



# The Decomposition of Similarity-Based Diversity and its Bias Correction

Eric Marcon, Zhiyi Zhang, Bruno Hérault

► **To cite this version:**

Eric Marcon, Zhiyi Zhang, Bruno Hérault. The Decomposition of Similarity-Based Diversity and its Bias Correction. 2014. <hal-00989454v3>

**HAL Id: hal-00989454**

**<https://hal-agroparistech.archives-ouvertes.fr/hal-00989454v3>**

Submitted on 18 Mar 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The Decomposition of Similarity-Based Diversity and its Bias Correction

Eric Marcon<sup>1\*</sup>, Zhiyi Zhang<sup>2</sup>, Bruno Hérault<sup>3</sup>

## Abstract

Similarity-based diversity defined by Leinster and Cobbold (2012) encompasses classical measures of diversity, including the number of species, Shannon and Simpson diversity, and also phylogenetic and functional diversity as long as distances between species are ultrametric. We derive two estimators to allow applying it to real, often under-sampled data, and its decomposition into alpha, beta and gamma diversity when an assemblage of communities is considered.

## Keywords

Entropy, Biodiversity, Phylogenetic diversity, Functional diversity

<sup>1</sup> AgroParisTech, UMR Écologie des Forêts de Guyane, BP 709, F-97310 Kourou, French Guiana.

<sup>2</sup> Department of Mathematics and Statistics, University of North Carolina at Charlotte, Charlotte, NC 28223

<sup>3</sup> CIRAD, UMR Écologie des Forêts de Guyane, F-97379 Kourou Cedex, France

\*Corresponding author: Eric.Marcon@ecofog.gf

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Methods</b>	<b>2</b>
2.1	Notations	2
2.2	Definition of ${}^qD^Z$ and ${}^qH^Z$	2
2.3	Similarity matrix	3
2.4	Estimation bias correction	3
	Chao-Shen estimator • Alternative estimator	
2.5	Decomposition of entropy	5
2.6	Decomposition of diversity	5
<b>3</b>	<b>Results</b>	<b>5</b>
3.1	Test of bias correction	5
3.2	Partitioning functional diversity	6
<b>4</b>	<b>Discussion and Conclusion</b>	<b>7</b>
<b>5</b>	<b>Acknowledgments</b>	<b>9</b>
	<b>References</b>	<b>9</b>

## 1. Introduction

It is still surprising that a so widely used ecological concept, such as the concept of diversity, is still so debated. Ecology is rich of such recurring debates on concepts (e.g. the “ecological resilience” concept [1], or the “functional trait” concept [2]) that could nevertheless be seen as part of a fundamental theoretical corpus. This may be why one may consider that ecology is still a “young” science [3]. Practically, the notion of diversity is more or less consensual among field ecologists. Very roughly, the diversity of an ecological system is even higher than the individuals in the system are “different”. Two ecological systems are even more different than their individuals are, too. Problems arise when we move from words to mathematical quantification. Entropy, viewed as

the average surprise provided by the data [4, 5], paved the way of a coherent theoretical framework [6] that respects a set of meaningful axioms [7, 8] including discarding species identities, continuity relatively to the probabilities of occurrence of species and the Pigou-Dalton [9] property (replacing an individual of a more abundant species by an individual of a less abundant species increases diversity). Generalized entropy, namely HCDT entropy [10–12], allows moving the cursor, the order of diversity  $q$ , from rare to abundant species. Conversion of entropy into Hill numbers [13–15] provides the effective numbers, *i.e.* the number of equally-frequent entities (e.g. species or communities) which would yield the same diversity value as the data. The product of the effective number of communities [16], namely  $\beta$  diversity, by  $\alpha$  diversity (the average diversity of communities) is equal to the  $\gamma$  diversity of a mixture of communities (*i.e.* a meta-community). HCDT entropy and Hill numbers have been applied to phylogenetic and functional diversity, and, though always perfectible, solutions have been proposed to most issues: the correspondence between HCDT entropy and Hill numbers [17, 18], the decomposition of diversity according to Jost [16], Chiu *et al.* [19] or Routledge [18, 20], providing different definitions of  $\alpha$  diversity, and robust estimators [21–24].

Recently, a major step forward was made by Leinster and Cobbold [25] who introduced a general measure of diversity denoted  ${}^qD^Z$ , the *similarity-based diversity*. Encompassing all previously cited measures (at least when distances measuring the difference between species are ultrametric, see [26], Appendix S1, for a discussion), it allows a direct measure of diversity as the inverse of the average ordinariness of species. The latter is defined as its average similarity with other species: the topology of species, *i.e.* the place of this species in a multivariate space, does rely on species-to-species similarity. Parameterizing the definition of the average [27] allows

