

Normalized Information-Based Divergences

J.-F. Coeurjolly, R. Drouilhet, and J.-F. Robineau

Université Pierre Mendès-France, Grenoble, France
Jean-Francois.Coeurjolly@upmf-grenoble.fr
Remy.Drouilhet@upmf-grenoble.fr

Received April 11, 2006; in final form, May 16, 2007

Abstract—This paper is devoted to the mathematical study of some divergences based on mutual information which are well suited to categorical random vectors. These divergences are generalizations of the “entropy distance” and “information distance.” Their main characteristic is that they combine a complexity term and the mutual information. We then introduce the notion of (normalized) information-based divergence, propose several examples, and discuss their mathematical properties, in particular, in some prediction framework.

DOI: 10.1134/S0032946007030015

1. INTRODUCTION

Shannon’s information theory, usually just called *information theory*, was introduced in 1948 [1]. The theory is aimed at providing means for measuring information. More precisely, the amount of information in an object may be measured by its *entropy* and may be interpreted as the length of the description of the object in some encoding way. In the Shannon approach, the objects to be encoded are assumed to be outcomes of a known source. Shannon’s theory also provides the notion of *mutual information* (related to two objects), which plays a central role in many applications, from lossy compression to machine learning methods.

Several authors noted that it would be useful to modify the mutual information such that the resulting quantity becomes a metric in a strict sense. As a first example, [2, 3] introduced the *entropy distance* defined as the sum of conditional entropies. Other interesting measures are the *information distance* [4] and its normalized version, the *similarity metric*, introduced in [5] in the context of the Kolmogorov complexity theory. More precisely, the information distance is defined as the maximum of the conditional Kolmogorov complexities. The similarity metric is universal in the sense defined by the authors and is not computable since it is based on an uncomputable notion of the Kolmogorov complexity.

Recent papers have demonstrated that applications of suitable versions of the similarity metric are of use in areas as diverse as genomics, virology, languages, literature, music, handwritten digits, and astronomy [6]. To apply the metric to real data, the authors have to replace the use of the non-computable Kolmogorov complexity with an approximation obtained by using standard real-world compressors: GenCompress for genomics [7], the *Normalized Compression Distance* (NCD) for music clustering [8], and the *Normalized Google Distance* (NGD) for automatic meaning discovery [9] are examples of effective compressors. To include the information distance and similarity metric in a framework based on information theory concepts, we make use of the principle that *the expected Kolmogorov complexity equals the Shannon entropy*; an interested reader is referred to [10–12] for more details.

Consequently, both the entropy distance and information distance are expressed in terms of conditional entropies: the first one as their sum and the second as their maximum. In [13] there is given a proof of the triangle inequality for these distances and their respective normalized versions.

In the supervised learning framework, the use of some method of selecting covariables among a large number is required when it is assumed that the data size is too small with respect to the number of available covariables (in order to apply any existing discriminant analysis method). Such a problem has been widely treated (see, e.g., [14]). The approach undertaken in [15] is mainly based on three kinds of methodological tools. The first one is a supervised quantization method consisting in the simplification of too complex covariables (in particular, with a too large number of possible values). Indeed, our main belief is that in order to predict the class variable generally representing a small number of categories of data, each possibly predictive covariable must be not too complex. The second one is a more usual step-by-step selection method combining the simplified covariables together in order to detect a cluster of data of the same class. The last one is aimed at detecting redundancy among the set of covariables. These three tasks can be realized using the entropy or information distances (or their normalized versions). Let us emphasize some properties allowing one to understand the usefulness of these criteria in such a context. The entropy and information distances, D^E and D^I , can be rewritten as the difference between some term (respectively, the joint entropy and the maximum of the marginal entropies) and the mutual information. The first term may be interpreted as a complexity term. Moreover, both are independence measures with the particular property to be minimal (in fact equal to 0) when random vectors share exactly the same information. In [15] it was then proposed to extend the definition of the entropy and information distances by introducing the notion of information-based divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ between two categorical random vectors \mathbf{X} and \mathbf{Y} , defined as the difference of some complexity term $C_{\mathbf{X},\mathbf{Y}}$ and the mutual information $I_{\mathbf{X},\mathbf{Y}}$ and such that $C_{\mathbf{X},\mathbf{Y}}$ is an upper bound for $I_{\mathbf{X},\mathbf{Y}}$ reached when \mathbf{X} and \mathbf{Y} share exactly the same information. The notion of the normalized information-based divergence $\delta_{\mathbf{X},\mathbf{Y}}$ is directly derived by dividing the associated information-based divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ by the complexity term $C_{\mathbf{X},\mathbf{Y}}$. The normalized versions d^E and d^I of D^E and D^I are particular examples. Other examples are given in [15]. Among them, one is of particular interest since its complexity term C^S is the mean of marginal entropies. The associated (unnormalized) information-based divergence Δ^S is not so different from D^E since it corresponds to its half. Nevertheless, the expression for its complexity term C^S really differs from the complexity term C^E of D^E (i.e., the joint entropy). For practical purposes, we may argue that D^I , D^E , and Δ^S are not well suited in the prediction framework since a small value of these distances means that both the explained and explicative variables have a good knowledge of each other. This is due to the fact that both conditional entropies have at least the same weight.

In this paper, this drawback is weakened by introducing a natural extension $C^{S,\alpha}$ of the complexity term C^S defined as a weighted mean (by α and $1 - \alpha$ for some $0 < \alpha \leq 1$) of the minimum and maximum of marginal entropies. This kind of complexity term leads to an expected I -divergence $\Delta^{S,\alpha}$, which is the weighted mean of the minimum and maximum of conditional entropies.

The paper is organized as follows. In Section 2 we recall the definitions and the main properties of the entropy and information distances (and their normalized versions). Similarly to [16], we extract the main characteristics to define some general concept of information divergence which could theoretically be applied in a more general setting (continuous, discrete, etc.). In Section 3 we concentrate on categorical data (and in particular discrete) random vectors, since this is usually the case in most applications that use entropy or information distance. We give the definition of the (normalized) information-based divergence and propose several examples. We study their mathematical properties in the general context and propose some sufficient conditions for these divergences to verify some triangle-type inequality. Finally, in Section 4 we exhibit some properties

of information-based divergences in the special prediction framework. In particular, we show that these divergences are useful to detect redundancy.

2. NORMALIZED ENTROPY DISTANCE
AND NORMALIZED INFORMATION DISTANCE

Let us denote by Γ the set of categorical random vectors, that is, discrete-valued random vectors with finite entropy. In the following, \mathbf{X} , \mathbf{Y} , and \mathbf{Z} are three elements of Γ .

2.1. Notation

We denote by $H_{\mathbf{X}}$ (when it exists) the Shannon entropy of \mathbf{X} given by

$$H_{\mathbf{X}} = - \sum_{\mathbf{x} \in \Omega_{\mathbf{X}}} p_{\mathbf{X}}(\mathbf{x}) \log(p_{\mathbf{X}}(\mathbf{x})), \quad \text{with} \quad p_{\mathbf{X}}(\mathbf{x}) = \mathbf{P}(\mathbf{X} = \mathbf{x}).$$

In the same way, one can define the joint entropy of \mathbf{X} and \mathbf{Y} , denoted by $H_{\mathbf{X},\mathbf{Y}}$, and the conditional entropy of \mathbf{X} by \mathbf{Y} (respectively, \mathbf{Y} by \mathbf{X}), denoted by $H_{\mathbf{X}|\mathbf{Y}}$ (respectively, $H_{\mathbf{Y}|\mathbf{X}}$). Finally, we denote by $I_{\mathbf{X},\mathbf{Y}}$ the mutual information between the random vectors \mathbf{X} and \mathbf{Y} . When these different quantities exist, there are the following relations (see, e.g., [17]):

$$H_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{X}} + H_{\mathbf{Y}|\mathbf{X}} = H_{\mathbf{Y}} + H_{\mathbf{X}|\mathbf{Y}}, \tag{1}$$

$$I_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{X}} - H_{\mathbf{X}|\mathbf{Y}} = H_{\mathbf{Y}} - H_{\mathbf{Y}|\mathbf{X}} = H_{\mathbf{X}} + H_{\mathbf{Y}} - H_{\mathbf{X},\mathbf{Y}}. \tag{2}$$

2.2. Definitions and Some Basic Properties

We now present some measures allowing us to overcome some drawbacks of the mutual information. As a first generalization, several authors noted that it would be useful to modify the mutual information such that the resulting quantity becomes a metric in a strict sense. Two such measures exist and are well known in the literature. The first one, called the ‘‘entropy distance,’’ is derived from the domain of information theory. The second one, called ‘‘information distance,’’ originates in works around the Kolmogorov complexity. These measures are defined (when they exist) for two random vectors \mathbf{X} and \mathbf{Y} as follows:

- Entropy distance

$$D_{\mathbf{X},\mathbf{Y}}^E = H_{\mathbf{X}|\mathbf{Y}} + H_{\mathbf{Y}|\mathbf{X}}; \tag{3}$$

- Information distance

$$D_{\mathbf{X},\mathbf{Y}}^I = \max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}}). \tag{4}$$

Both measures are indeed some modifications of the mutual information since from (1) and (2) we have

$$D_{\mathbf{X},\mathbf{Y}}^E = H_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}} \quad \text{and} \quad D_{\mathbf{X},\mathbf{Y}}^I = \max(H_{\mathbf{X}}, H_{\mathbf{Y}}) - I_{\mathbf{X},\mathbf{Y}}. \tag{5}$$

The quantities $H_{\mathbf{X},\mathbf{Y}}$ and $\max(H_{\mathbf{X}}, H_{\mathbf{Y}})$ are upper bounds for the mutual information $I_{\mathbf{X},\mathbf{Y}}$ and are reached when \mathbf{X} and \mathbf{Y} share exactly the same information. In other words, these two measures are nonnegative and vanish if and only if $H_{\mathbf{Y}|\mathbf{X}} = H_{\mathbf{X}|\mathbf{Y}} = 0$, expressing the fact that \mathbf{X} (respectively, \mathbf{Y}) predicts \mathbf{Y} (respectively, \mathbf{X}) with probability 1.

These measures satisfy

$$D_{\mathbf{X},\mathbf{Y}}^E \leq H_{\mathbf{X},\mathbf{Y}} \quad \text{and} \quad D_{\mathbf{X},\mathbf{Y}}^I \leq \max(H_{\mathbf{X}}, H_{\mathbf{Y}}), \tag{6}$$

where the equality holds if the vectors \mathbf{X} and \mathbf{Y} are independent. In [18, 19] it was noted that in bioinformatics an unnormalized distance may not be a proper evolutionary distance measure. To overcome this problem within the algorithmic framework, they form two normalized versions of distances, D^E and D^I . Their Shannon versions were proposed and studied in [13].

Definition 1. When they exist, one defines the following two measures:

- Normalized entropy distance

$$d_{\mathbf{X},\mathbf{Y}}^E = \frac{H_{\mathbf{X}|\mathbf{Y}} + H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{X},\mathbf{Y}}};$$

- Normalized information distance

$$d_{\mathbf{X},\mathbf{Y}}^I = \frac{\max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})}{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})}.$$

Since $H_{\mathbf{X},\mathbf{Y}} = 0 \Leftrightarrow H_{\mathbf{X}} = H_{\mathbf{Y}} = 0 \Leftrightarrow \max(H_{\mathbf{X}}, H_{\mathbf{Y}}) = 0$, we set by convention $d_{\mathbf{X},\mathbf{Y}}^E = 0$ (respectively, $d_{\mathbf{X},\mathbf{Y}}^I = 0$) when $H_{\mathbf{X}} = H_{\mathbf{Y}} = 0$.

We are encouraged to define the following equivalence class: the vectors \mathbf{X} and \mathbf{Y} are said to be equivalent if \mathbf{X} (respectively, \mathbf{Y}) predicts \mathbf{Y} (respectively, \mathbf{X}) with probability 1; we denote

$$\mathbf{X} \sim \mathbf{Y} \iff H_{\mathbf{Y}|\mathbf{X}} = H_{\mathbf{X}|\mathbf{Y}} = 0 \iff I_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{X}} = H_{\mathbf{Y}}. \quad (7)$$

Due to the previous convention,

$$d_{\mathbf{X},\mathbf{Y}}^E = 0 \iff d_{\mathbf{X},\mathbf{Y}}^I = 0 \iff \mathbf{X} \sim \mathbf{Y}.$$

From (1) and (2), one can obtain the following expressions for these two measures, allowing some new interpretations.

Proposition 1. *We have the following expressions for $d_{\mathbf{X},\mathbf{Y}}^E$ and $d_{\mathbf{X},\mathbf{Y}}^I$:*

$$d_{\mathbf{X},\mathbf{Y}}^E = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{H_{\mathbf{X},\mathbf{Y}}}, \quad (8)$$

$$d_{\mathbf{X},\mathbf{Y}}^I = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})} \quad (9)$$

$$= \max\left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}, \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}}\right). \quad (10)$$

Proposition 2. *The measures d^E and d^I constitute two distances bounded by 1.*

To our knowledge, these results were proved in [13]. The proofs are very similar to those of [20], where algorithmic version of these distances were considered. The proof is therefore omitted, but in Section 3.3 we propose a result extending this one in the sense that we give conditions for measures that can be written as (8) and (9) to constitute a metric.

2.3. Concept of Information Divergence

From the previous study related to D^I , D^E , d^I , and d^E , we can exhibit some characteristics useful for an attempt to define the concept of information divergence Δ in a more general setting. Let us first consider a similarity measure $\mathcal{I}_{\mathbf{X},\mathbf{Y}}$ (not necessarily the mutual information) which is minimal (in fact equal to 0) when \mathbf{X} and \mathbf{Y} are independent and maximal (in fact equal to $\mathcal{I}_{\mathbf{X},\mathbf{X}} = \mathcal{I}_{\mathbf{Y},\mathbf{Y}}$) when the distributions of \mathbf{X} given $\mathbf{Y} = \mathbf{y}$ and \mathbf{Y} given $\mathbf{X} = \mathbf{x}$ are trivial. An information divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ could satisfy the following properties:

- [P1] symmetry: $\Delta_{\mathbf{X},\mathbf{Y}} = \Delta_{\mathbf{Y},\mathbf{X}}$;
- [P2] nonnegativeness: $\Delta_{\mathbf{X},\mathbf{Y}} \geq 0$;
- [P3] $\Delta_{\mathbf{X},\mathbf{Y}}$ is minimal (i.e., $\Delta_{\mathbf{X},\mathbf{Y}} = 0$) if and only if \mathbf{X} and \mathbf{Y} share exactly the same information (i.e., $\mathcal{I}_{\mathbf{X},\mathbf{Y}}$ is maximal);
- [P4] $\Delta_{\mathbf{X},\mathbf{Y}}$ is maximal if and only if \mathbf{X} and \mathbf{Y} are independent (i.e., $\mathcal{I}_{\mathbf{X},\mathbf{Y}} = 0$).

