be two optimal minimal two-fans.

Recall that each two-fan of \mathbf{G} -diagrams is a \mathbf{G} -diagram of two-fans between the individual spaces, that is

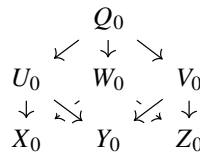
$$\begin{aligned} \mathcal{U} &= \{\mathcal{U}_i = (X_i \leftarrow U_i \rightarrow Y_i)\} \\ \mathcal{V} &= \{\mathcal{V}_i = (Y_i \leftarrow V_i \rightarrow Z_i)\}. \end{aligned}$$

We construct a coupling $\hat{\mathcal{W}}$ between \mathcal{X} and \mathcal{Z} in the following manner. Starting with the two-tents diagram between the initial spaces, we use adhesion to extend it to a full diagram, thus constructing a coupling between X_0 and Z_0 . This full diagram could then be “pushed down” and provides full extensions of two-tents on all lower levels. Thus we could “compose” couplings $\hat{\mathcal{U}}$ and $\hat{\mathcal{V}}$ and use a Shannon inequality to establish the triangle inequality for the intrinsic entropy distance. Details are as follows.

Consider a two-tents diagram

$$X_0 \leftarrow U_0 \rightarrow Y_0 \leftarrow V_0 \rightarrow Z_0$$

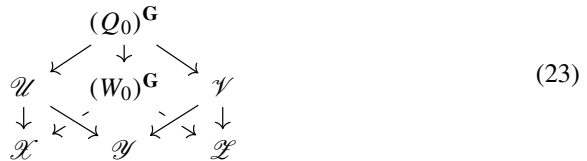
and extend it by adhesion, as described above, to a Λ_3 -diagram



Together with the reductions

$$\begin{array}{lll}
 (X_0)^{\mathbf{G}} \rightarrow \mathcal{X}, & (Y_0)^{\mathbf{G}} \rightarrow \mathcal{Y}, & (Z_0)^{\mathbf{G}} \rightarrow \mathcal{Z} \\
 (U_0)^{\mathbf{G}} \rightarrow \mathcal{U}, & (V_0)^{\mathbf{G}} \rightarrow \mathcal{V} &
 \end{array}$$

explained in Sect. 2.6, it gives rise to a Λ_3 -diagram of \mathbf{G} -diagrams



The diagram above is not necessarily minimal and we now consider the “minimization” of the Λ_3 -diagram (23), as provided by Lemma 2.2

$$\hat{\mathcal{Q}} = \left(\begin{array}{ccccc}
 & & \mathcal{Q} & & \\
 & \swarrow & \downarrow & \searrow & \\
 \mathcal{U} & & \mathcal{W} & & \mathcal{V} \\
 \downarrow & \swarrow & & \searrow & \downarrow \\
 \mathcal{X} & & \mathcal{Y} & & \mathcal{Z}
 \end{array} \right)$$

Applying Lemma 7.2 to each

$$\hat{\mathcal{Q}}_i = \begin{pmatrix} & & Q_i & & \\ & \swarrow & \downarrow & \searrow & \\ U_i & & W_i & & V_i \\ & \downarrow & \swarrow & \nwarrow & \downarrow \\ X_i & & Y_i & & Z_i \end{pmatrix}$$

which is minimal by Lemma 2.2(1), we obtain the required inequality, concluding the proof of the triangle inequality.

Finally, if $k(\mathcal{X}, \mathcal{Y}) = 0$, then there is a two-fan \mathcal{F} of \mathbf{G} -diagrams between \mathcal{X} and \mathcal{Y} with $\text{kd}(\mathcal{F}) = 0$, from which it follows that \mathcal{X} and \mathcal{Y} are isomorphic. \square