choosing the importance of ordinary species. The great improvement of this measure is to build on a species-to-species distance matrix directly. This is a very desirable property that allows taking into account not only evolutionary or functional distances but also any other kind of paired relationships [25]. Let's focus on functional diversity. While evolutionary relationships are naturally represented by a phylogenetic tree, functional diversity is more generally calculated from a matrix of multivariate distances, which is often highly distorted later when transformed into a functional dendrogram [28]. In addition to being problematic for the calculation of the functional diversity, this over-simplification of the raw trait matrix is extremely frustrating for field ecologists knowing the many difficulties encountered to measure these batteries of functional traits, especially in hyper-diverse ecosystems [29]. Knowing that single trait values or single trait variation axis are often inadequate predictors of the species fundamental niche [30], losing the multivariate information when estimating diversity prevents taking into account unique combinations of traits that influence ecosystems that are responsible for the ecosystem effect of a given species [31]. Up today, this *similarity-based diversity*  ${}^qD^Z$  still lacks both (i) a robust estimator and (ii) a decomposition framework. First, it has long been recognized that the observed diversity depends on the sample size [32], so that estimation-bias corrected estimators are required [33, 34]. Second, Leinster and Cobbold have said very little about  ${}^qD^Z$  decomposition into  $\alpha$ ,  $\beta$  and  $\gamma$  components. If the path is relatively well marked for diversity indices based on dendrograms, it remains to extend this analysis to the similarity-sensitive ones.

In this paper, we first recall the definitions and duality between similarity-based diversity and Ricotta and Szeidl's [35] entropy  ${}^qH^Z$ . We discuss the definition of the similarity matrix in depth. We explicit its parametric construction argued by Leinster and Cobbold and the differences with an alternative approach that lead Chiu and Chao [36] to propose a different measure of functional diversity. We propose two estimation-bias corrected estimators: the first one is an implementation of the Horvitz-Thompson [37] estimator, the other is built on the estimation of the ordinariness of unobserved species. Then we derive the decomposition of both  ${}^qD^Z$  and  ${}^qH^Z$  and provide a definition of  $\beta$  diversity and entropy.

## 2. Methods

### 2.1 Notations

Consider a random sample taken from a meta-community made of several local communities. Abundances of species in each local community is denoted  $n_{s,i}$  ( $s = 1, 2, \dots, S$  is the index of species,  $i$  the index of communities).  $n_s$  is the number of individuals of species  $s$  in the meta-community,  $n_i$  the number of individuals sampled in local community  $i$  and  $n$  the total number. The same notations are used for probabilities of occurrence  $p_{s,i}$  which are unknown but estimated with  $\hat{p}_{s,i} = n_{s,i}/n_i$ . Community weights are  $w_i$ : they may be equal to  $n_i/n$  but any positive values summing to 1 are allowed.

We assume that  $p_s = \sum_i w_i p_{s,i}$  for all species. This may be understood as the definition of our meta-community: the assemblage of communities whose species probabilities are the weighted average of those of communities. Diversity of the meta-community is  $\gamma$  diversity. Diversity of local communities is  $\alpha$  diversity.

Species similarity is introduced as a square matrix  $\mathbf{Z}$  of dimension  $S \times S$  whose elements  $z_{s,t}$  are the similarity between species  $s$  and  $t$ . Similarity is between 0 and 1 by definition, and  $z_{s,s} = 1$ : any species is completely similar to itself. A detailed presentation of the possible matrices can be found in [25] and is not repeated here. A matrix of particular interest is  $\mathbf{Z} = \mathbf{I}_S$ , the identity matrix where each species is completely different from the others, used to measure neutral diversity. The definition of  ${}^qD^Z$  can be extended to relatedness matrices, which only require that all terms  $z_{s,t}$  are positive, and the diagonal terms strictly positive (see appendix A5 of [25]). Our results remain valid for relatedness matrices except for our new estimator which explicitly supposes that a species similarity with itself is 1.

### 2.2 Definition of ${}^qD^Z$ and ${}^qH^Z$

A species ordinariness is defined as its average similarity in the community:

$$(\mathbf{Zp})_s = \sum_t p_t z_{s,t} \quad (1)$$

${}^qD^Z$  is the inverse of the generalized mean of order  $q-1$  [27] of the community's species ordinarienesses:

$${}^qD^Z = \left( \sum_s p_s (\mathbf{Zp})_s^{q-1} \right)^{\frac{1}{1-q}}; q \neq 1 \quad (2)$$

$${}^1D^Z = \frac{1}{\prod_s (\mathbf{Zp})_s^{p_s}} \quad (3)$$

In other words,  $1/{}^qD^Z$  is the average ordinariness of species of the community. The generalized mean allows giving more or less importance to *ordinary* species. Rarity is not concerned by the parameter:  $p_s$  is at power 1. When species-neutral diversity is addressed, ordinariness reduces to the species probability so both notions converge, but this not the general case. Chiu and Chao [36] have another approach: they define functional diversity as the weighted average distance between pairs of species, where the weight of the frequency of pairs is parameterized by its power  $q$ . Then,  $q$  allows giving more or less importance to rare pairs of species (hence to rare species). Both approaches are valid, but quite different. In similarity-based diversity, low values of  $q$  focus on original (understood as "not ordinary") species, high values on ordinary species.