Other supplementary properties could be that $\Delta_{\mathbf{X},\mathbf{Y}}$:

- [P5] is normalized: $\Delta_{\mathbf{X},\mathbf{Y}} \in [0, 1]$ and $\Delta_{\mathbf{X},\mathbf{Y}} = 1$ when \mathbf{X} and \mathbf{Y} are independent;
- [P6] satisfies the triangle inequality: $\Delta_{\mathbf{X},\mathbf{Y}} \leq \Delta_{\mathbf{X},\mathbf{Z}} + \Delta_{\mathbf{Z},\mathbf{Y}}$;
- [P7] is invariant under continuous and strictly increasing transformations $\varphi(\cdot)$ and $\psi(\cdot)$ of the vectors \mathbf{X} and \mathbf{Y} whenever they are quantitative random vectors.

There exists a vast literature discussing criteria that satisfy the stated properties. We may cite [21] or a recent work [16], where it is proposed to detect the dependence between two possibly nonlinear processes through the Bhattacharya–Matusita–Hellinger measure of dependence, given by

$$S_\rho = \frac{1}{2} \iint \left(\sqrt{f_1(\mathbf{x}, \mathbf{y})} - \sqrt{f_2(\mathbf{x}, \mathbf{y})} \right)^2 d\mathbf{x} d\mathbf{y},$$

where f_1 (respectively, f_2) is the joint density (respectively, product of marginal densities) of \mathbf{X} and \mathbf{Y} . This measure, which has another advantage to be applicable to both continuous or discrete variables, satisfies properties [P1]–[P7] (in fact, let us precise that [P7] is only valid if $\varphi(\cdot) = \psi(\cdot)$).

In some framework where the purpose is to predict some reference variable, one may find interesting to work with a divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ which combines the minimization of a nonnegative complexity term denoted by $\mathcal{C}_{\mathbf{X},\mathbf{Y}}$ and the maximization of a nonnegative information term $\mathcal{I}_{\mathbf{X},\mathbf{Y}}$. The quantity $\mathcal{C}_{\mathbf{X},\mathbf{Y}}$ is called a complexity term since it is assumed to be expressed as a function of $\mathcal{H}_{\mathbf{X}}$, $\mathcal{H}_{\mathbf{Y}}$, and $\mathcal{H}_{\mathbf{X},\mathbf{Y}}$ measuring in some way the complexity of the vectors \mathbf{X} , \mathbf{Y} , and (\mathbf{X}, \mathbf{Y}) , respectively. In other words, we may expect that an information divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ could also satisfy the following properties:

- [P8] When \mathbf{X}_1 and \mathbf{X}_2 have the same complexity (in the sense that $\mathcal{C}_{\mathbf{Y},\mathbf{X}_1} = \mathcal{C}_{\mathbf{Y},\mathbf{X}_2}$), we have $\Delta_{\mathbf{Y},\mathbf{X}_1} < \Delta_{\mathbf{Y},\mathbf{X}_2}$ whenever \mathbf{X}_1 has a better knowledge about \mathbf{Y} than \mathbf{X}_2 (i.e., $\mathcal{I}_{\mathbf{Y},\mathbf{X}_1} > \mathcal{I}_{\mathbf{Y},\mathbf{X}_2}$);
- [P9] When \mathbf{X}_1 and \mathbf{X}_2 have the same knowledge about \mathbf{Y} (i.e., $\mathcal{I}_{\mathbf{Y},\mathbf{X}_1} = \mathcal{I}_{\mathbf{Y},\mathbf{X}_2}$), we have $\Delta_{\mathbf{Y},\mathbf{X}_1} < \Delta_{\mathbf{Y},\mathbf{X}_2}$ whenever \mathbf{X}_1 is simpler than \mathbf{X}_2 in the sense that $\mathcal{C}_{\mathbf{Y},\mathbf{X}_1} < \mathcal{C}_{\mathbf{Y},\mathbf{X}_2}$. Moreover, in this particular situation
- [P10] $\mathcal{C}_{\mathbf{Y},\mathbf{X}_1} \leq \mathcal{C}_{\mathbf{Y},\mathbf{X}_2}$ must be equivalent to $\mathcal{H}_{\mathbf{X}_1} \leq \mathcal{H}_{\mathbf{X}_2}$;
- [P11] When \mathbf{X}_1 and \mathbf{X}_2 share almost exactly the same information (i.e., $\mathcal{I}_{\mathbf{X}_1,\mathbf{X}_2}$ is almost maximal and $\Delta_{\mathbf{X}_1,\mathbf{X}_2} \simeq 0$), the difference between the divergences $\Delta_{\mathbf{Y},\mathbf{X}_1}$ and $\Delta_{\mathbf{Y},\mathbf{X}_2}$ is almost zero (i.e., $\Delta_{\mathbf{Y},\mathbf{X}_1} \simeq \Delta_{\mathbf{Y},\mathbf{X}_2}$).

A class of candidates that satisfy [P8] and [P9] could be of the form

$$\Delta_{\mathbf{X},\mathbf{Y}} = \frac{\mathcal{C}_{\mathbf{X},\mathbf{Y}} - \mathcal{I}_{\mathbf{X},\mathbf{Y}}}{\mathcal{W}_{\mathbf{X},\mathbf{Y}}}, \tag{11}$$

where $\mathcal{W}_{\mathbf{X},\mathbf{Y}}$ is a positive term. When $\mathcal{W}_{\mathbf{X},\mathbf{Y}} = \mathcal{C}_{\mathbf{X},\mathbf{Y}}$, we obtain a normalized information divergence. Properties [P2] and [P3] and relation (11) imply that $\mathcal{C}_{\mathbf{X},\mathbf{Y}}$ is an upper bound for $\mathcal{I}_{\mathbf{X},\mathbf{Y}}$, reached when \mathbf{X} and \mathbf{Y} share exactly the same information.

In the rest of this paper we concentrate on divergences of the form (11) that are in addition well suited to categorical random variables (and in particular discrete random variables). In such a framework, we shall only describe some entropic-based criteria (i.e., $\mathcal{H}_{\mathbf{X}} = H_{\mathbf{X}}$), and so the information term will be set to the mutual information $\mathcal{I}_{\mathbf{X},\mathbf{Y}}$.

3. INFORMATION-BASED DIVERGENCES AND THEIR NORMALIZED VERSIONS

3.1. Definition and Examples

Definition 2. Two criteria Δ and δ are called, respectively, an information-based divergence and a normalized information-based divergence (for short, *I*-divergence and *NI*-divergence) if they

can be written as

$$\Delta_{\mathbf{X},\mathbf{Y}} = C_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}}, \tag{12}$$

$$\delta_{\mathbf{X},\mathbf{Y}} = \frac{C_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}} = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}}, \tag{13}$$

where the term $C_{\mathbf{X},\mathbf{Y}}$ constitutes a complexity term satisfying

- (i) $C_{\mathbf{X},\mathbf{Y}} = C_{\mathbf{Y},\mathbf{X}}$;
- (ii) $I_{\mathbf{X},\mathbf{Y}} \leq C_{\mathbf{X},\mathbf{Y}}$, and this bound is achieved if and only if the random vectors \mathbf{X} and \mathbf{Y} are equivalent, i.e., if and only if $\mathbf{X} \sim \mathbf{Y}$.

We set by convention $\delta_{\mathbf{X},\mathbf{Y}} = 0$ when $C_{\mathbf{X},\mathbf{Y}} = I_{\mathbf{X},\mathbf{Y}} = 0$.

This definition implies automatically that an I -divergence $\Delta_{\mathbf{X},\mathbf{Y}}$ (respectively, NI -divergence $\delta_{\mathbf{X},\mathbf{Y}}$) satisfies properties [P1]–[P4] (respectively, [P1]–[P5]). In the rest of the paper, the term $C_{\mathbf{X},\mathbf{Y}}$ is expressed as

$$C_{\mathbf{X},\mathbf{Y}} = f_C(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}}, I_{\mathbf{X},\mathbf{Y}}), \tag{14}$$

where $f_C(\cdot, \cdot, \cdot)$ is a nonnegative function. Under this expression for $C_{\mathbf{X},\mathbf{Y}}$, property [P7] is ensured since the conditional entropies and mutual information depend only on the joint probability distribution of the categorical random vectors \mathbf{X} and \mathbf{Y} .

From now on, we propose a series of examples where we adopt the following convention: an I -divergence (respectively, NI -divergence) satisfying the triangle inequality is denoted by D (respectively, d) rather than Δ (respectively, δ). Moreover, each example will be particularized by some discriminating additional letter in the same manner as D^E and D^I (respectively, d^E and d^I), which clearly constitute I -divergences (respectively, NI -divergences).

In [15], we investigate two new entropic criteria naturally expressed by

$$\delta_{\mathbf{X},\mathbf{Y}}^D = \frac{1}{2} \left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}} + \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}} \right) \quad \text{and} \quad \delta_{\mathbf{X},\mathbf{Y}}^S = \frac{H_{\mathbf{X}|\mathbf{Y}} + H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{X}} + H_{\mathbf{Y}}},$$

which can be rewritten as NI -divergences:

$$\delta_{\mathbf{X},\mathbf{Y}}^D = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^D} \quad \text{with} \quad C_{\mathbf{X},\mathbf{Y}}^D = \left(\frac{1}{2} \left(\frac{1}{H_{\mathbf{X}}} + \frac{1}{H_{\mathbf{Y}}} \right) \right)^{-1}, \tag{15}$$

$$\delta_{\mathbf{X},\mathbf{Y}}^S = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^S} \quad \text{with} \quad C_{\mathbf{X},\mathbf{Y}}^S = \frac{1}{2}(H_{\mathbf{X}} + H_{\mathbf{Y}}). \tag{16}$$

Their unnormalized versions are expressed as $\Delta_{\mathbf{X},\mathbf{Y}}^D = C_{\mathbf{X},\mathbf{Y}}^D - I_{\mathbf{X},\mathbf{Y}}$ and $D_{\mathbf{X},\mathbf{Y}}^S = C_{\mathbf{X},\mathbf{Y}}^S - I_{\mathbf{X},\mathbf{Y}}$.

In this paper, we are interested in a large family of I -divergences or NI -divergences with complexity terms of the form

$$C_{\mathbf{X},\mathbf{Y}}^\alpha = g^{-1}(\alpha g(m_{\mathbf{X},\mathbf{Y}}) + (1 - \alpha)g(M_{\mathbf{X},\mathbf{Y}})), \tag{17}$$

with $m_{\mathbf{X},\mathbf{Y}} = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and $M_{\mathbf{X},\mathbf{Y}} = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and where $0 \leq \alpha < 1$ and $g(\cdot)$ is any monotone function on \mathbb{R}^+ . When this is not ambiguous, we set $m = m_{\mathbf{X},\mathbf{Y}}$ and $M = M_{\mathbf{X},\mathbf{Y}}$. To be convinced that I -divergences and NI -divergences with complexity terms of the form (17) satisfy Definition 2(ii), let us note that

$$I_{\mathbf{X},\mathbf{Y}} = g^{-1}(\alpha g(I_{\mathbf{X},\mathbf{Y}}) + (1 - \alpha)g(I_{\mathbf{X},\mathbf{Y}})) \leq g^{-1}(\alpha g(m) + (1 - \alpha)g(M)).$$

When $\alpha = 0$, the complexity term C^α corresponds to C^I . When $\alpha = 1$, the complexity term defined as $\min(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and denoted by $C_{\mathbf{X},\mathbf{Y}}^{\min}$ does not satisfy Definition 2(ii) and therefore [P3]. The associated Δ^{\min} (respectively, δ^{\min}) is not an I -divergence (respectively, NI -divergence).

Now we pay particular attention to the complexity terms $C^{D,\alpha}$, $C^{S,\alpha}$, $C^{R,\alpha}$, and $C^{P,\alpha}$ of the form (17) with, respectively, $g^D(\cdot) = 1/\cdot$, $g^S(\cdot) = \cdot$, $g^R(\cdot) = \sqrt{\cdot}$, and $g^P(\cdot) = \log(\cdot)$:

$$C_{\mathbf{X},\mathbf{Y}}^{D,\alpha} = \left(\alpha \frac{1}{\min(H_{\mathbf{X}}, H_{\mathbf{Y}})} + (1 - \alpha) \frac{1}{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})} \right)^{-1}, \tag{18}$$

$$C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} = \alpha \min(H_{\mathbf{X}}, H_{\mathbf{Y}}) + (1 - \alpha) \max(H_{\mathbf{X}}, H_{\mathbf{Y}}), \tag{19}$$

$$C_{\mathbf{X},\mathbf{Y}}^{R,\alpha} = \left(\alpha \sqrt{\min(H_{\mathbf{X}}, H_{\mathbf{Y}})} + (1 - \alpha) \sqrt{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})} \right)^2, \tag{20}$$

$$C_{\mathbf{X},\mathbf{Y}}^{P,\alpha} = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})^\alpha \max(H_{\mathbf{X}}, H_{\mathbf{Y}})^{1-\alpha}. \tag{21}$$

The previous measures Δ^S , δ^S , Δ^D , and δ^D are particular examples of such a family since the value of $\alpha = \frac{1}{2}$ leads to $C_{\mathbf{X},\mathbf{Y}}^{1/2} = g^{-1}\left(\frac{1}{2}g(H_{\mathbf{X}}) + \frac{1}{2}g(H_{\mathbf{Y}})\right)$. When $\alpha = \frac{1}{2}$, $\Delta^{\bullet,\alpha}$ and $\delta^{\bullet,\alpha}$ will simply be denoted by Δ^\bullet and δ^\bullet , where \bullet stands for S, R, P , and D .