Proposition 3.3 *Let \mathbf{G} be a complete poset category. Then with respect to the Kolmogorov distance on $\mathbf{Prob}(\mathbf{G})$ the tensor product*

$$\otimes : (\mathbf{Prob}(\mathbf{G}), k)^2 \rightarrow (\mathbf{Prob}(\mathbf{G}), k)$$

is 1-Lipschitz in each variable, that is, for every triple $\mathcal{X}, \mathcal{Y}, \mathcal{Y}' \in \mathbf{Prob}(\mathbf{G})$ the following bound holds

$$k(\mathcal{X} \otimes \mathcal{Y}, \mathcal{X} \otimes \mathcal{Y}') \leq k(\mathcal{Y}, \mathcal{Y}').$$

Proof The claim follows easily from the additivity of entropy in equation (6). Suppose that $\mathcal{X} = \{X_i; f_{ij}\}$, $\mathcal{Y} = \{Y_i; g_{ij}\}$ and $\mathcal{Y}' = \{Y'_i; g'_{ij}\}$ are three \mathbf{G} -diagrams and

$$\mathcal{F} = (\mathcal{Y} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}')$$

is an optimal fan, so that

$$k(\mathcal{Y}, \mathcal{Y}') = \sum_i [2\text{Ent}(Z_i) - \text{Ent}(Y_i) - \text{Ent}(Y'_i)].$$

Consider the fan

$$\mathcal{G} = (\mathcal{X} \otimes \mathcal{Y} \leftarrow \mathcal{X} \otimes \mathcal{Z} \rightarrow \mathcal{X} \otimes \mathcal{Y}').$$

Then, by additivity of entropy, in equation (6), we have

$$\begin{aligned} \text{kd}(\mathcal{F}) &= \sum_i [2\text{Ent}(X_i \otimes Z_i) - \text{Ent}(X_i \otimes Y_i) - \text{Ent}(X_i \otimes Y'_i)] \\ &= \sum_i [2\text{Ent}(Z_i) - \text{Ent}(Y_i) - \text{Ent}(Y'_i)] \\ &= \text{kd}(\mathcal{G}) \end{aligned}$$

and, therefore,

$$\mathbf{k}(\mathcal{X} \otimes \mathcal{Y}, \mathcal{X} \otimes \mathcal{Y}') \leq \text{kd}(\mathcal{G}) = \text{kd}(\mathcal{F}) = \mathbf{k}(\mathcal{Y}, \mathcal{Y}').$$

Thus, the tensor product of probability spaces is 1-Lipschitz with respect to each argument. □

Proposition 3.7 *For any triple of diagrams $\mathcal{X}, \mathcal{Y}, \mathcal{U}$ holds*

$$\kappa(\mathcal{X} \otimes \mathcal{U}, \mathcal{Y} \otimes \mathcal{U}) = \kappa(\mathcal{X}, \mathcal{Y}).$$

Proof Define a translation-invariant bivariate function κ' on $\mathbf{Prob}\langle \mathbf{G} \rangle$ by setting for a pair $\mathcal{A}, \mathcal{B} \in \mathbf{Prob}\langle \mathbf{G} \rangle$

$$\kappa'(\mathcal{A}, \mathcal{B}) := \inf_{\mathcal{C} \in \mathbf{Prob}\langle \mathbf{G} \rangle} \kappa(\mathcal{A} \otimes \mathcal{C}, \mathcal{B} \otimes \mathcal{C}).$$

Clearly, function κ' is translation invariant and satisfies $\kappa' \leq \kappa$. Our task now is to show that $\kappa' \geq \kappa$.

Fix $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}\langle \mathbf{G} \rangle$. Choose some $\varepsilon > 0$ and take $\mathcal{Z} \in \mathbf{Prob}\langle \mathbf{G} \rangle$ such that

$$\kappa(\mathcal{X} \otimes \mathcal{Z}, \mathcal{Y} \otimes \mathcal{Z}) \leq \kappa'(\mathcal{X}, \mathcal{Y}) + \varepsilon.$$

Since translations are non-expanding, it holds also for any $\mathcal{U} \in \mathbf{Prob}\langle \mathbf{G} \rangle$

$$\kappa(\mathcal{X} \otimes \mathcal{Z} \otimes \mathcal{U}, \mathcal{Y} \otimes \mathcal{Z} \otimes \mathcal{U}) \leq \kappa'(\mathcal{X}, \mathcal{Y}) + \varepsilon.$$