${}^qD^Z$  is the deformed exponential of order  $q$  of the entropy of Ricotta and Szeidl [35]:

$${}^qH^Z = \frac{1 - \sum_s p_s (\mathbf{Zp})_s^{q-1}}{q-1} = \sum_s p_s \ln_q \frac{1}{(\mathbf{Zp})_s}; q \neq 1 \quad (4)$$

$${}^1H^Z = - \sum_s p_s \ln(\mathbf{Zp})_s \quad (5)$$

Diversity is a monotonic, increasing function of entropy, which is more easily written with the formalism of deformed logarithms [38]. For  $q \neq 1$ :

$$\ln_q x = \frac{x^{1-q} - 1}{1-q} \quad (6)$$

$$e_q^x = [1 + (1-q)x]^{1/(1-q)} \quad (7)$$

The relation between entropy and diversity is the same as that of HCDT entropy and Hill numbers:

$$\ln_q {}^qD^Z = {}^qH^Z \quad (8)$$

$${}^qD^Z = e_q^{{}^qH^Z} \quad (9)$$

### 2.3 Similarity matrix

A euclidean distance matrix  $\Delta$  between species, whose elements  $\delta_{s,t}$  are the distance between species  $s$  and  $t$ , can be obtained by a generalization of Gower's distance [39] from the values of traits measured on each species. An example is given below. The simplest way to obtain a similarity matrix from a distance matrix is to normalize it and take its complement to 1:  $z_{s,t} = 1 - \delta_{s,t}/\max(\Delta)$ . We denote  $d_{s,t}$  the normalized distances.

We follow Leinster and Cobbold to suggest  $z_{s,t} = e^{-ud_{s,t}}$  where  $u$  is a positive constant that can be interpreted as a scaling factor of the distance matrix. The negative exponential transformation means that the similarity between species decreases exponentially with respect to distance rather than linearly. From a theoretical point of view, this definition is supported by the theory of categories [40] which provides several useful results. One of them is that the maximum value of  ${}^qD^Z$  (given  $Z$ ) is obtained for a unique distribution of probabilities whatever  $q$  [41], and this value is the magnitude of the subset of euclidean space defined by the distance matrix. Magnitude is an invariant of metric spaces, analogous to the cardinality of sets. Although this theory is not necessary to define what diversity is, it is particularly appealing because of the evident analogy with the effective number of species. Empirically,  $u$  controls the convexity of the transformation

of distances into similarities. When  $u = 0$ , similarity is 1 whatever the distance.  $u = 1$  is close to a linear transformation. When  $u$  increases beyond 1, similarities are more and more packed close to 0. When  $u \rightarrow +\infty$ , the similarity matrix tends to the identity matrix and  ${}^qD^Z$  tends to species-neutral diversity. So  $u$  is a parameter that controls for the distribution of similarity values. This may be seen as an issue, since it adds one degree of freedom to the measure of diversity. On the other hand, building a distance matrix relies on arbitrary choices of methods. The ordination of distances clearly has an ecological meaning (more distant species are less similar) but the distribution of distances has less support. The parameter  $u$  allows to chose how similar to each other species are.

### 2.4 Estimation bias correction

As many other measures of diversity,  ${}^qD^Z$  suffers estimation bias [42]. The plug-in estimator, built by simply plugging  $\hat{p}_s$  into the formula of  ${}^qD^Z$ , is biased for two reasons: unobserved species [22] and the nonlinearity of diversity with respect to probabilities [43].

Bias-corrected estimators exist for HCDT entropy [20], including the Horvitz and Thompson [37] estimator first adapted by Chao and Shen [22] for Shannon entropy. We use it here to estimate  ${}^qH^Z$ . We then propose a new estimator of  ${}^qD^Z$ .

#### 2.4.1 Chao-Shen estimator

${}^qH^Z$  is the sum over species of the measure  $h_s = p_s \ln_q \frac{1}{(\mathbf{Zp})_s}$ . Unobserved species are responsible for unobserved terms of the sum. An unbiased estimator of such a sum has been derived by Horvitz and Thompson [37]: each observed term is divided by the probability for the species to be sampled: the lower it is, the higher  $h_s$  is weighted in the sum.

Chao and Shen [22] proposed to combine it with the estimator of the sample coverage. The total probability of occurrence of unobserved species is by definition 1 minus the sample coverage [44], which can be estimated from the data following Zhang and Huang [45]:

$$\hat{C} = 1 - \sum_{r=1}^n (-1)^{r+1} \binom{n}{r}^{-1} f_r \quad (10)$$

$n$  is the sample size and  $f_r$  the number of species observed  $r$  times in the sample. This is an improvement of the well-known Turing's formula  $\hat{C} = 1 - f_1/n$ . The observed species probabilities can be better estimated by  $\{\hat{C}\hat{p}_1, \dots, \hat{C}\hat{p}_T\}$  where  $T$  is the number of observed species.  $1 - \hat{C}$  is the estimated total probability of unobserved species.