Let us first comment on the particular expressions of the divergences $\Delta^{S,\alpha}$ and $\delta^{D,\alpha}$ associated to $C^{D,\alpha}$ and $C^{S,\alpha}$ given by

$$\begin{aligned} \Delta_{\mathbf{X},\mathbf{Y}}^{S,\alpha} &= \alpha \min(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}}) + (1 - \alpha) \max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}}) \\ &= \alpha \Delta_{\mathbf{X},\mathbf{Y}}^{\min} + (1 - \alpha) D_{\mathbf{X},\mathbf{Y}}^I, \\ \delta_{\mathbf{X},\mathbf{Y}}^{D,\alpha} &= \alpha \min\left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}, \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}}\right) + (1 - \alpha) \max\left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}, \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}}\right) \\ &= \alpha \delta_{\mathbf{X},\mathbf{Y}}^{\min} + (1 - \alpha) d_{\mathbf{X},\mathbf{Y}}^I. \end{aligned}$$

Clearly, the previous representation of $\Delta_{\mathbf{X},\mathbf{Y}}^{S,\alpha}$ (respectively, $\delta_{\mathbf{X},\mathbf{Y}}^{D,\alpha}$) as a convex combination of $\Delta_{\mathbf{X},\mathbf{Y}}^{\min}$ and $D_{\mathbf{X},\mathbf{Y}}^I$ (respectively, $\delta_{\mathbf{X},\mathbf{Y}}^{\min}$ and $d_{\mathbf{X},\mathbf{Y}}^I$) introduces a degree of freedom that could be useful for practical purposes in the prediction framework, where \mathbf{Y} could represent some class variable. According to the parameter α , one may choose between one or two prediction terms, $H_{\mathbf{X}|\mathbf{Y}}$ and $H_{\mathbf{Y}|\mathbf{X}}$ (respectively, $\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}$ and $\frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}}$). This possibility to introduce a nonuniform mixing of the entropic contributions in the expression of the complexity terms seems to be not feasible by a direct adaptation of $C_{\mathbf{X},\mathbf{Y}}^I$.

Remark 1. By choosing $g(\cdot) = (\cdot)^\gamma$ for some $\gamma > 0$, the complexity term is given by $C_{\mathbf{X},\mathbf{Y}}^{\gamma,\alpha} = \|\left(\alpha^{\frac{1}{\gamma}} m, (1 - \alpha)^{\frac{1}{\gamma}} M\right)\|_\gamma$, where $\|\mathbf{x}\|_\gamma = \left(\sum_{i=1}^2 |x_i|^\gamma\right)^{1/\gamma}$ denotes the norm of some vector \mathbf{x} of length 2. Note that for any $0 \leq \alpha \leq 1$, we have

$$(\alpha^\wedge)^{\frac{1}{\gamma}} \|(H_{\mathbf{X}}, H_{\mathbf{Y}})\|_\gamma \leq C_{\mathbf{X},\mathbf{Y}}^{\gamma,\alpha} \leq (\alpha^\vee)^{\frac{1}{\gamma}} \|(H_{\mathbf{X}}, H_{\mathbf{Y}})\|_\gamma,$$

with $\alpha^\wedge = \min(\alpha, 1 - \alpha)$ and $\alpha^\vee = \max(\alpha, 1 - \alpha)$. When γ goes to infinity, $C_{\mathbf{X},\mathbf{Y}}^{\gamma,\alpha}$ converges towards $C_{\mathbf{X},\mathbf{Y}}^I$.

Remark 2. The complexity term C^α is invariant under linear transformations of g . In particular, g and $-g$ provide the same complexity term. Consequently, without loss of generality we could restrict g to be an increasing function.

Let us now propose a result to arrange these different examples considered in this paper. First, some preliminary result is given.

Lemma 1. *Let $C^{(1)}$ and $C^{(2)}$ be two complexity terms of the form (17) with functions g_1 and g_2 . Assume that either the function $g_1 \circ g_2^{-1}$ is concave or $g_2 \circ g_1^{-1}$ is convex. Then $C_{\mathbf{X},\mathbf{Y}}^{(1)} \leq C_{\mathbf{X},\mathbf{Y}}^{(2)}$.*

Proof. By rewriting $g_1 = (g_1 \circ g_2^{-1}) \circ g_2$ when $g_1 \circ g_2^{-1}$ is concave and $g_1^{-1} = g_2^{-1} \circ (g_2 \circ g_1^{-1})$ when $(g_2 \circ g_1^{-1})$ is convex, one gets

$$\begin{aligned} g_1^{-1}(\alpha g_1(m) + (1 - \alpha)g_1(M)) &\leq \begin{cases} g_2^{-1}(\alpha(g_2 \circ g_1^{-1}) \circ g_1(m) + (1 - \alpha)(g_2 \circ g_1^{-1}) \circ g_1(M)) \\ g_1^{-1}(g_1 \circ g_2^{-1}(\alpha g_2(m) + (1 - \alpha)g_2(M))) \end{cases} \\ &\leq g_2^{-1}(\alpha g_2(m) + (1 - \alpha)g_2(M)), \end{aligned}$$

where $m = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and $M = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$. \triangle

Proposition 3. For any I -divergences $\Delta^{(1)}$ or $\Delta^{(2)}$ or any NI -divergences $\delta^{(1)}$ or $\delta^{(2)}$ with respective complexity terms $C^{(1)}$ and $C^{(2)}$, we have the following equivalence:

$$\Delta_{\mathbf{X},\mathbf{Y}}^{(1)} \leq \Delta_{\mathbf{X},\mathbf{Y}}^{(2)} \iff \delta_{\mathbf{X},\mathbf{Y}}^{(1)} \leq \delta_{\mathbf{X},\mathbf{Y}}^{(2)} \iff C_{\mathbf{X},\mathbf{Y}}^{(1)} \leq C_{\mathbf{X},\mathbf{Y}}^{(2)}. \tag{22}$$

Since for any $0 \leq \alpha \leq \alpha' \leq 1$ we have

$$C_{\mathbf{X},\mathbf{Y}}^{\alpha'} \leq C_{\mathbf{X},\mathbf{Y}}^{\alpha} \leq C_{\mathbf{X},\mathbf{Y}}^I, \tag{23}$$

the associated I -divergences and NI -divergences are ordered according to equation (22). Furthermore, a similar result holds for the main examples of this paper since

$$C_{\mathbf{X},\mathbf{Y}}^{D,\alpha} \leq C_{\mathbf{X},\mathbf{Y}}^{P,\alpha} \leq C_{\mathbf{X},\mathbf{Y}}^{R,\alpha} \leq C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \leq C_{\mathbf{X},\mathbf{Y}}^I \leq C_{\mathbf{X},\mathbf{Y}}^E. \tag{24}$$

Proof. Equation (22) is direct. The left-hand side of (23) comes from

$$\begin{aligned} C_{\mathbf{X},\mathbf{Y}}^{\alpha} &= g^{-1}(\alpha g(\min(H_{\mathbf{X}}, H_{\mathbf{Y}})) + (1 - \alpha)g(\max(H_{\mathbf{X}}, H_{\mathbf{Y}}))) \\ &\leq g^{-1}(g(\max(H_{\mathbf{X}}, H_{\mathbf{Y}}))) = C_{\mathbf{X},\mathbf{Y}}^I, \end{aligned}$$

and the right-hand side is direct. Since $g^P \circ (g^D)^{-1}(\cdot) = -\log(\cdot)$, $g^R \circ (g^P)^{-1}(\cdot) = \exp(\frac{1}{2} \cdot)$, and $g^S \circ (g^R)^{-1}(\cdot) = (\cdot)^2$ are convex functions, (24) is a direct consequence of Lemma 1. \triangle

Remark 3. By assuming that either $g(\cdot)$ is a concave function or $g^{-1}(\cdot)$ is a convex function, we have the inequality

$$C_{\mathbf{X},\mathbf{Y}}^{\alpha} \leq \alpha m + (1 - \alpha)M = C_{\mathbf{X},\mathbf{Y}}^{S,\alpha},$$

which means that any Δ^{α} (respectively, δ^{α}) (satisfying the previous assumption) is upper bounded by $\Delta^{S,\alpha}$ (respectively, $\delta^{S,\alpha}$).

The following proposition gives a larger class of examples of I -divergences and NI -divergences.

Proposition 4. Let $(\alpha^{(j)})_{j=1,\dots,J}$ be some vector of probability weights for some $J \geq 1$.

(i) Let $\delta^{(1)}, \dots, \delta^{(J)}$ be NI -divergences. Then the measure defined by

$$\delta_{\mathbf{X},\mathbf{Y}} = \sum_{j=1}^J \alpha^{(j)} \delta_{\mathbf{X},\mathbf{Y}}^{(j)} \tag{25}$$

is an NI -divergence with complexity term given by

$$C_{\mathbf{X},\mathbf{Y}} = \left(\sum_{j=1}^J \frac{\alpha^{(j)}}{C_{\mathbf{X},\mathbf{Y}}^{(j)}} \right)^{-1}. \tag{26}$$

(ii) Let $\Delta^{(1)}, \dots, \Delta^{(j)}$ be I -divergences and $\delta^{(1)}, \dots, \delta^{(j)}$ be NI -divergences with complexity terms $C_{\mathbf{X},\mathbf{Y}}^{(1)}, \dots, C_{\mathbf{X},\mathbf{Y}}^{(J)}$. Then the measures defined by

$$\Delta_{\mathbf{X},\mathbf{Y}} = C_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}} \quad \text{and} \quad \delta_{\mathbf{X},\mathbf{Y}} = 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}}, \quad \text{with} \quad C_{\mathbf{X},\mathbf{Y}} = \sum_{j=1}^J \alpha^{(j)} C_{\mathbf{X},\mathbf{Y}}^{(j)}, \tag{27}$$

are also, respectively, an I -divergence and NI -divergence.

The proof is immediate.

3.2. Around Property [P3]

The fact that an I -divergence Δ (respectively, NI -divergence δ) satisfies property [P3] may be expressed as follows: $\Delta_{\mathbf{X},\mathbf{Y}} = 0 \Leftrightarrow D_{\mathbf{X},\mathbf{Y}}^I = 0$ (respectively, $\delta_{\mathbf{X},\mathbf{Y}} = 0 \Leftrightarrow d_{\mathbf{X},\mathbf{Y}}^I = 0$). In fact, [P3] should be extended to a more useful assumption: $\Delta_{\mathbf{X},\mathbf{Y}}$ (or $\delta_{\mathbf{X},\mathbf{Y}}$) is near from the minimum 0 if and only if \mathbf{X} and \mathbf{Y} share almost the same information. This may be translated by the following implications related to an I -divergence Δ (respectively, NI -divergence δ):

- For all $\gamma > 0$ there exists $\varepsilon > 0$ such that for all $(\mathbf{X}, \mathbf{Y}) \in \Upsilon$

$$\Delta_{\mathbf{X},\mathbf{Y}} \leq \varepsilon \implies D_{\mathbf{X},\mathbf{Y}}^I \leq \gamma \quad (\text{respectively, } \delta_{\mathbf{X},\mathbf{Y}} \leq \varepsilon \implies d_{\mathbf{X},\mathbf{Y}}^I \leq \gamma);$$

- For all $\varepsilon > 0$ there exists $\gamma > 0$ such that for all $(\mathbf{X}, \mathbf{Y}) \in \Upsilon$

$$D_{\mathbf{X},\mathbf{Y}}^I \leq \gamma \implies \Delta_{\mathbf{X},\mathbf{Y}} \leq \varepsilon \quad (\text{respectively, } d_{\mathbf{X},\mathbf{Y}}^I \leq \gamma \implies \delta_{\mathbf{X},\mathbf{Y}} \leq \varepsilon).$$

An I -divergence Δ (respectively, NI -divergence δ) inherits the previous property if it satisfies:

[P3'(Υ, k_1, k_2)] There exist some positive constants k_1 and k_2 ($k_1 \leq k_2$) such that for all pairs $(\mathbf{X}, \mathbf{Y}) \in \Upsilon \subset \Gamma^2$ we have

$$k_1 D_{\mathbf{X},\mathbf{Y}}^I \leq \Delta_{\mathbf{X},\mathbf{Y}} \leq k_2 D_{\mathbf{X},\mathbf{Y}}^I \quad (\text{respectively, } k_1 d_{\mathbf{X},\mathbf{Y}}^I \leq \delta_{\mathbf{X},\mathbf{Y}} \leq k_2 d_{\mathbf{X},\mathbf{Y}}^I). \quad (28)$$

Among our examples, we assert that D^E and d^E both satisfy [P3'($\Gamma^2, 1, 2$)], that is,

$$D_{\mathbf{X},\mathbf{Y}}^I \leq D_{\mathbf{X},\mathbf{Y}}^E \leq 2D_{\mathbf{X},\mathbf{Y}}^I \quad (\text{respectively, } d_{\mathbf{X},\mathbf{Y}}^I \leq d_{\mathbf{X},\mathbf{Y}}^E \leq 2d_{\mathbf{X},\mathbf{Y}}^I).$$

Most of complexity terms considered in this paper are of the particular form (17) where the function $g(\cdot)$ is a monotone function on \mathbb{R}^+ . From (23), we can point out that for such complexity terms (expressed in terms of Δ or δ), the constant k_2 is equal to 1. Moreover, we assert that if Δ satisfies [P3'($\Upsilon, k_1, 1$)], then the associated δ also satisfies [P3'($\Upsilon, k_1, 1$)] since

$$k_1 d_{\mathbf{X},\mathbf{Y}}^I = \frac{k_1 D_{\mathbf{X},\mathbf{Y}}^I}{C_{\mathbf{X},\mathbf{Y}}^I} \leq \frac{\Delta_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}} = \delta_{\mathbf{X},\mathbf{Y}}.$$

Therefore, the results presented hereafter in the rest of this section will be expressed for I -divergences only.

Furthermore, we now consider only complexity terms of the form (17) defined through a function $g(\cdot)$ continuously differentiable on some set $\mathcal{D}^g \subset \mathbb{R}^+$. Let us first introduce the following two subsets of \mathcal{D}^g :

$$\mathcal{E}_1^g = \{\Theta \subset \mathcal{D}^g : 0 < \varkappa_{\text{inf},\Theta}^g < \varkappa_{\text{sup},\Theta}^g < +\infty\} \quad \text{and} \quad \mathcal{E}_2^{g,\alpha} = \left\{ \Theta \subset \mathcal{E}_1^g : \frac{\alpha \varkappa_{\text{sup},\Theta}^g}{\varkappa_{\text{inf},\Theta}^g} < 1 \right\},$$

with $\varkappa_{\text{inf},\Theta}^g = \inf_{x \in \Theta} |g'(x)|$ and $\varkappa_{\text{sup},\Theta}^g = \sup_{x \in \Theta} |g'(x)|$. Denote also $\alpha^\wedge = \min(\alpha, 1 - \alpha)$.

In the following, two results ensuring that an I -divergence Δ^α of the form (17) satisfies property [P3'(Υ, k_1, k_2)] are proposed. The difference relies upon the framework: the constants k_1 and k_2 differ whenever the set Υ differs.