Then for any $n \in \mathbb{N}$ we can estimate

$$\begin{aligned} \kappa(\mathcal{X}, \mathcal{Y}) &= \frac{1}{n} \kappa(\mathcal{X}^n, \mathcal{Y}^n) \\ &\leq \frac{1}{n} [2\kappa(\{\bullet\}^{\mathbf{G}}, \mathcal{Z}) + \kappa(\mathcal{X}^n \otimes \mathcal{Z}, \mathcal{Y}^n \otimes \mathcal{Z})]. \end{aligned}$$

For $i = 0, \dots, n$ we set $\mathcal{T}_i = \mathcal{X}^{n-i} \otimes \mathcal{Y}^i \otimes \mathcal{Z}$. Then we have $\mathcal{T}_0 = \mathcal{X}^n \otimes \mathcal{Z}$ and $\mathcal{T}_n = \mathcal{Y}^n \otimes \mathcal{Z}$. Also for each $i = 0, \dots, n - 1$ the pair $(\mathcal{T}_i, \mathcal{T}_{i+1})$ is a translation of the pair $(\mathcal{X} \otimes \mathcal{Z}, \mathcal{Y} \otimes \mathcal{Z})$ by $\mathcal{X}^{n-i-1} \otimes \mathcal{Y}^i$, therefore

$$\kappa(\mathcal{T}_i, \mathcal{T}_{i+1}) \leq \kappa(\mathcal{X} \otimes \mathcal{Z}, \mathcal{Y} \otimes \mathcal{Z}) \leq \kappa'(\mathcal{X}, \mathcal{Y}) + \varepsilon.$$

Now we continue the estimate

$$\begin{aligned} \kappa(\mathcal{X}, \mathcal{Y}) &\leq \frac{1}{n} [2\kappa(\{\bullet\}^{\mathbf{G}}, \mathcal{Z}) + \kappa(\mathcal{X}^n \otimes \mathcal{Z}, \mathcal{Y}^n \otimes \mathcal{Z})] \\ &\leq \frac{1}{n} \left[2\kappa(\{\bullet\}^{\mathbf{G}}, \mathcal{Z}) + \sum_{i=0}^{n-1} \kappa(\mathcal{T}_i, \mathcal{T}_{i+1}) \right] \end{aligned}$$

$$\leq \frac{2}{n} \kappa(\{\bullet\}^{\mathbf{G}}, \mathcal{X}) + \kappa'(\mathcal{X}, \mathcal{Y}) + \varepsilon.$$

Since n is arbitrarily large, the first summand in the right-hand-side could be dropped and then by choosing $\varepsilon > 0$ arbitrarily small we obtain the required inequality. \square

Proposition 3.8 *Suppose \mathbf{G} is a complete poset category and $\delta = \mathbf{k}, \kappa$ is either intrinsic entropy distance or asymptotic entropy distance on $\mathbf{Prob}\langle \mathbf{G} \rangle$. Then the entropy function*

$$\text{Ent}_* : (\mathbf{Prob}\langle \mathbf{G} \rangle, \delta) \rightarrow (\mathbb{R}^{\llbracket \mathbf{G} \rrbracket}, |\cdot|_1), \quad \mathcal{X} = \{X_i, f_{ij}\} \mapsto (\text{Ent} X_i)_i \in \mathbb{R}^{\llbracket \mathbf{G} \rrbracket}$$

is 1-Lipschitz.

Proof Let $\mathcal{X}, \mathcal{Y} \in \mathbf{Prob}\langle \mathbf{G} \rangle$ and let

$$\mathcal{G} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$$

be an optimal fan with components

$$\mathcal{G}_i = (X_i \leftarrow Z_i \rightarrow Y_i).$$

For a fixed index i we can estimate the difference of entropies

$$\text{Ent}(X_i) - \text{Ent}(Y_i) = 2(\text{Ent}(X_i) - \text{Ent}(Z_i)) + \text{kd}(\mathcal{G}_i) \leq \text{kd}(\mathcal{G}_i).$$