We also have to estimate  $(\mathbf{Zp})_s$ . The observed  $Z$  matrix lacks the  $S - T$  columns of unobserved species, and the vector of probabilities lacks the corresponding probabilities  $\{\hat{p}_{T+1}, \dots, \hat{p}_S\}$ . Nothing is known about the similarity of unobserved species with observed ones. A reasonable assumption is that their average similarity with any other species is identical to the average similarity between observed species,

$\bar{z} = (\sum_{s,t \neq s} z_{s,t}) / [T(T-1)]$ . Since the sum of missing probabilities of species is known to be  $1 - \hat{C}$ , an estimator of  $(\mathbf{Zp})_s$  is:

$$(\widetilde{\mathbf{Zp}})_s = \sum_t \hat{C} \hat{p}_t z_{s,t} + (1 - \hat{C}) \bar{z} \quad (11)$$

The first term is the sum over observed species of their similarity with species  $s$  and the second term is the expected sum over unobserved species. Plugging the estimator of  $p_s$  and that of  $(\mathbf{Zp})_s$  into the Horvitz-Thompson estimator, an estimation-bias corrected estimator of  ${}^q H^Z$  is:

$${}^q \tilde{H}^Z = \sum_s \frac{\hat{C} \hat{p}_s \ln_q \frac{1}{(\widetilde{\mathbf{Zp}})_s}}{1 - (1 - \hat{C} \hat{p}_s)^n} \quad (12)$$

The estimator of diversity is  ${}^q \tilde{D}^Z = e^{q \tilde{H}^Z}$ .

#### 2.4.2 Alternative estimator

Instead of correcting the estimator for unobserved species, we now want to explicitly estimate their contribution to the value of diversity.

The observed  $\mathbf{Z}$  matrix lacks  $S - T$  lines and columns. The actual  $\mathbf{Z}$  matrix is (only bold elements are known):

$$\mathbf{Z} = \begin{pmatrix} \mathbf{1} & \cdots & \mathbf{z}_{1,T} & z_{1,T+1} = \bar{z} & \cdots & z_{1,S} = \bar{z} \\ \vdots & \ddots & \vdots & \vdots & & \vdots \\ \mathbf{z}_{T,1} & \cdots & \mathbf{1} & \bar{z} & & \bar{z} \\ z_{T+1,1} = \bar{z} & \cdots & \bar{z} & 1 & & \bar{z} \\ \vdots & & & & \ddots & \vdots \\ z_{S,1} = \bar{z} & \cdots & \bar{z} & \bar{z} & \cdots & 1 \end{pmatrix} \quad (13)$$

If  $q \neq 1$ , diversity can split into two terms:

$${}^q D^Z = \left[ \sum_{s \leq T} p_s (\mathbf{Zp})_s^{q-1} + \sum_{t > T} p_t (\mathbf{Zp})_t^{q-1} \right]^{\frac{1}{1-q}} = (K + U)^{\frac{1}{1-q}} \quad (14)$$

The first term  $K$  is easy to estimate, as in the Chao-Shen estimator:

$$\hat{K} = \sum_{s \leq T} \hat{C} \hat{p}_s \left[ \left( \sum_{t \leq T} \hat{C} \hat{p}_t z_{s,t} \right) + (1 - \hat{C}) \bar{z} \right]^{q-1} \quad (15)$$

The second one is more problematic:

$$U = \sum_{t > T} p_t [\bar{z}(1 - p_t) + p_t]^{q-1} \quad (16)$$

The number of terms is unknown, as is each  $p_t$ , but  $U$  can be estimated following Zhang and Grabchak [24]. We define:

$$V = \sum_s p_s [\bar{z}(1 - p_s) + p_s]^{q-1} \quad (17)$$

$V$  is a linear diversity index, using the terminology of Zhang and Grabchak. It can be estimated (derivation in Appendix S1) by:

$$\hat{V} = 1 + \sum_{s \leq T} \frac{n_s}{n} \sum_{v=1}^{n-n_s} (1 - \bar{z})^v \left[ \prod_{i=1}^v \frac{i-q}{i} \right] \left[ \prod_{j=1}^v \left( 1 - \frac{n_s-1}{n-j} \right) \right] \quad (18)$$

$U$  is the sum of the terms concerning unobserved species, so observed terms must be subtracted from  $\hat{V}$ . The improved estimator of  ${}^q D^Z$  is thus:

$${}^q \tilde{D}^Z = \left\{ \hat{K} + \hat{V} - \sum_{s \leq T} \hat{C} \hat{p}_s [\bar{z}(1 - \hat{C} \hat{p}_s) + \hat{C} \hat{p}_s]^{q-1} \right\}^{\frac{1}{1-q}} \quad (19)$$

The estimator of  ${}^q H^Z$  is:

$${}^q \tilde{H}^Z = \frac{\hat{K} + \hat{V} - \sum_{s \leq T} \hat{C} \hat{p}_s [\bar{z}(1 - \hat{C} \hat{p}_s) + \hat{C} \hat{p}_s]^{q-1} - 1}{1 - q} \quad (20)$$

If  $q = 1$ , estimating entropy is easier. It can be split into two terms:

$${}^1 H^Z = - \sum_{s \leq T} p_s \ln (\mathbf{Zp})_s - \sum_{t > T} p_t \ln (\mathbf{Zp})_t = L + X \quad (21)$$