Proposition 5. For any $\Theta \in \mathcal{E}_1^g$, the I -divergence Δ^α satisfies [P3'($\Upsilon_\Theta, \alpha^\wedge \frac{\varkappa_{\text{inf},\Theta}^g}{\varkappa_{\text{sup},\Theta}^g}, 1$)] with $\Upsilon_\Theta = \{(\mathbf{X}, \mathbf{Y}) \in \Gamma^2 : H_{\mathbf{X}}, H_{\mathbf{Y}}, I_{\mathbf{X},\mathbf{Y}} \in \Theta\}$.

Proof. Denote $x = \min(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})$, $y = \max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})$, and $z = I_{\mathbf{X},\mathbf{Y}}$. There exist c_1, c_2 , and c_3 such that

$$\begin{aligned} g^{-1}(\alpha g(x+z) + (1-\alpha)g(y+z)) - z &= (\alpha(g(x+z) - g(z)) + (1-\alpha)(g(y+z) - g(z)))(g^{-1})'(c_1) \\ &= \alpha|g'(c_2)||g^{-1})'(c_1)|x + (1-\alpha)|g'(c_3)||g^{-1})'(c_1)|y, \end{aligned}$$

with $c_1 \in [\min(g(z), \alpha g(x+z) + (1-\alpha)g(y+z)), \max(g(z), \alpha g(x+z) + (1-\alpha)g(y+z))]$, $c_2 \in [z, x+z]$, and $c_3 \in [z, y+z]$. Then for all x, y , and z we obtain

$$g^{-1}(\alpha g(x+z) + (1-\alpha)g(y+z)) - z \geq \alpha^\wedge \frac{\varkappa_{\text{inf},\Theta}^g}{\varkappa_{\text{sup},\Theta}^g} \max(x, y),$$

which means that $\alpha^\wedge \frac{\varkappa_{\text{inf},\Theta}^g}{\varkappa_{\text{sup},\Theta}^g} D_{\mathbf{X},\mathbf{Y}}^I \leq \Delta_{\mathbf{X},\mathbf{Y}}^\alpha \cdot \Delta$

Proposition 6. For any $\Theta \in \mathcal{E}_2^g$, the I-divergence Δ^α satisfies $[\mathbf{P3}'(\Gamma_\Theta^2, 1 - \alpha \frac{\varkappa_{\text{sup},\Theta}^g}{\varkappa_{\text{inf},\Theta}^g}, \mathbf{1})]$ with $\Gamma_\Theta = \{\mathbf{Z} \in \Gamma : H_{\mathbf{Z}} \in \Theta\}$.

Proof. We have

$$\begin{aligned} D_{\mathbf{X},\mathbf{Y}}^I - \Delta_{\mathbf{X},\mathbf{Y}}^\alpha &= C_{\mathbf{X},\mathbf{Y}}^I - C_{\mathbf{X},\mathbf{Y}} = \alpha(g^{-1})'(c_1)(g(\max(H_{\mathbf{X}}, H_{\mathbf{Y}})) - g(\min(H_{\mathbf{X}}, H_{\mathbf{Y}}))) \\ &= \alpha|(g^{-1})'(c_1)||g'(c_2)||H_{\mathbf{X}} - H_{\mathbf{Y}}|, \end{aligned}$$

with $c_1 \in [g(\min(H_{\mathbf{X}}, H_{\mathbf{Y}})), g(\max(H_{\mathbf{X}}, H_{\mathbf{Y}}))]$ and $c_2 \in [\min(H_{\mathbf{X}}, H_{\mathbf{Y}}), \max(H_{\mathbf{X}}, H_{\mathbf{Y}})]$. Then we obtain

$$D_{\mathbf{X},\mathbf{Y}}^I - \Delta_{\mathbf{X},\mathbf{Y}}^\alpha \leq \alpha \frac{\varkappa_{\text{sup},\Theta}^g}{\varkappa_{\text{inf},\Theta}^g} D_{\mathbf{X},\mathbf{Y}}^I,$$

which leads to the result. Δ

For the sake of simplicity, we use the notation $\varkappa_{\text{inf},\Theta}^\bullet$ and $\varkappa_{\text{sup},\Theta}^\bullet$ instead of $\varkappa_{\text{inf},\Theta}^{g^\bullet}$ and $\varkappa_{\text{sup},\Theta}^{g^\bullet}$.

The following result is devoted to our different examples. We apply the previous two propositions and present a new result obtained by taking into account the specific form of each example.

Proposition 7. The I-divergence $\Delta^{\bullet,\alpha}$ satisfies properties $[\mathbf{P3}'(\Upsilon_\Theta, k_1^{a,\bullet}, \mathbf{1})]$ (from Proposition 5), $[\mathbf{P3}'(\Gamma_\Theta^2, k_1^{b,\bullet}, \mathbf{1})]$ (from Proposition 6), and $[\mathbf{P3}'(\Gamma_\Theta^2, k_1^{c,\bullet}, \mathbf{1})]$, where \bullet stands for S, R, P, and D, and where

\bullet	Θ	$\varkappa_{\text{inf},\Theta}^\bullet$	$\varkappa_{\text{sup},\Theta}^\bullet$	$k_1^{a,\bullet} = \alpha^\wedge \frac{\varkappa_{\text{inf},\Theta}^\bullet}{\varkappa_{\text{sup},\Theta}^\bullet}$	$k_1^{b,\bullet} = 1 - \alpha \frac{\varkappa_{\text{sup},\Theta}^\bullet}{\varkappa_{\text{inf},\Theta}^\bullet}$	$k_1^{c,\bullet}$
S	\mathbb{R}^+	1	1	α^\wedge	$1 - \alpha$	
R	$[c_1, c_2]$	$\frac{1}{2\sqrt{c_2}}$	$\frac{1}{2\sqrt{c_1}}$	$\frac{\alpha^\wedge}{\sqrt{\rho}}$	$1 - \alpha\sqrt{\rho}$ (if $\rho < \frac{1}{\alpha^2}$)	$(1 - \alpha) \left(1 - \frac{\alpha}{(1 + \frac{1}{\sqrt{\rho}})^2} \right)$
R	\mathbb{R}^+					$(1 - \alpha)^2$
P	$[c_1, c_2]$	$\frac{1}{c_2}$	$\frac{1}{c_1}$	$\frac{\alpha^\wedge}{\rho}$	$1 - \alpha\rho$ (if $\rho < \frac{1}{\alpha}$)	$\frac{\rho^{1-\alpha} - 1}{\rho - 1}$
D	$[c_1, c_2]$	$\frac{1}{c_2^2}$	$\frac{1}{c_1^2}$	$\frac{\alpha^\wedge}{\rho^2}$	$1 - \alpha\rho^2$ (if $\rho < \frac{1}{\sqrt{\alpha}}$)	$\frac{1}{1 + \frac{\alpha}{1 - \alpha}\rho}$

with $0 < c_1 \leq c_2 < +\infty$ and $\rho = \frac{c_2}{c_1}$.

Proof. The computation of $k_1^{a,\bullet}$ and $k_1^{b,\bullet}$ follows from Propositions 5 and 6. Hence, let us only concentrate on $k_1^{c,\bullet}$ for the complexity terms $C^{R,\alpha}$, $C^{P,\alpha}$, and $C^{D,\alpha}$. We denote $m = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and $M = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$.

- Complexity term $C^{R,\alpha}$:

$$\begin{aligned} D_{\mathbf{X},\mathbf{Y}}^I - \Delta_{\mathbf{X},\mathbf{Y}}^{R,\alpha} &= \alpha(1 - \alpha)(\sqrt{M} - \sqrt{m})^2 + \alpha(M - m) \\ &= \alpha(1 - \alpha)\frac{(M - m)^2}{(\sqrt{M} + \sqrt{m})^2} + \alpha(M - m) \\ &\leq \alpha(1 - \alpha)\frac{(D_{\mathbf{X},\mathbf{Y}}^I)^2}{(\sqrt{M} + \sqrt{m})^2} + \alpha D_{\mathbf{X},\mathbf{Y}}^I. \end{aligned}$$

Thus,

$$\Delta_{\mathbf{X},\mathbf{Y}}^{R,\alpha} \geq (1 - \alpha)D_{\mathbf{X},\mathbf{Y}}^I \left(1 - \alpha \frac{D_{\mathbf{X},\mathbf{Y}}^I}{(\sqrt{M} + \sqrt{m})^2} \right).$$

The result is obtained by noting that

$$\frac{D_{\mathbf{X},\mathbf{Y}}^I}{(\sqrt{m} + \sqrt{M})^2} \leq \frac{M}{(\sqrt{m} + \sqrt{M})^2} = \frac{1}{\left(1 + \sqrt{\frac{m}{M}}\right)^2} \leq \frac{1}{\left(1 + \sqrt{\frac{c_1}{c_2}}\right)^2} \leq 1.$$

- Complexity term $C^{P,\alpha}$: by using the Taylor expansion with integral rest, one obtains

$$\begin{aligned} D_{\mathbf{X},\mathbf{Y}}^I - \Delta_{\mathbf{X},\mathbf{Y}}^{P,\alpha} &= M^{1-\alpha}(M^\alpha - m^\alpha) \\ &= M^{1-\alpha}(M - m) \int_0^1 \frac{\alpha}{(m + t(M - m))^{1-\alpha}} dt \\ &\leq (M - m) \int_0^1 \frac{\alpha}{\left(\frac{1}{\rho} + t\left(1 - \frac{1}{\rho}\right)\right)^{1-\alpha}} dt \\ &\leq D_{\mathbf{X},\mathbf{Y}}^I \frac{1}{1 - \frac{1}{\rho}} \left[\left(\frac{1}{\rho} + t\left(1 - \frac{1}{\rho}\right)\right)^\alpha \right]_0^1 = D_{\mathbf{X},\mathbf{Y}}^I \frac{1 - \left(\frac{1}{\rho}\right)^\alpha}{1 - \frac{1}{\rho}}. \end{aligned}$$

Thus,

$$\Delta_{\mathbf{X},\mathbf{Y}}^{P,\alpha} \geq D_{\mathbf{X},\mathbf{Y}}^I \left(1 - \frac{1 - \left(\frac{1}{\rho}\right)^\alpha}{1 - \frac{1}{\rho}} \right) = D_{\mathbf{X},\mathbf{Y}}^I \frac{\rho^{1-\alpha} - 1}{\rho - 1}.$$

- Complexity term $C^{D,\alpha}$:

$$D_{\mathbf{X},\mathbf{Y}}^I - \Delta_{\mathbf{X},\mathbf{Y}}^{D,\alpha} = M - \frac{mM}{\alpha M + (1 - \alpha)m} = \frac{\alpha M}{\alpha M + (1 - \alpha)m}(M - m) \leq \frac{1}{1 + \frac{1 - \alpha c_1}{\alpha c_2}} D_{\mathbf{X},\mathbf{Y}}^I. \quad \Delta$$

3.3. Around the Triangle Inequality

The question arises now whether an I -divergence or NI -divergence satisfies property [P6], i.e., the triangle inequality. The following proposition establishes sufficient conditions for such measures to constitute a metric.

Lemma 2. *We have*

$$H_{\mathbf{X},\mathbf{Y}} \leq H_{\mathbf{X},\mathbf{Z}} + H_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}}, \quad (29)$$

$$I_{\mathbf{X},\mathbf{Y}} \geq I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}}. \quad (30)$$

Proof. From the general properties on entropy, one can obtain

$$H_{\mathbf{X},\mathbf{Y}} \leq H_{\mathbf{X},\mathbf{Y},\mathbf{Z}} = H_{\mathbf{X},\mathbf{Z}} + H_{\mathbf{Y}|\mathbf{X},\mathbf{Z}} \leq H_{\mathbf{X},\mathbf{Z}} + H_{\mathbf{Y}|\mathbf{Z}} = H_{\mathbf{X},\mathbf{Z}} + H_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}}. \quad (31)$$

Equation (30) directly follows from (2). Δ

Proposition 8. *Assume that the complexity term defining an I-divergence satisfies the following property:*

$$C_{\mathbf{X},\mathbf{Y}} \leq C_{\mathbf{X},\mathbf{Z}} + C_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}}. \quad (32)$$

Then the associated I-divergence satisfies the triangle inequality, that is,

$$\Delta_{\mathbf{X},\mathbf{Y}} \leq \Delta_{\mathbf{X},\mathbf{Z}} + \Delta_{\mathbf{Y},\mathbf{Z}}. \quad (33)$$

In addition, if C satisfies

$$C_{\mathbf{X},\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Z}}), \quad (34)$$

then the associated NI-divergence also satisfies the triangle inequality, that is,

$$\delta_{\mathbf{X},\mathbf{Y}} \leq \delta_{\mathbf{X},\mathbf{Z}} + \delta_{\mathbf{Y},\mathbf{Z}}. \quad (35)$$

Proof. Since the quantity

$$A = -(C_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}}) + (C_{\mathbf{X},\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}}) + (C_{\mathbf{Y},\mathbf{Z}} - I_{\mathbf{Y},\mathbf{Z}}),$$

is nonnegative from (30) and (32), equation (33) immediately follows. Moreover, we have the following equation:

$$\delta_{\mathbf{X},\mathbf{Y}} \leq 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}} + A}. \quad (36)$$

Now it is also easy to see from (34) that

$$A + C_{\mathbf{X},\mathbf{Y}} \geq C_{\mathbf{X},\mathbf{Z}} + C_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}} \geq \max(C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}}).$$

From (36) it follows that

$$\delta_{\mathbf{X},\mathbf{Y}} \leq \frac{C_{\mathbf{X},\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}} + C_{\mathbf{Y},\mathbf{Z}} - I_{\mathbf{Y},\mathbf{Z}}}{\max(C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}})} \leq \frac{C_{\mathbf{X},\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}}}{C_{\mathbf{X},\mathbf{Z}}} + \frac{C_{\mathbf{Y},\mathbf{Z}} - I_{\mathbf{Y},\mathbf{Z}}}{C_{\mathbf{Y},\mathbf{Z}}} = \delta_{\mathbf{X},\mathbf{Z}} + \delta_{\mathbf{Y},\mathbf{Z}}. \quad \Delta$$

Remark 4. In Proposition 8 there is no implication between (32) and (34). Indeed, one may check that the NI-divergence δ^S (with $\alpha = 1/2$ for example) satisfies the first inequality but not the second one. Now consider an NI-divergence with complexity term $C_{\mathbf{X},\mathbf{Y}} = \max(H_{\mathbf{X}}, H_{\mathbf{Y}}) + H_{\mathbf{X}|\mathbf{Y}}H_{\mathbf{Y}|\mathbf{X}}$. By choosing \mathbf{X} , \mathbf{Y} , and \mathbf{Z} such that $H_{\mathbf{X}|\mathbf{Z}} = H_{\mathbf{Y}|\mathbf{Z}} = 0$ and $H_{\mathbf{X}|\mathbf{Y}} = H_{\mathbf{Y}|\mathbf{X}} = I_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{Z}}/3 = H_{\mathbf{X},\mathbf{Y}}/3 > 2$, one sees that (34) is satisfied but (32) is not.