By symmetry we then have

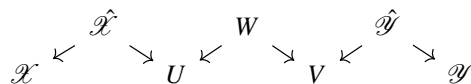
$$|\text{Ent}(X_i) - \text{Ent}(Y_i)| \leq \text{kd}(\mathcal{G}_i).$$

Adding above inequalities for all i we have

$$|\text{Ent}_*(\mathcal{X}) - \text{Ent}_*(\mathcal{Y})|_1 \leq \text{kd}(\mathcal{G}) = \mathbf{k}(\mathcal{X}, \mathcal{Y}).$$

By the additivity of entropy we also obtain the 1-Lipschitz property of the entropy function with respect to the asymptotic entropy distance κ . \square

Proposition 3.9 (Slicing Lemma) *Suppose \mathbf{G} is a complete poset category and we are given $\mathcal{X}, \hat{\mathcal{X}}, \mathcal{Y}, \hat{\mathcal{Y}} \in \hat{\mathbf{Prob}}\langle \mathbf{G} \rangle$ —four \mathbf{G} -diagrams and $U, V, W \in \mathbf{Prob}$ —three probability spaces, that are included into the following three-tents diagram*



such that the two-fan $(U \leftarrow W \rightarrow V)$ is minimal. Then the following estimate holds

$$\mathbf{k}(\mathcal{X}, \mathcal{Y}) \leq \int_W \mathbf{k}(\mathcal{X} \llbracket u, \mathcal{Y} \llbracket v) d p_W(u, v)$$

$$+ \llbracket \mathbf{G} \rrbracket \cdot kd(U \leftarrow W \rightarrow V) + \sum_i [Ent(U \lfloor X_i) + Ent(V \lfloor Y_i)].$$

Proof Since the two-fan $(U \leftarrow W \rightarrow V)$ is minimal the probability space W could be considered having underlying set to be a subset of the Cartesian product of the underlying sets of U and V . For any pair $(u, v) \in \underline{W}$ with a positive weight consider an optimal two-fan

$$\mathcal{G}_{uv} = \left(\mathcal{X} \lfloor u \xleftarrow{\pi_{\mathcal{X}}} \mathcal{Z}_{uv} \xrightarrow{\pi_{\mathcal{Y}}} \mathcal{Y} \lfloor v \right) \tag{24}$$

where $\mathcal{Z}_{uv} = \{Z_{uv,i}; \rho_{ij}\}$. Let $p_{uv,i}$ be the probability distributions on $Z_{uv,i}$ —the individual spaces in the diagram \mathcal{Z}_{uv} . The next step is to take a convex combination of distributions $p_{uv,i}$ weighted by p_W to construct a coupling $\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y}$.

First we extend the 7-vertex diagram to a full Λ_4 -diagram of \mathbf{G} -diagrams, such that the top vertex has the distribution

$$p_i(x, y, u, v) := p_{uv,i}(x, y)p_W(u, v)$$

as described in the Sect. 2.5.7.

If we integrate over y , we obtain

$$\int_{y \in Y} d p_i(x, u, v, y) = ((\pi_{\mathcal{X},i})_* p_{uv,i})(x) p_W(u, v).$$

Equation (24) implies that $(\pi_{\mathcal{X},i})_* p_{uv,i} = p_{X_i}(\cdot \lfloor u)$ and therefore

$$\int_{y \in Y} d p_i(x, y, u, v) = p_{X_i}(x \lfloor u) p_W(u, v).$$

In the same way,

$$\int_{x \in X} d p_i(x, y, u, v) = p_{Y_i}(y \lfloor v) p_W(u, v).$$

It follows that

$$\mathcal{X} \lfloor uv = \mathcal{X} \lfloor u \quad \text{and} \quad \mathcal{Y} \lfloor uv = \mathcal{Y} \lfloor v \tag{25}$$

and

$$Ent(X_i \lfloor UV) = Ent(X_i \lfloor U) \quad \text{and} \quad Ent(Y_i \lfloor UV) = Ent(Y_i \lfloor V). \tag{26}$$

The extended diagram contains a two-fan of diagrams $\mathcal{F} = (\mathcal{X} \leftarrow \mathcal{Z} \rightarrow \mathcal{Y})$ with sink vertices \mathcal{X} and \mathcal{Y} . We call its initial vertex $\mathcal{Z} = \{XY_i, f_{ij}\}$.