As above, the first term is estimated from the data:

$$\hat{L} = - \sum_{s \leq T} \hat{C} \hat{p}_s \ln \left[ \left( \sum_{t \leq T} \hat{C} \hat{p}_t z_{s,t} \right) + (1 - \hat{C}) \bar{z} \right] \quad (22)$$

The second term addresses unobserved species:

$$X = - \sum_{t > T} p_t \ln [\bar{z}(1 - p_t) + p_t] \quad (23)$$

We define the linear diversity index:

$$W = - \sum_s p_s \ln [\bar{z}(1 - p_s) + p_s] \quad (24)$$

$W$  can be estimated by (Appendix S1):

$$\hat{W} = \sum_{s \leq T} \frac{n_s}{n} \sum_{v=1}^{n-n_s} \frac{(1-\bar{z})^v}{v} \left[ \prod_{j=1}^v \left( 1 - \frac{n_s-1}{n-j} \right) \right] \quad (25)$$

Finally:

$${}^1\hat{H}^Z = \hat{L} + \hat{W} + \sum_{s \leq T} \hat{C}\hat{p}_s \ln [\bar{z}(1-\hat{C}\hat{p}_s) + \hat{C}\hat{p}_s] \quad (26)$$

$${}^1\hat{D}^Z = e^{1\hat{H}^Z} \quad (27)$$

## 2.5 Decomposition of entropy

We decompose entropy and diversity following Marcon *et al.* [18]. Entropy and diversity can be measured in each community. When communities are pooled, a meta-community (whose species probabilities are the weighted average of those of communities) is defined. Entropy and diversity of the meta-community are called  $\gamma$ . The  $\alpha$  entropy of the meta-community is the weighted average entropy of communities, *i.e.* we follow Routledge's [46] definition of  $\alpha$  entropy. Its deformed exponential is  $\alpha$  diversity (which is thus not the average diversity of communities). Entropy is decomposed additively, according to Patil and Taillie's [8] concept of diversity of a mixture. Diversity is decomposed multiplicatively, and both decompositions are equivalent. Our purpose here is to characterize  $\beta$  entropy, beyond defining it only as the difference between  $\gamma$  and  $\alpha$  entropy.

The  $\alpha$  entropy of community  $i$  is:

$${}^q H_{\alpha}^Z = \frac{1 - \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1}}{q-1} \quad (28)$$

The  $\alpha$  entropy of the meta-community is the weighted average entropy of communities:

$${}^q H_{\alpha}^Z = \sum_i w_i {}^q H_{\alpha}^Z = \frac{1 - \sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1}}{q-1} \quad (29)$$

Algebra detailed in Appendix S2 gives the expression of  $\beta$  entropy:

$${}^q H_{\beta}^Z = \sum_i w_i \sum_s p_{si} \left( \ln_q \frac{1}{(\mathbf{Zp})_s} - \ln_q \frac{1}{(\mathbf{Zp})_{si}} \right) \quad (30)$$

$\beta$  entropy is the generalized Jensen-Shannon divergence [20] between the distribution of  $\mathbf{Zp}$  in each community and in the meta-community.

The  $\alpha$  diversity of the meta-community is the deformed exponential of its  $\alpha$  entropy:

$${}^q D_{\alpha}^Z = e_q^{q H_{\alpha}^Z} = \left[ \sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1} \right]^{\frac{1}{1-q}} \quad (31)$$

$\beta$  diversity can be obtained by taking the deformed exponential of the decomposition of entropy:

$${}^q D_{\beta}^Z = e_q^{\frac{{}^q H_{\beta}^Z}{1+(1-q){}^q H_{\alpha}^Z}} \quad (32)$$

## 2.6 Decomposition of diversity

It is also interesting to directly decompose diversity to characterize  $\beta$  diversity. It is usually only defined as the ratio between  $\gamma$  and  $\alpha$  diversity [47] or as a transformation of  $\beta$  entropy [18, 20, 48].

The inverse of  ${}^q D_{\gamma}^Z$  is the generalized mean of order  $q-1$  of  $(\mathbf{Zp})_s$ :

$$\frac{1}{{}^q D_{\gamma}^Z} = \left[ \sum_i p_s(\mathbf{Zp})_s^{q-1} \right]^{\frac{1}{q-1}} \quad (33)$$

From  ${}^q D_{\gamma} = {}^q D_{\alpha} {}^q D_{\beta}$  and Routledge weighting of alpha diversity, simple algebra (Appendix S3) yields:

$${}^q D_{\beta}^Z = \left[ \sum_i w_i \left( \frac{1/q {}^q D_{\alpha}^Z}{1/q {}^q D_{\gamma}^Z} \right)^{q-1} \right]^{\frac{1}{q-1}} \quad (34)$$

As  $1/q {}^q D_{\gamma}^Z$  is the average ordinariness of species,  ${}^q D_{\beta}^Z$  is the average normalized ordinariness of communities where the normalized ordinariness of community  $i$  is defined as  $\frac{1/q {}^q D_{\alpha}^Z}{1/q {}^q D_{\gamma}^Z}$ , that is to say its average species ordinariness divided by that of the meta-community.