Remark 5. Let us consider an NI-divergence δ with the complexity term

$$C_{\mathbf{X},\mathbf{Y}} = C'_{\mathbf{X},\mathbf{Y}} + \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$$

such that $C'_{\mathbf{X},\mathbf{Y}} \geq 0$ (necessarily $C'_{\mathbf{X},\mathbf{Y}} = 0$ whenever $\mathbf{X} \sim \mathbf{Y}$). Then, Δ and δ satisfy the triangle inequality if C' also satisfies the triangle inequality. However, this is not a necessary condition. Indeed, the triangle inequality is not satisfied for the same example of the previous remark with $C'_{\mathbf{X},\mathbf{Y}} = H_{\mathbf{X}|\mathbf{Y}}H_{\mathbf{Y}|\mathbf{X}}$, for which $C'_{\mathbf{X},\mathbf{Z}} = C'_{\mathbf{Y},\mathbf{Z}} = 0$, whereas $C'_{\mathbf{X},\mathbf{Y}} > 0$.

Let us now propose some examples and consequences through the following corollary.

Corollary 1. (i) *The measures D^E and D^I satisfy condition (32) and so are metrics.*

(ii) *The measures d^E and d^I satisfy conditions (32) and (34) and so are metrics.*

(iii) *The measure $D^{S,\alpha}$ for $\alpha \leq \frac{1}{2}$ satisfies condition (32) and so is a metric. Moreover, when $\alpha > \frac{1}{2}$, this measure does not satisfy (32).*

(iv) *Let $(\alpha^{(j)})_{j=1,\dots,J}$ be some vector of probability weights for some $J \geq 1$. Let $\Delta^{(1)}, \dots, \Delta^{(J)}$ be I -divergences (respectively, $\delta^{(1)}, \dots, \delta^{(J)}$ be NI -divergences) with complexity terms $C_{\mathbf{X},\mathbf{Y}}^{(1)}, \dots, C_{\mathbf{X},\mathbf{Y}}^{(J)}$ satisfying (32) (respectively, (32) and (34)). Then the measures defined by (27) satisfy the triangle inequality.*

Proof. (i), (ii) Equation (29) corresponds exactly to (32) for $C_{\mathbf{X},\mathbf{Y}}^E = H_{\mathbf{X},\mathbf{Y}}$. Since $H_{\mathbf{X},\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Z}})$, we have proved that D^E and d^E are metrics. Concerning D^I and d^I , the complexity term corresponds to $C_{\mathbf{X},\mathbf{Y}}^I = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$. Thus, it suffices to prove (32), which is quite obvious. Indeed,

$$\max(H_{\mathbf{X}}, H_{\mathbf{Z}}) + \max(H_{\mathbf{Y}}, H_{\mathbf{Z}}) - H_{\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Y}}).$$

(iii) Let $m = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and $M = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$. We distinguish three cases.

- $H_{\mathbf{Z}} < m$:

$$C_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + C_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} - H_{\mathbf{Z}} = (2\alpha - 1)H_{\mathbf{Z}} + (1 - \alpha)(m + M).$$

If $\alpha > \frac{1}{2}$ and $H_{\mathbf{X}} = H_{\mathbf{Y}}$, the right-hand side of the previous equation equals $(1 - 2\alpha) \times (C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} - H_{\mathbf{Z}}) + C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} < C_{\mathbf{X},\mathbf{Y}}^{S,\alpha}$. Thus, (32) can never be satisfied for $\alpha > \frac{1}{2}$. Now, if $\alpha \leq \frac{1}{2}$, we have

$$C_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + C_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} - H_{\mathbf{Z}} > (1 - \alpha)(m + M) \geq C_{\mathbf{X},\mathbf{Y}}^{S,\alpha}.$$

- $H_{\mathbf{Z}} > M$:

$$C_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + C_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} - H_{\mathbf{Z}} = (2\alpha - 1)H_{\mathbf{Z}} + (1 - \alpha)(m + M) \geq \alpha m + (1 - \alpha)M = C_{\mathbf{X},\mathbf{Y}}^{S,\alpha}.$$

- $m \leq H_{\mathbf{Z}} \leq M$:

$$C_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + C_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} - H_{\mathbf{Z}} = \alpha m + (1 - \alpha)M = C_{\mathbf{X},\mathbf{Y}}^{S,\alpha}.$$

(iv) This is trivial. Δ

We claim that the measures $\Delta^{R,\alpha}$, $\Delta^{P,\alpha}$, and $\Delta^{D,\alpha}$ (and so $\delta^{R,\alpha}$, $\delta^{P,\alpha}$, and $\delta^{D,\alpha}$) do not satisfy condition (32). Consider, for example, $\Delta^{D,\alpha}$. Let us choose \mathbf{X} , \mathbf{Y} , and \mathbf{Z} such that $H_{\mathbf{Z}} > \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and $H_{\mathbf{Z}} = \frac{2 - \alpha}{1 - \alpha} H_{\mathbf{X}} = \frac{2 - \alpha}{1 - \alpha} H_{\mathbf{Y}}$. This leads to

$$C_{\mathbf{X},\mathbf{Z}}^{D,\alpha} + C_{\mathbf{Y},\mathbf{Z}}^{D,\alpha} - H_{\mathbf{Z}} = H_{\mathbf{Z}} \left(\frac{H_{\mathbf{X}}}{\alpha H_{\mathbf{X}} + (1 - \alpha)H_{\mathbf{Z}}} + \frac{H_{\mathbf{Y}}}{\alpha H_{\mathbf{Y}} + (1 - \alpha)H_{\mathbf{Z}}} - 1 \right) = 0 < C_{\mathbf{X},\mathbf{Y}}^{D,\alpha},$$

which contradicts (32).

Concerning these divergences and the measures $\Delta^{S,\alpha}$ (for $\alpha > \frac{1}{2}$) and $\delta^{S,\alpha}$, we do not know if they satisfy the triangle inequality. We can only say that our tool cannot be applied to prove it. We suggest to weaken property [P6] in the following way in order to obtain more results. Let an I -divergence or NI -divergence satisfy

[P6'(Υ , c)] There exists $c \geq 1$ such that for all $(\mathbf{X}, \mathbf{Y}), (\mathbf{Y}, \mathbf{Z}), (\mathbf{X}, \mathbf{Z}) \in \Upsilon$

$$\Delta_{\mathbf{X},\mathbf{Y}} \leq c(\Delta_{\mathbf{X},\mathbf{Z}} + \Delta_{\mathbf{Y},\mathbf{Z}}).$$

Property **[P6]** is then equivalent to **[P6'($\Gamma^2, \mathbf{1}$)]**, and we already know that D^E, d^E, D^I, d^I , and $D^{S,\alpha}$ (for $\alpha \leq \frac{1}{2}$) satisfy **[P6'($\Gamma^2, \mathbf{1}$)]**. When $\Upsilon \subsetneq \Gamma^2$, property **[P6']** is in a sense local, whereas it is global (as the classical triangle inequality) when $\Upsilon = \Gamma^2$.

Note that if an I -divergence (or NI -divergence) satisfies **[P3'($\Upsilon, \mathbf{k}_1, \mathbf{k}_2$)]**, then **[P6'($\Upsilon, \frac{\mathbf{k}_2}{\mathbf{k}_1}$)]** is satisfied since

$$\Delta_{\mathbf{X},\mathbf{Y}} \leq k_2 D_{\mathbf{X},\mathbf{Y}}^I \leq k_2 (D_{\mathbf{X},\mathbf{Z}}^I + D_{\mathbf{Y},\mathbf{Z}}^I) \leq \frac{k_2}{k_1} (\Delta_{\mathbf{X},\mathbf{Z}} + \Delta_{\mathbf{Y},\mathbf{Z}}).$$

We then inherit a lot of results from Proposition 7 related to our examples. In particular, $\Delta^{\bullet,\alpha}$ and $\delta^{\bullet,\alpha}$ (where \bullet stands for S, R, P , and D) both satisfy **[P6'($\Upsilon_\Theta, \frac{\mathbf{1}}{\mathbf{k}_1^a}$)]**, **[P6'($\Gamma_\Theta^2, \frac{\mathbf{1}}{\mathbf{k}_1^b}$)]**, and **[P6'($\Gamma_\Theta^2, \frac{\mathbf{1}}{\mathbf{k}_1^c}$)]**.

In the rest of this section, we attempt to ensure the global property **[P6'(Γ^2, \mathbf{c})]**. From Proposition 7 (with $\Theta = \mathbb{R}^+$) we see that the divergences $\Delta^{S,\alpha}$ (when $\alpha > \frac{1}{2}$) and $\delta^{S,\alpha}$ (respectively, $\Delta^{R,\alpha}$ and $\delta^{R,\alpha}$) satisfy **[P6'($\Gamma^2, \frac{\mathbf{1}}{1-\alpha}$)]** (respectively, **[P6'($\Gamma^2, \frac{\mathbf{1}}{(1-\alpha)^2}$)]**).

When $\alpha \leq \frac{1}{2}$, we could improve the previous result on $\Delta^{R,\alpha}$ by proving that it satisfies property **[P6'($\Gamma^2, \frac{\mathbf{1}}{\alpha^2 + (1-\alpha)^2}$)]**, in the same spirit as in the proof leading to **[P3']**. Indeed,

$$\begin{aligned} D_{\mathbf{X},\mathbf{Y}}^{S,\alpha} - \Delta_{\mathbf{X},\mathbf{Y}}^{R,\alpha} &= \alpha(1-\alpha)(m+M-2\sqrt{mM}) \\ &\leq 2\alpha(1-\alpha)(D_{\mathbf{X},\mathbf{Y}}^{S,\alpha} + I_{\mathbf{X},\mathbf{Y}} - \sqrt{mM}) \\ &\leq 2\alpha(1-\alpha)D_{\mathbf{X},\mathbf{Y}}^{S,\alpha}, \end{aligned}$$

which leads to $\Delta_{\mathbf{X},\mathbf{Y}}^{R,\alpha} \geq (\alpha^2 + (1-\alpha)^2)D_{\mathbf{X},\mathbf{Y}}^{S,\alpha}$. Finally, let us note that

$$\Delta_{\mathbf{X},\mathbf{Y}}^{R,\alpha} \leq D_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \leq D_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + D_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} \leq \frac{1}{\alpha^2 + (1-\alpha)^2} (\Delta_{\mathbf{X},\mathbf{Z}}^{R,\alpha} + \Delta_{\mathbf{Y},\mathbf{Z}}^{R,\alpha}).$$

Now we give a more general result, allowing us, in particular, to improve **[P6'($\Gamma^2, \frac{\mathbf{1}}{1-\alpha}$)]** for $\Delta^{S,\alpha}$ when $\alpha > \frac{1}{2}$.

Proposition 9. *Let us consider the following assumptions on a complexity term: there exists a constant $c \geq 1$ such that*

$$cC_{\mathbf{X},\mathbf{Z}} + cC_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}) \geq C_{\mathbf{X},\mathbf{Y}}, \quad (37)$$

$$cC_{\mathbf{X},\mathbf{Z}} + cC_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}) \geq \max(C_{\mathbf{X},\mathbf{Y}}, C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}}). \quad (38)$$

*If an I -divergence satisfies (37) or an NI -divergence satisfies (38), then they satisfy property **[P6'(Γ^2, \mathbf{c})]**.*

Proof. Let us introduce

$$A = -(C_{\mathbf{X},\mathbf{Y}} - I_{\mathbf{X},\mathbf{Y}}) + c(C_{\mathbf{X},\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}}) + c(C_{\mathbf{Y},\mathbf{Z}} - I_{\mathbf{Y},\mathbf{Z}}).$$

From (30) and (37), one may assert that

$$A \geq cC_{\mathbf{X},\mathbf{Z}} + cC_{\mathbf{Y},\mathbf{Z}} - C_{\mathbf{X},\mathbf{Y}} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}) \geq 0,$$

which implies that the result is valid for Δ . Now from (38) one may write

$$A + C_{\mathbf{X},\mathbf{Y}} \geq \max(C_{\mathbf{X},\mathbf{Y}}, C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}}) \geq \max(C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}}),$$

which leads to

$$\delta_{\mathbf{X},\mathbf{Y}} \leq \frac{c(C_{\mathbf{X},\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}}) + c(C_{\mathbf{Y},\mathbf{Z}} - I_{\mathbf{Y},\mathbf{Z}})}{\max(C_{\mathbf{X},\mathbf{Z}}, C_{\mathbf{Y},\mathbf{Z}})} \leq c\delta_{\mathbf{X},\mathbf{Z}} + c\delta_{\mathbf{Y},\mathbf{Z}}. \quad \Delta$$

Corollary 2. *The measures $\Delta^{S,\alpha}$ for $\alpha > \frac{1}{2}$ satisfy $[\mathbf{P6}'(\Gamma^2, \frac{\alpha}{1-\alpha})]$.*

Proof. Let us concentrate on $\Delta^{S,\alpha}$ for $\alpha > \frac{1}{2}$. Let

$$A = cC_{\mathbf{X},\mathbf{Z}}^{S,\alpha} + cC_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}).$$

Without loss of generality, we assume that $H_{\mathbf{X}} \leq H_{\mathbf{Y}}$. We distinguish three cases.

- $H_{\mathbf{Z}} \leq H_{\mathbf{X}} \leq H_{\mathbf{Y}}$: since $C_{\mathbf{Y},\mathbf{Z}}^{S,\alpha} \geq I_{\mathbf{Y},\mathbf{Z}}$, we have

$$A \geq c(1-\alpha)H_{\mathbf{X}} + (1-\alpha)H_{\mathbf{Y}} + (c\alpha + \alpha - 1)H_{\mathbf{Z}} - (c-1)I_{\mathbf{X},\mathbf{Z}}.$$

Then

$$A - C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \geq (c(1-\alpha) - \alpha)H_{\mathbf{X}} + (c\alpha + \alpha - 1)H_{\mathbf{Z}} - (c-1)I_{\mathbf{X},\mathbf{Z}} \geq (c-1)(H_{\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}}) \geq 0$$

as soon as $c \geq \frac{\alpha}{1-\alpha}$.