The following estimates conclude the proof the Slicing Lemma. First we use the definitions of intrinsic entropy distance \mathbf{k} and of $\text{kd}(\mathcal{F})$ to estimate

$$\begin{aligned} \mathbf{k}(\mathcal{X}, \mathcal{Y}) &\leq \text{kd}(\mathcal{F}) \\ &= \sum_i \text{kd}(\mathcal{F}_i) \\ &= \sum_i [2\text{Ent}(XY_i) - \text{Ent}(X_i) - \text{Ent}(Y_i)]. \end{aligned}$$

Next, we apply the definition of the conditional entropy to rewrite the right-hand side

$$\begin{aligned} \mathbf{k}(\mathcal{X}, \mathcal{Y}) &\leq \sum_i [2\text{Ent}(XY_i|UV) + 2\text{Ent}(UV) - 2\text{Ent}(UV|XY_i) \\ &\quad - \text{Ent}(X_i|U) - \text{Ent}(U) + \text{Ent}(U|X_i) \\ &\quad - \text{Ent}(Y_i|V) - \text{Ent}(V) + \text{Ent}(V|Y_i)]. \end{aligned}$$

We now use (26) and rearrange terms to obtain

$$\begin{aligned} \mathbf{k}(\mathcal{X}, \mathcal{Y}) &\leq \sum_i [2\text{Ent}(XY_i|UV) - \text{Ent}(X_i|UV) - \text{Ent}(Y_i|UV) \\ &\quad + 2\text{Ent}(UV) - \text{Ent}(U) - \text{Ent}(V) \\ &\quad - 2\text{Ent}(UV|XY_i) + \text{Ent}(U|X_i) + \text{Ent}(V|Y_i)]. \end{aligned}$$

By the integral formula for conditional entropy (7) applied to the first three terms we get

$$\begin{aligned} &\sum_i [2\text{Ent}(XY_i|UV) - \text{Ent}(X_i|UV) - \text{Ent}(Y_i|UV)] \\ &= \int_{UV} \mathbf{k}(\mathcal{X}|uv, \mathcal{Y}|uv) d p_W(u, v) \end{aligned}$$

However, because of (25) this simplifies to

$$\int_{UV} \mathbf{k}(\mathcal{X}|uv, \mathcal{Y}|uv) d p_W(u, v) = \int_{UV} \mathbf{k}(\mathcal{X}|u, \mathcal{Y}|v) d p_W(u, v).$$

Therefore,

$$\begin{aligned} \mathbf{k}(\mathcal{X}, \mathcal{Y}) &\leq \int_{UV} \mathbf{k}(\mathcal{X}|u, \mathcal{Y}|v) d p_W(u, v) + \mathbf{[G]} \cdot \text{kd}(U \leftarrow W \rightarrow V) \\ &\quad + \sum_i [\text{Ent}(U|X_i) + \text{Ent}(V|Y_i)]. \end{aligned}$$

□

Acknowledgements Open access funding provided by Max Planck Society. We would like to thank Tobias Fritz, František Matúš, Misha Movshev and Johannes Rauh for inspiring discussions. We are grateful to the participants of the Wednesday Morning Session at the CASA group at the Eindhoven University of Technology for valuable feedback on the introduction of the article. Finally, we thank the Max Planck Institute for Mathematics in the Sciences, Leipzig, for its hospitality.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