We may consider each community as an assemblage of monospecific populations. Then, each population's diversity is  ${}^q D_{\alpha}^Z = 1$ , and weights are  $p_s$ . Introducing these values into eq. (34) yields  ${}^q D_{\gamma}^Z = {}^q D_{\beta}^Z$ : as shown by Rao and Nayak [49] for entropy, the diversity of an assemblage of species is the  $\beta$  diversity between its monospecific assemblages.

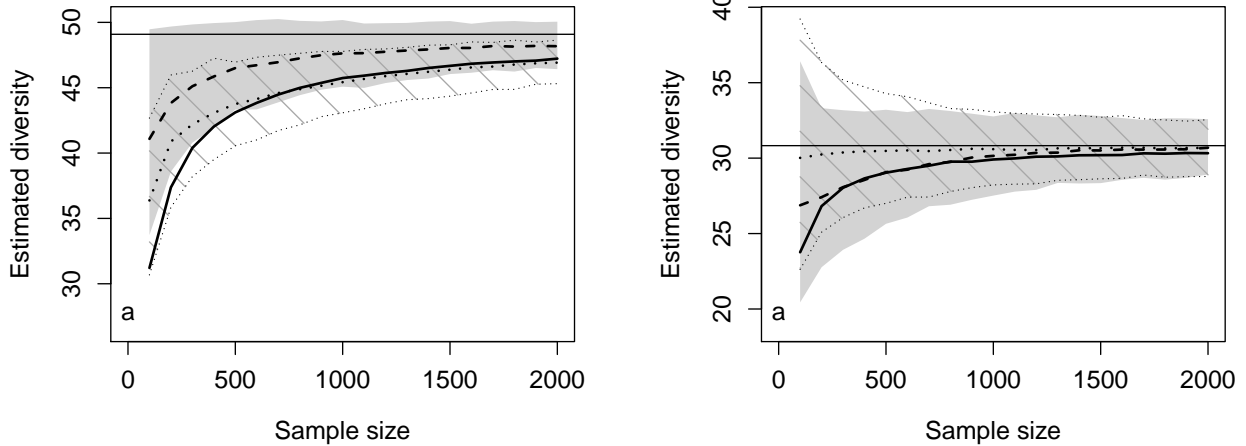
## 3. Results

### 3.1 Test of bias correction

We used the Barro-Colorado Island (BCI) 50-ha plot forest inventory [50–52]. Year 2005 census contains 20852 individuals from 229 species, among which 24 have been observed only once. The sample coverage is close to 99.9%, allowing to consider that the inventory is almost exhaustive and to use it as a reference to test the efficiency of estimators applied to subsamples. We set species similarity equal to 2/3 inside a genus, 1/3 inside a family, and 0 outside family. The average similarity between pairs of distinct species is  $\bar{z} \approx 0.01$ . The similarity matrix is somehow rough but no better data was available for such a number of species.

We simulated inventories between 100 and 2000 individuals in a multinomial distribution respecting same species





**Figure 1.** Estimation of BCI diversity  $^{0.5}D^Z$  (a) and  $^{1.5}D^Z$  (b) depending on sample size.

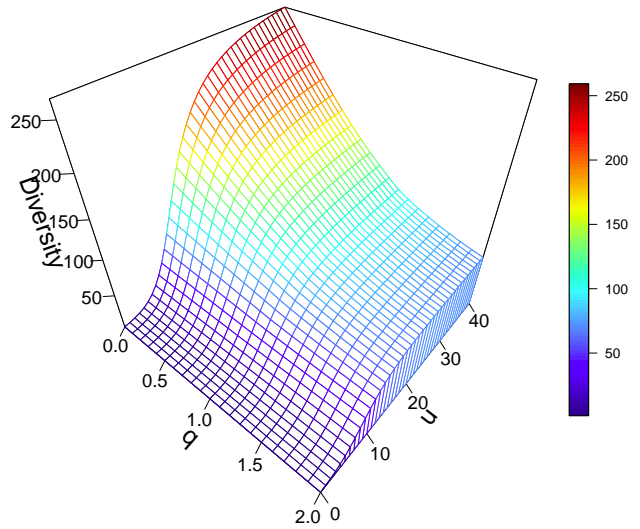
The plug-in estimator (plain line) is the least effective. The Chao-Shen estimator (dashed line, grey-shaded envelope) performs better than our new estimator (dotted line, dashed envelope) for small values of  $q$  but our new estimator is very effective for high values of  $q$ . Envelopes are the 2.5% and 97.5% quantiles of simulated values. The horizontal line is the diversity of the whole data set.

probabilities. Simulations were repeated 1000 times, estimated entropy averaged and finally transformed into diversity for plotting.

Simulated estimations are plotted in Figure 1. 95% confidence envelopes are shown: simulated samples are realizations of the assumed multinomial distribution of the community; stochasticity is not due to the estimators but to sampling. The best estimator depends on the value of  $q$ . Their correction is almost identical at  $q = 1.2$  (not shown). At  $q = 1.5$ , diversity can be estimated very accurately by the new estimator with a sample whose size is less than the number of species of the community. Following Marcon *et al.* [20]; we choose a pragmatic estimation-bias correction using the maximum value of the two estimators.