- $H_{\mathbf{X}} \leq H_{\mathbf{Y}} \leq H_{\mathbf{Z}}$: we have

$$A \geq \alpha H_{\mathbf{X}} + c\alpha H_{\mathbf{Y}} + ((1-\alpha) + c(1-\alpha) - 1)H_{\mathbf{Z}} - (c-1)I_{\mathbf{Y},\mathbf{Z}}.$$

Then

$$A - C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \geq (c\alpha - (1-\alpha))H_{\mathbf{Y}} + ((1-\alpha) + c(1-\alpha) - 1)H_{\mathbf{Z}} - (c-1)I_{\mathbf{Y},\mathbf{Z}} \geq (c-1)(H_{\mathbf{Y}} - I_{\mathbf{Y},\mathbf{Z}}) \geq 0$$

as soon as $c \geq \frac{\alpha}{1-\alpha}$.

- $H_{\mathbf{X}} < H_{\mathbf{Z}} < H_{\mathbf{Y}}$: we have

$$A \geq c\alpha H_{\mathbf{X}} + (1-\alpha)H_{\mathbf{Y}} + (c(1-\alpha)H_{\mathbf{Z}} + \alpha - 1) - (c-1)I_{\mathbf{X},\mathbf{Z}}.$$

Then

$$A - C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \geq (c-1)\alpha H_{\mathbf{X}} + (c-1)(1-\alpha)H_{\mathbf{Z}} - I_{\mathbf{X},\mathbf{Z}} \geq 0.$$

Hence, we obtain $A - C_{\mathbf{X},\mathbf{Y}}^{S,\alpha} \geq 0$ for $c = \frac{\alpha}{1-\alpha}$. Δ

Remark 6. The tool presented in Proposition 9 cannot be applied to the I -divergence $\Delta^{D,\alpha}$ and NI -divergence $\delta^{D,\alpha}$. Indeed, let us be given some $c \geq 1$ and consider the quantity

$$A = cC_{\mathbf{X},\mathbf{Z}}^{D,\alpha} + cC_{\mathbf{Y},\mathbf{Z}}^{D,\alpha} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}).$$

In fact, one can always find \mathbf{X} , \mathbf{Y} , and \mathbf{Z} such that for all $c \geq 1$ the quantity A is negative. Indeed, let us choose \mathbf{Z} independent of \mathbf{X} and \mathbf{Y} and such that $\alpha H_{\mathbf{Z}} + (1-\alpha)H_{\mathbf{X}} = 3cH_{\mathbf{X}}$ and $\alpha H_{\mathbf{Z}} + (1-\alpha)H_{\mathbf{Y}} = 3cH_{\mathbf{Y}}$. It is easy to see that $H_{\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and therefore $A = H_{\mathbf{Z}}\left(\frac{1}{3} + \frac{1}{3} - 1\right) < 0$. In the same manner, the tool is inapplicable to the I -divergence $\Delta^{P,\alpha}$ and NI -divergence $\delta^{P,\alpha}$. Indeed, let us take \mathbf{Z} independent of \mathbf{X} and \mathbf{Y} and such that $H_{\mathbf{X}} = H_{\mathbf{Y}} = \left(\frac{1}{3c}\right)^{1/\alpha} H_{\mathbf{Z}}$; then $H_{\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$ and

$$A = cC_{\mathbf{X},\mathbf{Z}}^{P,\alpha} + cC_{\mathbf{Y},\mathbf{Z}}^{P,\alpha} - H_{\mathbf{Z}} - (c-1)(I_{\mathbf{X},\mathbf{Z}} + I_{\mathbf{Y},\mathbf{Z}}) = -\frac{1}{3}H_{\mathbf{Z}} < 0.$$

The following result is an extension of Proposition 9, well suited to be applied to $\delta^{D,\alpha}$.

Proposition 10. *Let us assume that there exist two positive integers, I and J such that an NI-divergence $\delta_{\mathbf{X},\mathbf{Y}}$ can be expressed as*

$$\delta_{\mathbf{X},\mathbf{Y}} = \sum_{i=1}^I \frac{S_{\mathbf{X},\mathbf{Y}}^{(i)}}{U_{\mathbf{X},\mathbf{Y}}^{(i)}} = \sum_{j=1}^J \alpha^{(j)} \left(1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^{(j)}} \right),$$

where $(\alpha^{(j)})_{j=1,\dots,J}$ is some vector of probability weights. Denote $S_{\mathbf{X},\mathbf{Y}} = \sum_{i=1}^I S_{\mathbf{X},\mathbf{Y}}^{(i)}$ and $U_{\mathbf{X},\mathbf{Y}} = \max_{i=1,\dots,I} U_{\mathbf{X},\mathbf{Y}}^{(i)}$. Then, if there exists some real number $c \geq 1$ such that for any $j = 1, \dots, J$ we have

- (i) $A^{(j)} = I_{\mathbf{X},\mathbf{Y}} - C_{\mathbf{X},\mathbf{Y}}^{(j)} + c(S_{\mathbf{X},\mathbf{Z}} + S_{\mathbf{Z},\mathbf{Y}}) \geq 0$,
- (ii) $A^{(j)} + C_{\mathbf{X},\mathbf{Y}}^{(j)} \geq \max(U_{\mathbf{X},\mathbf{Z}}, U_{\mathbf{Z},\mathbf{Y}})$,

then δ satisfies [P6'(Γ^2, c)].

Proof. Using assumptions (i) and (ii), one can prove that for all $j = 1, \dots, J$ we have

$$1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^{(j)}} \leq 1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^{(j)} + A^{(j)}} \leq c \frac{S_{\mathbf{X},\mathbf{Z}} + S_{\mathbf{Y},\mathbf{Z}}}{\max(U_{\mathbf{X},\mathbf{Z}}, U_{\mathbf{Z},\mathbf{Y}})} \leq c(\delta_{\mathbf{X},\mathbf{Z}} + \delta_{\mathbf{Y},\mathbf{Z}}).$$

It follows that

$$\delta_{\mathbf{X},\mathbf{Y}} = \sum_{j=1}^J \alpha^{(j)} \left(1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{C_{\mathbf{X},\mathbf{Y}}^{(j)}} \right) \leq \sum_{j=1}^J \alpha^{(j)} c(\delta_{\mathbf{X},\mathbf{Z}} + \delta_{\mathbf{Y},\mathbf{Z}}) = c(\delta_{\mathbf{X},\mathbf{Z}} + \delta_{\mathbf{Y},\mathbf{Z}}). \quad \triangle$$

Corollary 3. *The measure $\delta^{D,\alpha}$ satisfies [P6'($\Gamma^2, \frac{1}{\alpha^\wedge}$)].*

Proof. We have

$$\begin{aligned} \delta_{\mathbf{X},\mathbf{Y}}^{D,\alpha} &= \alpha \min \left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}, \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}} \right) + (1 - \alpha) \max \left(\frac{H_{\mathbf{X}|\mathbf{Y}}}{H_{\mathbf{X}}}, \frac{H_{\mathbf{Y}|\mathbf{X}}}{H_{\mathbf{Y}}} \right) \\ &= \alpha \frac{\min(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})}{\min(H_{\mathbf{X}}, H_{\mathbf{Y}})} + (1 - \alpha) \frac{\max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})}{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})} \\ &= \alpha \left(1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{\min(H_{\mathbf{X}}, H_{\mathbf{Y}})} \right) + (1 - \alpha) \left(1 - \frac{I_{\mathbf{X},\mathbf{Y}}}{\max(H_{\mathbf{X}}, H_{\mathbf{Y}})} \right). \end{aligned}$$

Using the notation introduced in Proposition 10, we have $I = J = 2$, $S_{\mathbf{X},\mathbf{Y}}^{(1)} = \alpha \min(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})$, $S_{\mathbf{X},\mathbf{Y}}^{(2)} = (1 - \alpha) \max(H_{\mathbf{X}|\mathbf{Y}}, H_{\mathbf{Y}|\mathbf{X}})$, $U_{\mathbf{X},\mathbf{Y}}^{(1)} = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$, $U_{\mathbf{X},\mathbf{Y}}^{(2)} = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$, $C_{\mathbf{X},\mathbf{Y}}^{(1)} = \min(H_{\mathbf{X}}, H_{\mathbf{Y}})$, and $C_{\mathbf{X},\mathbf{Y}}^{(2)} = \max(H_{\mathbf{X}}, H_{\mathbf{Y}})$. Let us fix c to the value $\frac{1}{\alpha^\wedge}$. We have

$$\begin{aligned} A^{(1)} &= I_{\mathbf{X},\mathbf{Y}} - \min(H_{\mathbf{X}}, H_{\mathbf{Y}}) + \frac{1}{\alpha^\wedge} \left(\alpha \min(H_{\mathbf{X}|\mathbf{Z}}, H_{\mathbf{Z}|\mathbf{X}}) + (1 - \alpha) \max(H_{\mathbf{X}|\mathbf{Z}}, H_{\mathbf{Z}|\mathbf{X}}) \right. \\ &\quad \left. + \alpha \min(H_{\mathbf{Y}|\mathbf{Z}}, H_{\mathbf{Z}|\mathbf{Y}}) + (1 - \alpha) \max(H_{\mathbf{Y}|\mathbf{Z}}, H_{\mathbf{Z}|\mathbf{Y}}) \right). \end{aligned}$$

Clearly, from (29) we have

$$\begin{aligned} A^{(1)} &\geq \max(H_{\mathbf{X}}, H_{\mathbf{Y}}) - H_{\mathbf{X},\mathbf{Y}} + 2H_{\mathbf{X},\mathbf{Z}} + 2H_{\mathbf{Y},\mathbf{Z}} - H_{\mathbf{X}} - H_{\mathbf{Y}} - 2H_{\mathbf{Z}} \\ &\geq H_{\mathbf{X},\mathbf{Z}} + H_{\mathbf{Y},\mathbf{Z}} - \min(H_{\mathbf{X}}, H_{\mathbf{Y}}) - H_{\mathbf{Z}} \geq 0. \end{aligned}$$

Furthermore,

$$\min(H_{\mathbf{X}}, H_{\mathbf{Y}}) + A^{(1)} \geq H_{\mathbf{X}, \mathbf{Z}} + H_{\mathbf{Y}, \mathbf{Z}} - H_{\mathbf{Z}} \geq \max(H_{\mathbf{X}}, H_{\mathbf{Y}}, H_{\mathbf{Z}}) = \max(U_{\mathbf{X}, \mathbf{Z}}, U_{\mathbf{Y}, \mathbf{Z}}).$$

It follows that $A^{(1)}$ satisfies conditions (i) and (ii) of Proposition 10 with $c = \frac{1}{\alpha^\wedge}$. The proof is strictly similar for $A^{(2)}$. \triangle

4. PREDICTION FRAMEWORK

Let us consider properties related to the prediction of some fixed random vector \mathbf{Y} .

4.1. Prediction Framework

Recall that our purpose is to find a random vector \mathbf{X} that minimizes $\Delta_{\mathbf{Y}, \mathbf{X}}$ (respectively, $\delta_{\mathbf{Y}, \mathbf{X}}$), which combines a complexity term $C_{\mathbf{X}, \mathbf{Y}}$ (to minimize) and an information term $I_{\mathbf{X}, \mathbf{Y}}$ (to maximize). Let us assume that we have already got some \mathbf{X}_1 and its associated measure $\Delta_{\mathbf{Y}, \mathbf{X}_1}$ (respectively, $\delta_{\mathbf{Y}, \mathbf{X}_1}$). After evaluating $\Delta_{\mathbf{Y}, \mathbf{X}_2}$ (respectively, $\delta_{\mathbf{Y}, \mathbf{X}_2}$), we may be interested in describing conditions under which \mathbf{X}_2 is better (or worse) than \mathbf{X}_1 .

Proposition 11. *The following cases may occur.*

Case 1. *We choose \mathbf{X}_2 instead of \mathbf{X}_1 when*

$$\begin{aligned} \Delta_{\mathbf{Y}, \mathbf{X}_2} < \Delta_{\mathbf{Y}, \mathbf{X}_1} &\iff C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} < I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}, \\ \delta_{\mathbf{Y}, \mathbf{X}_2} < \delta_{\mathbf{Y}, \mathbf{X}_1} &\iff \frac{C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1}}{C_{\mathbf{Y}, \mathbf{X}_1}} < \frac{I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}}{I_{\mathbf{Y}, \mathbf{X}_1}}. \end{aligned} \tag{39}$$

Case 2. *We keep \mathbf{X}_1 and reject \mathbf{X}_2 when*

$$\begin{aligned} \Delta_{\mathbf{Y}, \mathbf{X}_2} \geq \Delta_{\mathbf{Y}, \mathbf{X}_1} &\iff C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} \geq I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}, \\ \delta_{\mathbf{Y}, \mathbf{X}_2} \geq \delta_{\mathbf{Y}, \mathbf{X}_1} &\iff \frac{C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1}}{C_{\mathbf{Y}, \mathbf{X}_1}} \geq \frac{I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}}{I_{\mathbf{Y}, \mathbf{X}_1}}. \end{aligned}$$

This result implies automatically that properties [P8] and [P9] are satisfied. Let us comment more precisely on the previous proposition.

- Case 1 holds when
 1. \mathbf{X}_2 is simpler than \mathbf{X}_1 (i.e., $C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} < 0$) and \mathbf{X}_2 is at least as informative as \mathbf{X}_1 (i.e., $I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1} \geq 0$).
 2. \mathbf{X}_2 and \mathbf{X}_1 have the same complexity (i.e., $C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} = 0$) and \mathbf{X}_2 is more informative than \mathbf{X}_1 (i.e., $I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1} > 0$).
 3. \mathbf{X}_2 is simpler and less informative than \mathbf{X}_1 and such that the absolute (respectively, relative) excess of complexity is less than the absolute (respectively, relative) gain of information, i.e., $C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} < I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1} < 0$ (respectively, $\frac{C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1}}{C_{\mathbf{Y}, \mathbf{X}_1}} < \frac{I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}}{I_{\mathbf{Y}, \mathbf{X}_1}} < 0$).
 4. \mathbf{X}_2 is more complex and more informative than \mathbf{X}_1 and such that the absolute (respectively, relative) excess of complexity is less than the absolute (respectively, relative) gain of information, i.e., $0 < C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} < I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}$ (respectively, $0 < \frac{C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1}}{C_{\mathbf{Y}, \mathbf{X}_1}} < \frac{I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1}}{I_{\mathbf{Y}, \mathbf{X}_1}}$).
- Case 2 holds when
 1. \mathbf{X}_2 is at least as complex as \mathbf{X}_1 (i.e., $C_{\mathbf{Y}, \mathbf{X}_2} - C_{\mathbf{Y}, \mathbf{X}_1} \geq 0$) and \mathbf{X}_2 is at most as informative as \mathbf{X}_1 (i.e., $I_{\mathbf{Y}, \mathbf{X}_2} - I_{\mathbf{Y}, \mathbf{X}_1} \leq 0$).