1. Ay, N., Bertschinger, N., Der, R., Güttler, F., Olbrich, E.: Predictive information and explorative behavior of autonomous robots. *Eur. Phys. J. B* **63**(3), 329–339 (2008)
2. Baez, J.C., Fritz, T., Leinster, T.: A characterization of entropy in terms of information loss. *Entropy* **13**(11), 1945–1957 (2011)
3. Boltzmann, L.: über die mechanische Bedeutung des zweiten Hauptsatzes der Wärmetheorie. *Wiener Berichte* **53**, 195–220 (1866)
4. Boltzmann, L.: Vorlesungen über Gastheorie, vols. I, II. J.A. Barth, Leipzig (1896)
5. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., Ay, N.: Quantifying unique information. *Entropy* **16**(4), 2161–2183 (2014)
6. Cicalese, F., Gargano, L., Vaccaro, U.: How to find a joint probability distribution of minimum entropy (almost) given the marginals. In: 2017 IEEE International Symposium on Information Theory (ISIT), pp. 2173–2177. IEEE, New York (2017)
7. Csiszár, I.: The method of types. *IEEE Trans. Inf. Theory*, 44(6):2505–2523 (1998) (**Information theory: 1948–1998**)
8. Cover, T.M., Thomas, J.A.: *Elements of Information Theory*. Wiley Series in Telecommunications. A Wiley-Interscience Publication, Wiley, New York (1991)
9. Friston, K.: The free-energy principle: a rough guide to the brain? *Trends Cogn. Sci.* **13**(7), 293–301 (2009)
10. Gromov, M.: In a search for a structure, part 1: on entropy (2012). Preprint. <https://www.ihes.fr/~gromov/wp-content/uploads/2018/08/structre-serch-entropy-july5-2012.pdf> and <https://math.mit.edu/~dspivak/teaching/sp13/gromov--EntropyViaCT.pdf>. Accessed 08 Oct 2018
11. Kocaoglu, M., Dimakis, A.G., Vishwanath, S., Hassibi, B.: Entropic causality and greedy minimum entropy coupling (2017). arXiv preprint. [arXiv:1701.08254](https://arxiv.org/abs/1701.08254)
12. Kolmogorov, A.N.: New metric invariant of transitive dynamical systems and endomorphisms of Lebesgue spaces. *Dokl. Russ. Acad. Sci.* **119**(5), 861–864 (1958)
13. Kolmogorov, A.N.: New metric invariant of transitive dynamical systems and endomorphisms of Lebesgue spaces. *Dokl. Russ. Acad. Sci.* **124**, 754–755 (1959)
14. Kovacevic, M., Stanojevic, I., Senk, V.: On the hardness of entropy minimization and related problems. In: *Information Theory Workshop (ITW)*. 2012 IEEE, pp. 512–516. IEEE, New York (2012)
15. MacLane, S.: *Categories for the Working Mathematician*. Graduate Texts in Mathematics, vol. 5. Springer, New York (1971)
16. Matus, F.: Infinitely many information inequalities. In: *IEEE International Symposium on Information Theory, 2007. ISIT 2007*, pp. 41–44. IEEE, New York (2007)
17. Ornstein, D.: Bernoulli shifts with the same entropy are isomorphic. *Adv. Math.* **4**, 337–352 (1970)
18. Steudel, B., Ay, N.: Information-theoretic inference of common ancestors. *Entropy* **17**(4), 2304–2327 (2015)
19. Shannon, C.E.: A mathematical theory of communication. *Bell Syst. Tech. J.* **27**(3), 379–423 (1948)
20. Shannon, C.: The lattice theory of information. *Trans. IRE Prof. Group Inf. Theory* **1**(1), 105–107 (1953)
21. Sinai, Ya G.: On the notion of entropy of a dynamical system. *Dokl. Russ. Acad. Sci.* **124**, 768–771 (1959)

22. Sinai, Y.G.: Introduction to Ergodic Theory. Princeton University Press, Princeton (1976) (**Translated by V. Scheffer Mathematical Notes, vol. 18**)
23. Van Dijk, S.G., Polani, D.: Informational constraints-driven organization in goal-directed behavior. *Adv. Complex Syst.* **16**(02n03), 1350016 (2013)
24. Vidyasagar, M.: A metric between probability distributions on finite sets of different cardinalities and applications to order reduction. *IEEE Trans. Autom. Control* **57**(10), 2464–2477 (2012)
25. Yeung, R.W.: A First Course in Information Theory. Springer Science & Business Media, Berlin (2012)