### 3.2 Partitioning functional diversity

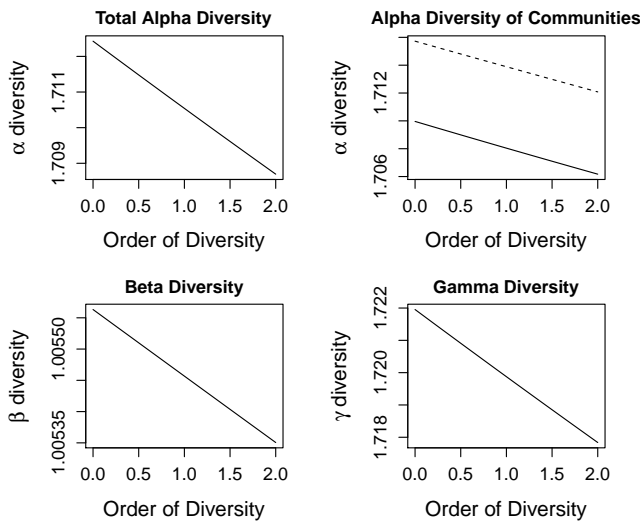
We used the same tropical forest dataset as Marcon and Hérault [18], made of two 1-ha fully inventoried plots in the Paracou field station in French Guiana. 1124 individual trees (diameter at breast height over 10 cm) have been sampled among 229 species. Four key functional traits were addressed: seed mass and tree maximum height [53], and specific leaf area and wood specific gravity [29]. A dissimilarity matrix was first built using the Gower metric by the *daisy* function of the cluster package [54] for R [55]. Our distance matrix is the square root of the output of the *daisy* function: it is euclidean [39]. The similarity matrix  $Z$  was defined as  $z_{s,t} = e^{-ud_{s,t}}$ . Figure 2 shows the diversity profile of the whole data set ( $\gamma$  diversity) depending on both parameters  $q$  and  $u$ . The effective number of species varies little whatever the order of diversity when  $u$  is small. When  $u$  increases, species are considered more



**Figure 2.** Bivariate profile of the Paracou meta-community. Parameter  $q$  controls the importance of original species, parameter  $u$  the similarity between species. High values of  $u$  make functional diversity converge to species-neutral diversity.

and more dissimilar and the diversity profile with respect to  $q$  converges to that of species-neutral diversity.

Among all possible values of  $u$ , we think that one is of particular interest: the one which maximizes the variance of similarity. We will discuss it in the next section. Considering our distance matrix, it is  $u = 2.48$ . We calculated diversity of order 1. Neutral  $\gamma$  diversity is 134 effective species. Similarity-based functional  $\gamma$  diversity is 3.68 effective species. It can be interpreted as: the average species ordinariness is  $1/3.68 = 0.27$ . It is also a Hill number: 3.68 completely different species with equal frequencies would have the same diversity as the observed meta-community made of 229 species. This is a very small value; due to the low functional distances between species in the similarity matrix:  $\bar{z} \approx 0.26$ .



**Figure 3.** Functional diversity profile of the Paracou forest communities.

Top left:  $\alpha$  diversity of the meta-community; top right:  $\alpha$  diversity of each community (P006: plain line, P018: dotted line); bottom left:  $\beta$  diversity; bottom right:  $\gamma$  diversity. Diversity is the number of effective species (effective communities for diversity) against the order  $q$ . It is estimated with bias correction.

Neutral  $\beta$  diversity is 1.46 effective communities: the species distributions are very different between the communities. Yet, functional  $\beta$  diversity is 1.02: almost no functional  $\beta$  diversity is detected between the communities. Figure 3 shows the functional diversity profile of the meta-community: its values are almost constant whatever  $q$ .

## 4. Discussion and Conclusion

In this paper, we derived the decomposition of similarity-based diversity  ${}^qD^Z$  and proposed two estimators to reduce estimation bias from the level of the plug-in estimator. The Chao-Shen estimator is the best for small orders of diversity but our alternative estimator outperforms it for higher orders.

Similarity-based diversity is preferred to estimate functional diversity from a distance matrix because it does not require building a dendrogram and so it preserves the topology of species in the space of functional traits. Even if we focused mainly on the functional diversity, the method works equally well with any similarity measures. As highlighted by Leinster and Cobbold, similarity can be measured in any meaningful way: a genetic similarity will lead to a decomposition of genetic diversity; a molecular similarity will lead to a decomposition of molecular diversity, and so on. Results are effective number of basic entities and entity assemblages which have the same desirable properties as classical diversity measures.