2. \mathbf{X}_2 is simpler and less informative than \mathbf{X}_1 and such that the absolute (respectively, relative) excess of complexity is greater than or equal to the absolute (respectively, relative) gain of information, i.e., $0 > C_{\mathbf{Y},\mathbf{X}_2} - C_{\mathbf{Y},\mathbf{X}_1} \geq I_{\mathbf{Y},\mathbf{X}_2} - I_{\mathbf{Y},\mathbf{X}_1}$ (respectively, $0 > \frac{C_{\mathbf{Y},\mathbf{X}_2} - C_{\mathbf{Y},\mathbf{X}_1}}{C_{\mathbf{Y},\mathbf{X}_1}} \geq \frac{I_{\mathbf{Y},\mathbf{X}_2} - I_{\mathbf{Y},\mathbf{X}_1}}{I_{\mathbf{Y},\mathbf{X}_1}}$).
3. \mathbf{X}_2 is more complex and more informative than \mathbf{X}_1 and such that the absolute (respectively, relative) excess of complexity is greater than or equal to the absolute (respectively, relative) gain of information, i.e., $C_{\mathbf{Y},\mathbf{X}_2} - C_{\mathbf{Y},\mathbf{X}_1} \geq I_{\mathbf{Y},\mathbf{X}_2} - I_{\mathbf{Y},\mathbf{X}_1} > 0$ (respectively, $\frac{C_{\mathbf{Y},\mathbf{X}_2} - C_{\mathbf{Y},\mathbf{X}_1}}{C_{\mathbf{Y},\mathbf{X}_1}} \geq \frac{I_{\mathbf{Y},\mathbf{X}_2} - I_{\mathbf{Y},\mathbf{X}_1}}{I_{\mathbf{Y},\mathbf{X}_1}} > 0$).

Proposition 12. Any complexity term C^α of the form (17) with $\alpha \in]0, 1[$ satisfies [P10].

Proof. Without loss of generality, the function $g(\cdot)$ that defines C^α is assumed to be increasing. It is clear that $H_{\mathbf{X}_2} \geq H_{\mathbf{X}_1} \Rightarrow C_{\mathbf{Y},\mathbf{X}_2}^\alpha \geq C_{\mathbf{Y},\mathbf{X}_1}^\alpha$ and respectively (by symmetry) $H_{\mathbf{X}_1} \geq H_{\mathbf{X}_2} \Rightarrow C_{\mathbf{Y},\mathbf{X}_1}^\alpha \geq C_{\mathbf{Y},\mathbf{X}_2}^\alpha$, which is equivalent to $C_{\mathbf{Y},\mathbf{X}_2}^\alpha > C_{\mathbf{Y},\mathbf{X}_1}^\alpha \Rightarrow H_{\mathbf{X}_2} > H_{\mathbf{X}_1}$. Thus, it remains to check that $C_{\mathbf{Y},\mathbf{X}_2}^\alpha = C_{\mathbf{Y},\mathbf{X}_1}^\alpha \Rightarrow H_{\mathbf{X}_2} = H_{\mathbf{X}_1}$:

$$\begin{aligned}
C_{\mathbf{Y},\mathbf{X}_2}^\alpha - C_{\mathbf{Y},\mathbf{X}_1}^\alpha = 0 &\iff \alpha(g(m_1) - g(m_2)) + (1 - \alpha)(g(M_1) - g(M_2)) = 0 \\
&\iff g(m_1) = g(m_2) \quad \text{and} \quad g(M_1) = g(M_2) \\
&\iff m_1 = m_2 \quad \text{and} \quad M_1 = M_2 \\
&\iff H_{\mathbf{X}_2} = H_{\mathbf{X}_1}.
\end{aligned} \tag{40}$$

Equation (40) is obtained by observing that $g(m_1) - g(m_2)$ and $g(M_1) - g(M_2)$ are of the same sign. \triangle

Remark 7. The complexity terms C^E and C^I (which corresponds to the case of any C^α with $\alpha = 0$) do not satisfy property [P10] in the general case. Indeed, there is no implication for C^E , and one can only prove that $H_{\mathbf{X}_1} \geq H_{\mathbf{X}_2} \Rightarrow C_{\mathbf{Y},\mathbf{X}_1}^I \geq C_{\mathbf{Y},\mathbf{X}_2}^I$. However, we point out that when $I_{\mathbf{Y},\mathbf{X}_1} = I_{\mathbf{Y},\mathbf{X}_2}$, then both C^E and C^I satisfy [P10].

More specifically, two frameworks may be of special interest.

- \mathbf{X}_2 is as informative as \mathbf{X}_1 (i.e., $I_{\mathbf{Y},\mathbf{X}_1} = I_{\mathbf{Y},\mathbf{X}_2}$): we expect to select the random variable with the smallest entropy. This is effectively what happens when [P10] is satisfied; from Proposition 12 and Remark 7 (in this framework), this is the case for C^\bullet with $\bullet = I, S, R, P, D$ in the general case and for C^E in this framework, since $H_{\mathbf{Y},\mathbf{X}_2} - H_{\mathbf{Y},\mathbf{X}_1} = H_{\mathbf{X}_2} - H_{\mathbf{X}_1}$.
- $\mathbf{X}_1 = g(\mathbf{X}_2)$ with g some surjective (but not injective) mapping: \mathbf{X}_2 is more complex than \mathbf{X}_1 , and \mathbf{X}_2 is at least as informative as \mathbf{X}_1 . Consequently, this case is nontrivial since both the absolute (respectively, relative) excess of complexity and absolute (respectively, relative) gain of information are competing. Let us give two important examples of such a context.
 1. Quantization problem: given a quantized version \mathbf{X}_1 of some (continuous) random variable with its associated partition \mathcal{A}_1 , the problem is to know whether some new quantized version \mathbf{X}_2 with an associated partition \mathcal{A}_2 finer than \mathcal{A}_1 should be preferred to predict \mathbf{Y} .
 2. Variable selection problem: assume that one wants to construct an ascending selection method. The vector \mathbf{X}_1 could represent some selected set of covariables and $\mathbf{X}_2 = (\mathbf{X}_1, \mathbf{X}'_2)$ a larger set of covariables. The aim is therefore to know if \mathbf{X}'_2 should be integrated to the selected set or not.

Some simple algorithms of quantization and selection methods are proposed in [15] using these results.

4.2. Around the Redundancy of Two Random Vectors \mathbf{X}_1 and \mathbf{X}_2

In the future use of an I -divergence or NI -divergence, one would expect that if two discrete-valued random vectors \mathbf{X}_1 and \mathbf{X}_2 have the same (or almost the same) information with respect to an I -divergence or NI -divergence, then both have the same effect on the prediction of another vector \mathbf{Y} . This requirement, expressed by property [P11], could be used for example in a variable selection problem in the context of discrimination to detect redundant variables.

In order to make property [P11] applicable for practical purpose, we may find interesting to have a bound on the difference $|\Delta_{\mathbf{Y},\mathbf{X}_1} - \Delta_{\mathbf{Y},\mathbf{X}_2}|$ (respectively, $|\delta_{\mathbf{Y},\mathbf{X}_1} - \delta_{\mathbf{Y},\mathbf{X}_2}|$) expressed in terms of $D_{\mathbf{X}_1,\mathbf{X}_2}^I$ (respectively, $d_{\mathbf{X}_1,\mathbf{X}_2}^I$). More precisely, the question may arise whether there exists a function $h(\cdot)$ satisfying $h(x) \rightarrow 0$ as $x \rightarrow 0$ and such that $|\Delta_{\mathbf{Y},\mathbf{X}_1} - \Delta_{\mathbf{Y},\mathbf{X}_2}| \leq h(D_{\mathbf{X}_1,\mathbf{X}_2}^I)$ (respectively, $|\delta_{\mathbf{Y},\mathbf{X}_1} - \delta_{\mathbf{Y},\mathbf{X}_2}| \leq h(d_{\mathbf{X}_1,\mathbf{X}_2}^I)$). Here, according to our examples, we only concentrate on linear functions $h(\cdot)$.

We then propose to translate property [P11] to an I -divergence Δ (respectively, NI -divergence δ) as follows:

[P11'(Υ, k)] There exists some positive constant k such that for all $(\mathbf{X}_1, \mathbf{X}_2) \in \Upsilon \subset \Gamma^2$ we have

$$|\Delta_{\mathbf{Y},\mathbf{X}_1} - \Delta_{\mathbf{Y},\mathbf{X}_2}| \leq kD_{\mathbf{X}_1,\mathbf{X}_2}^I. \tag{41}$$

As a first answer, let us note that if the I -divergence (respectively, NI -divergence) satisfies the triangle inequality [P6'($\Gamma^2, 1$)] or [P3'(Υ, k_1, k_2)], then it satisfies [P11'(Υ, k_2)] due to the following equivalent expression of the triangle inequality:

$$|D_{\mathbf{Y},\mathbf{X}_1} - D_{\mathbf{Y},\mathbf{X}_2}| \leq D_{\mathbf{X}_1,\mathbf{X}_2} \quad (\text{respectively, } |d_{\mathbf{Y},\mathbf{X}_1} - d_{\mathbf{Y},\mathbf{X}_2}| \leq d_{\mathbf{X}_1,\mathbf{X}_2}).$$

A priori, if an I -divergence or NI -divergence only satisfies [P6'(Γ^2, c)] with some $c > 1$, then this property does no more seem to be true: indeed, for all \mathbf{Y} , \mathbf{X}_1 , and \mathbf{X}_2 , one may prove for an I -divergence, for instance, that

$$|\Delta_{\mathbf{Y},\mathbf{X}_1} - \Delta_{\mathbf{Y},\mathbf{X}_2}| \leq c\Delta_{\mathbf{X}_1,\mathbf{X}_2} + (c - 1) \min(\Delta_{\mathbf{Y},\mathbf{X}_1}, \Delta_{\mathbf{Y},\mathbf{X}_2}) \not\leq c\Delta_{\mathbf{X}_1,\mathbf{X}_2}.$$

Actually, this apparently disappointing result only expresses that the “redundancy” property cannot (always) be derived from a triangle-type inequality.

The following proposition gives some sufficient conditions on a complexity term ensuring that the associated Δ and δ satisfy property [P11'].

Proposition 13. (i) Assume that there exists some positive constant \varkappa_1 such that the complexity term of an I -divergence for all $(\mathbf{X}_1, \mathbf{X}_2) \in \Upsilon$ satisfies

$$|C_{\mathbf{Y},\mathbf{X}_1} - C_{\mathbf{Y},\mathbf{X}_2}| \leq \varkappa_1 |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|. \tag{42}$$

Then Δ satisfies [P11'($\Upsilon, 1 + \varkappa_1$)].

(ii) If in addition there exists some positive constant \varkappa_2 such that for all $(\mathbf{X}_1, \mathbf{X}_2) \in \Upsilon$

$$\max(C_{\mathbf{Y},\mathbf{X}_1}, C_{\mathbf{Y},\mathbf{X}_2}) \geq \varkappa_2 C_{\mathbf{X}_1,\mathbf{X}_2}^I, \tag{43}$$

then the associated NI -divergence satisfies [P11'($\Upsilon, \frac{1 + \varkappa_1}{\varkappa_2}$)].

Proof. (i) Let us first write

$$|\Delta_{\mathbf{Y},\mathbf{X}_1} - \Delta_{\mathbf{Y},\mathbf{X}_2}| \leq |I_{\mathbf{Y},\mathbf{X}_1} - I_{\mathbf{Y},\mathbf{X}_2}| + |C_{\mathbf{Y},\mathbf{X}_1} - C_{\mathbf{Y},\mathbf{X}_2}|. \tag{44}$$

Now note that

$$I_{\mathbf{Y},\mathbf{X}_1} \geq I_{\mathbf{Y},\mathbf{X}_2} + I_{\mathbf{X}_1,\mathbf{X}_2} - H_{\mathbf{X}_2},$$

from which one can deduce

$$|I_{\mathbf{Y},\mathbf{X}_1} - I_{\mathbf{Y},\mathbf{X}_2}| \leq \max(H_{\mathbf{X}_1}, H_{\mathbf{X}_2}) - I_{\mathbf{X}_1,\mathbf{X}_2} = \max(H_{\mathbf{X}_1|\mathbf{X}_2}, H_{\mathbf{X}_2|\mathbf{X}_1}) = D_{\mathbf{X}_1,\mathbf{X}_2}^I. \tag{45}$$

The result is then obtained by combining (42), (44), and (45).

(ii) We can obtain the following result:

$$\begin{aligned} |\delta_{\mathbf{Y},\mathbf{X}_1} - \delta_{\mathbf{Y},\mathbf{X}_2}| &\leq \frac{\min(C_{\mathbf{Y},\mathbf{X}_1}, C_{\mathbf{Y},\mathbf{X}_2})(|I_{\mathbf{Y},\mathbf{X}_1} - I_{\mathbf{Y},\mathbf{X}_2}| + |C_{\mathbf{Y},\mathbf{X}_1} - C_{\mathbf{Y},\mathbf{X}_2}|)}{C_{\mathbf{Y},\mathbf{X}_1} C_{\mathbf{Y},\mathbf{X}_2}} \\ &\leq \frac{|I_{\mathbf{Y},\mathbf{X}_1} - I_{\mathbf{Y},\mathbf{X}_2}| + |C_{\mathbf{Y},\mathbf{X}_1} - C_{\mathbf{Y},\mathbf{X}_2}|}{\max(C_{\mathbf{Y},\mathbf{X}_1}, C_{\mathbf{Y},\mathbf{X}_2})}. \end{aligned}$$

The result then comes from (42), (43), and (45). Δ

Let us apply the previous result to our different examples.

Corollary 4. Let $\mathbf{X}_1, \mathbf{X}_2 \in \Gamma_\Theta$ with $\Theta = [c_1, c_2]$; define $\gamma_i, i = 1, 2$, such that $c_i = \gamma_i H_{\mathbf{Y}}$. Then

$$|\Delta_{\mathbf{Y},\mathbf{X}_1}^\bullet - \Delta_{\mathbf{Y},\mathbf{X}_2}^\bullet| \leq (1 + \varkappa_{1,\Theta}^\bullet) D_{\mathbf{X}_1,\mathbf{X}_2}^I \quad \text{and} \quad |\delta_{\mathbf{Y},\mathbf{X}_1}^\bullet - \delta_{\mathbf{Y},\mathbf{X}_2}^\bullet| \leq \frac{1 + \varkappa_{1,\Theta}^\bullet}{\varkappa_{2,\Theta}^\bullet} d_{\mathbf{X}_1,\mathbf{X}_2}^I, \tag{46}$$

where \bullet stands for S, R, P , and D , and where the different constants are expressed by

\bullet	$\varkappa_{1,\Theta}^\bullet$	$\varkappa_{2,\Theta}^\bullet$
S	α^\vee	$(1 - \alpha) + \alpha\gamma_{1,2}$
R	$\alpha^{\vee 2} + \frac{\alpha(1 - \alpha)}{\sqrt{\gamma_1}}$	$((1 - \alpha) + \alpha\sqrt{\gamma_{1,2}})^2$
P	$\max\left(\frac{1 - \alpha}{\gamma_1^\alpha}, \frac{\alpha}{\gamma_1^{1 - \alpha}}, \mathbf{1}_{]0,1]}(\gamma_1)\right)$	$\gamma_{1,2}^\alpha$
D	$\frac{\alpha^\vee}{(\alpha^\wedge)^2} \frac{1}{(1 + \gamma_{1,2})^2}$	$\left(\frac{\alpha}{\gamma_{1,2}} + (1 - \alpha)\right)^{-1}$

with $\gamma_{1,2} = \min\left(\gamma_1, \frac{1}{\gamma_2}\right)$.

Proof. For the sake of simplicity, let us denote $m_i = \min(H_{\mathbf{Y}}, H_{\mathbf{X}_i})$ and $M_i = \max(H_{\mathbf{Y}}, H_{\mathbf{X}_i})$ for $i = 1, 2$. The proofs mainly rely upon the two following tools:

- Since $m_1 - m_2$ and $M_1 - M_2$ are of the same sign, we have

$$|(m_1 - m_2)| + |M_1 - M_2| = |(m_1 - m_2) + (M_1 - M_2)| = |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|$$

and then for any $f_1, f_2 \geq 0$

$$|f_1(m_1 - m_2) + f_2(M_1 - M_2)| \leq f^\vee |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|, \tag{47}$$

with $f^\vee = \max(f_1, f_2)$.

- For $i = 1, 2$, we have

$$m_i \geq \left\{ \begin{array}{l} \min(1, \gamma_1) M_i \\ \min\left(1, \frac{1}{\gamma_2}\right) M_i \end{array} \right\} \geq \min\left(1, \gamma_1, \frac{1}{\gamma_2}\right) M_i = \gamma_{1,2} M_i,$$

since $\gamma_{1,2} < 1$ as a direct consequence of $\gamma_1 < \gamma_2$.

For different cases, we need in particular to check (43), which holds whenever there exists some positive constant \varkappa such that $C_{\mathbf{Y},\mathbf{X}_i}^{\bullet,\alpha} \geq \varkappa M_i$ for $i = 1, 2$, since $\max(C_{\mathbf{Y},\mathbf{X}_1}^{\bullet,\alpha}, C_{\mathbf{Y},\mathbf{X}_2}^{\bullet,\alpha}) \geq \varkappa \max(M_1, M_2) \geq \varkappa C_{\mathbf{X}_1, \mathbf{X}_2}^I$.

- Complexity term $C^{S,\alpha}$: we have

$$\begin{aligned} |C_{\mathbf{Y},\mathbf{X}_1}^{S,\alpha} - C_{\mathbf{Y},\mathbf{X}_2}^{S,\alpha}| &= |\alpha m_1 + (1 - \alpha)M_1 - \alpha m_2 - (1 - \alpha)M_2| \\ &= |\alpha(m_1 - m_2) + (1 - \alpha)(M_1 - M_2)| \\ &\leq \alpha^\vee |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}| \end{aligned}$$

from (47), with $f_1 = \alpha$ and $f_2 = (1 - \alpha)$. Moreover, for $i = 1, 2$,

$$C_{\mathbf{Y},\mathbf{X}_i}^{S,\alpha} = \alpha m_i + (1 - \alpha)M_i \geq ((1 - \alpha) + \alpha\gamma_{1,2})M_i.$$

- Complexity term $C^{R,\alpha}$: we have

$$\left| C_{\mathbf{Y},\mathbf{X}_1}^{R,\alpha} - C_{\mathbf{Y},\mathbf{X}_2}^{R,\alpha} \right| = \left| \alpha^2(m_1 - m_2) + (1 - \alpha)^2(M_1 - M_2) + 2\alpha(1 - \alpha)\sqrt{H_{\mathbf{Y}}}\left(\sqrt{H_{\mathbf{X}_1}} - \sqrt{H_{\mathbf{X}_2}}\right) \right|.$$

Furthermore, we may obtain

$$\left| \alpha^2(m_1 - m_2) + (1 - \alpha)^2(M_1 - M_2) \right| \leq \alpha^{\vee 2} |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|$$

and

$$\left| \sqrt{H_{\mathbf{Y}}}\left(\sqrt{H_{\mathbf{X}_1}} - \sqrt{H_{\mathbf{X}_2}}\right) \right| = \frac{\sqrt{H_{\mathbf{Y}}}}{2\sqrt{\min(H_{\mathbf{X}_1}, H_{\mathbf{X}_2})}} |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}| \leq \frac{1}{2\sqrt{\gamma_1}} |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|.$$

Hence,

$$\left| C_{\mathbf{Y},\mathbf{X}_1}^{R,\alpha} - C_{\mathbf{Y},\mathbf{X}_2}^{R,\alpha} \right| \leq \left(\alpha^{\vee 2} + \frac{\alpha(1 - \alpha)}{\sqrt{\gamma_1}} \right) |H_{\mathbf{X}_1} - H_{\mathbf{X}_2}|$$

from (47), with $f_1 = \alpha^2$ and $f_2 = (1 - \alpha)^2$. Moreover, one can prove for $i = 1, 2$

$$C_{\mathbf{Y},\mathbf{X}_i}^{R,\alpha} = (\alpha\sqrt{m_i} + (1 - \alpha)\sqrt{M_i})^2 \geq ((1 - \alpha) + \alpha\sqrt{\gamma_{1,2}})^2 M_i.$$

- Complexity term $C^{P,\alpha}$: we have (by assuming that $H_{\mathbf{X}_2} > H_{\mathbf{X}_1}$)

$$\begin{aligned} \left| C_{\mathbf{Y},\mathbf{X}_1}^{P,\alpha} - C_{\mathbf{Y},\mathbf{X}_2}^{P,\alpha} \right| &= \left| m_1^\alpha M_1^{1-\alpha} - m_2^\alpha M_2^{1-\alpha} \right| \\ &= \begin{cases} H_{\mathbf{Y}}^\alpha (H_{\mathbf{X}_2}^{1-\alpha} - H_{\mathbf{X}_1}^{1-\alpha}) & \text{if } H_{\mathbf{Y}} \leq \min(H_{\mathbf{X}_1}, H_{\mathbf{X}_2}), \\ H_{\mathbf{Y}}^{1-\alpha} (H_{\mathbf{X}_2}^\alpha - H_{\mathbf{X}_1}^\alpha) & \text{if } H_{\mathbf{Y}} \geq \max(H_{\mathbf{X}_1}, H_{\mathbf{X}_2}), \\ H_{\mathbf{Y}}^\alpha H_{\mathbf{X}_2}^{1-\alpha} - H_{\mathbf{X}_1}^\alpha H_{\mathbf{Y}}^{1-\alpha} & \text{otherwise.} \end{cases} \end{aligned}$$

Note that the third case cannot occur if $\gamma_1 \geq 1$. We have

$$\begin{aligned} \left| C_{\mathbf{Y},\mathbf{X}_1}^{P,\alpha} - C_{\mathbf{Y},\mathbf{X}_2}^{P,\alpha} \right| &\leq \begin{cases} \frac{1 - \alpha}{\gamma_1^\alpha} (H_{\mathbf{X}_2} - H_{\mathbf{X}_1}) & \text{if } H_{\mathbf{Y}} \leq \min(H_{\mathbf{X}_1}, H_{\mathbf{X}_2}), \\ \frac{\alpha}{\gamma_1^{1-\alpha}} (H_{\mathbf{X}_2} - H_{\mathbf{X}_1}) & \text{if } H_{\mathbf{Y}} \geq \max(H_{\mathbf{X}_1}, H_{\mathbf{X}_2}), \\ H_{\mathbf{X}_2} - H_{\mathbf{X}_1} & \text{otherwise} \end{cases} \\ &\leq \max\left(\frac{1 - \alpha}{\gamma_1^\alpha}, \frac{\alpha}{\gamma_1^{1-\alpha}}, \mathbf{1}_{]0,1[}(\gamma_1)\right) |H_{\mathbf{X}_2} - H_{\mathbf{X}_1}|. \end{aligned}$$

Moreover, we may obtain for $i = 1, 2$

$$C_{\mathbf{Y},\mathbf{X}_i}^{P,\alpha} = m_i^\alpha M_i^{1-\alpha} \geq \gamma_{1,2}^\alpha M_i.$$

- Complexity term $C^{D,\alpha}$: we have

$$\begin{aligned} \left| C_{\mathbf{Y}, \mathbf{X}_1}^{D,\alpha} - C_{\mathbf{Y}, \mathbf{X}_2}^{D,\alpha} \right| &= \frac{|\alpha M_1 M_2 (m_1 - m_2) + (1 - \alpha) m_1 m_2 (M_1 - M_2)|}{(\alpha M_1 + (1 - \alpha) m_1)(\alpha M_2 + (1 - \alpha) m_2)} \\ &\leq \frac{\alpha^\vee}{(\alpha^\wedge)^2} \frac{M_1 M_2}{(m_1 + M_1)(m_2 + M_2)} |H_{\mathbf{X}_2} - H_{\mathbf{X}_1}| \\ &\leq \frac{\alpha^\vee}{(\alpha^\wedge)^2} \frac{1}{(1 + \gamma_{1,2})^2} |H_{\mathbf{X}_2} - H_{\mathbf{X}_1}|. \end{aligned} \quad (48)$$

Equation (48) is obtained using (47) with $f_1 = \alpha M_1 M_2$ and $f_2 = (1 - \alpha) m_1 m_2$. Finally, we also have for $i = 1, 2$

$$C_{\mathbf{Y}, \mathbf{X}_i}^{D,\alpha} = \left(\frac{\alpha}{m_i} + \frac{1 - \alpha}{M_i} \right)^{-1} \geq \left(\frac{\alpha}{\gamma_{1,2}} + (1 - \alpha) \right)^{-1} M_i. \quad \triangle$$

Remark 8. Note that when $\alpha \leq \frac{1}{2}$, the measure $\Delta^{S,\alpha}$ is a metric and so we derive (46) directly from [P3'].

The authors would like to thank the referee whose suggestions and fruitful comments helped them to improve the paper.

REFERENCES

1. Shannon, C.E., A Mathematical Theory of Communication (continued), *Bell Syst. Tech. J.*, 1948, vol. 27, no. 4, pp. 623–656.
2. Crutchfield J.P., Information and Its Metric, *Nonlinear Structures in Physical Systems: Pattern Formation, Chaos, and Waves (Proc. 2nd Woodward Conf., San Jose State Univ., 1989)*, Lam, L. and Morris, H.C., Eds., New York: Springer, 1990, pp. 119–130.
3. Hillman, C., A Formal Theory of Information: I. Statics, *Preprint*, 1998. Available from <http://citeseer.ist.psu.edu/89520.html>.
4. Bennett, C.H., Gács, P., Li, M., Vitányi, P.M.B., and Zurek, W.H., Information Distance, *IEEE Trans. Inform. Theory*, 1998, vol. 44, no. 4, pp. 1407–1423.
5. Li, M., Chen, X., Li, X., Ma, B., and Vitányi, P.M.B., The Similarity Metric, *IEEE Trans. Inform. Theory*, 2004, vol. 50, no. 12, pp. 3250–3264.
6. Cilibrasi, R. and Vitányi, P.M.B., Clustering by Compression, *IEEE Trans. Inform. Theory*, 2005, vol. 51, no. 4, pp. 1523–1545.
7. Li, M., Badger, J.H., Chen, X., Kwong, S., Kearney, P., and Zhang, H., An Information-Based Sequence Distance and Its Application to Whole Mitochondrial Genome Phylogeny, *Bioinformatics*, 2001, vol. 17, no. 1, pp. 149–154.
8. Cilibrasi, R., Vitányi, P., and de Wolf, R., Algorithmic Clustering of Music, 2003, e-print arxiv.org/cs.SD/0303025.
9. Cilibrasi, R. and Vitányi, P., Automatic Meaning Discovery Using Google, 2005, e-print arxiv.org/cs.CL/0412098.
10. Grünwald, P. and Vitányi, P., Shannon Information and Kolmogorov Complexity, 2004, e-print arxiv.org/cs.IT/0410002.
11. Leung-Yan-Cheong, S.K. and Cover, T.M., Some Equivalences between Shannon Entropy and Kolmogorov Complexity, *IEEE Trans. on Inf. Theory*, 1978, vol. 24, no. 3, pp. 331–338.
12. Hammer, D., Romashchenko, A., Shen, A., and Vereshchagin, N., Inequalities for Shannon Entropy and Kolmogorov Complexity, *J. Comput. Syst. Sci.*, 2000, vol. 60, no. 2, pp. 442–464.

13. Kraskov, A., Stögbauer, H., Andrzejak, R.G., and Grassberger, P., Hierarchical Clustering Based on Mutual Information, 2003, e-print arxiv.org/q-bio/0311039.
14. Liu, H. and Motoda, H., *Feature Selection for Knowledge Discovery and Data Mining*, Boston: Kluwer, 1998.
15. Robineau, J.-F., Méthodes de sélection de variables (parmi un grand nombre) dans un cadre de discrimination, *PhD Thesis*, Grenoble, France: Univ. Joseph Fourier, 2004.
16. Granger, C.W., Maasoumi, E., and Racine, J., A Dependence Metric for Possibly Nonlinear Processes, *J. Time Ser. Anal.*, 2004, vol. 25, no. 5, pp. 649–669.
17. Cover, T.M. and Thomas, J.A., *Elements of Information Theory*, New York: Wiley, 1991.
18. Kaltchenko, A., Algorithms for Estimation of Information Distance with Application to Bioinformatics and Linguistics, in *Proc. 2004 Canadian Conf. on Electrical and Computer Engineering*, 2004, vol. 4, p. 2255. Available from arxiv.org/cs.CC/0404039.
19. Li, M. and Vitányi, P., *An Introduction to Kolmogorov Complexity and Its Applications*, New York: Springer, 1997, 2nd ed.
20. Li, M, Chen, X., Li, X., Ma, B., and Vitányi, P., The Similarity Metric, in *Proc. 14th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA-03)*, New York, 2003, pp. 863–872.
21. Ullah, A., Entropy, Divergence and Distance Measures with Econometric Applications, *J. Stat. Plan. Inference*, 1996, vol. 49, no. 1, pp. 137–162.