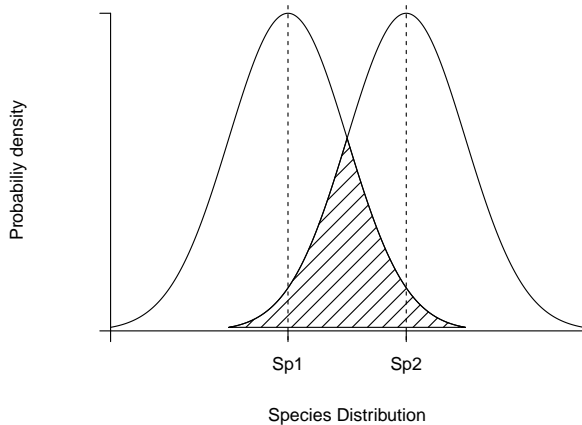
A critical issue is the definition of the scale parameter  $u$ . Chiu and Chao [36] question the legitimacy of  ${}^qD^Z$  to measure diversity in terms of effective number of species since many examples show that  ${}^qD^Z$  almost does not vary with  $q$ . The reason may be quite simple: all these examples define similarity as 1 minus the normalized distances, resulting in diversity profiles similar to that of Figure 2 for  $u = 1$ . Numerically, when  $u$  is small, similarities are higher, so are the values of species ordinariness  $(\mathbf{Zp})_s$ .  $(\mathbf{Zp})_s^{q-1}$  can be approached by its Taylor expansion of order 1 when it is close to 1. Simple algebra shows that the consequence of this linear approximation is that the generalized mean of order  $q - 1$  reduces to the arithmetic mean, allowing further simplifications in the calculation of  ${}^qD^Z$  which can be approximated by the inverse of the average similarity between species pairs:  $1/(\sum_s \sum_t p_s p_t z_{s,t})$ . Numerical simulations show that the approximation remains quite good for ordinariness values quite far from 1. In other words, if a significant proportion (in terms of sum of probabilities) of species have an ordinariness over say 0.4, this rough approximation appears to work very well on real examples. Another reason may be that diversity actually varies little because species ordinarinesses are very similar [56]. Then there are few original species and they do not contribute much to the value of diversity, whether they are given more importance by  $q$  or not.

This problem is partially solved when the appropriate definition of similarity is applied. Yet, another one emerges since no decisive criterion exists to choose a value of  $u$ , so a profile is required. The comparison of the diversity of two communities must be done along the bivariate profile allowing positive values of both  $q$  and  $u$ , as recommended by Leinster and Cobbold. The effective number of species given by  ${}^qD^Z$  for a specific couple of values  $q$  and  $u$  is perfectly defined but difficult to interpret since it varies substantially (from 1 to the species-neutral diversity) according to  $u$ . It is clearly useful to be able to choose a particular value of  $u$  to be able to synthesize the diversity of a community in a few effective numbers. We propose two solutions.

First, a classical approach in ordination methods consists of finding the best point of view to observe the highest possible variability of the data. We follow it here: considering the distance matrix and the constraints of its transformation into



a similarity matrix, there exists a value of  $u$  that maximizes the variance of  $\mathbf{Z}$ . It appears to be a reasonable choice because lower values make species excessively similar to each other, and greater values make them more different (up to the complete dissimilarity defining species-neutral diversity).



**Figure 4.** Representation of species overlap.

Species take place in the multidimensional space of traits. Intraspecific variability causes some uncertainty. Individuals are located around the expected position of each species, following a Gaussian distribution. Variance is assumed to be the same in all dimensions (whatever the trait) and for all species, for simplicity of the conceptual model. Two species are represented with their density of probability projected onto the axis that joins them. The standard deviation of the density of probability on the figure is half the distance between the two species.

The other possible solution consists of considering  $u$  as a proxy for species overlap. [57] represent the niche of species in a one-dimension space. It can be generalized to the multidimensional space of traits assuming variability is the same for all species and all traits (Figure 4). Similarity between species may be defined as the area of overlap of the distribution curves. In other words, a species' own niche can be defined as the region of the space of traits where its probability of occurrence is the greatest. Species overlap is the probability for a species to be in the other's niche. Mathematically, supposing both species have a normal distribution with the same variance, their overlap is:

$$O_{s,t} = 2 \int_{d_{s,t}/2}^{+\infty} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{d_{s,t}}{\sigma}\right)^2} \quad (35)$$

The overlap does not change as long as  $d_{s,t}/\sigma$  remains constant. As distances have been normalized, the variable of interest is  $\sigma$  which represents the intra-specific variability of traits. Let's consider the two most distant species, with

$d_{s,t} = 1$ . If  $\sigma = 1/6$ , their similarity is very small (around 1%). Higher values of  $\sigma$  mean that no species can be considered completely dissimilar, so  $z_{s,t} > 0$  whatever the species. Smaller intra-specific variability allows  $z_{s,t} \approx 0$  for closer couples of species.

A quite good approximation of  $O_{s,t}$  is  $e^{-\frac{d}{\sigma\sqrt{2}}}$ . Taking  $u = \frac{1}{\sigma\sqrt{2}}$ , the similarity  $z_{s,t} = e^{-ud_{s,t}}$  represents species overlap. The value of  $u$  corresponding to  $\sigma = 1/6$  is about 4. Higher scales correspond to less intra-specific variability.

We established the link between the mathematical properties of the scale parameter  $u$  and its biological meaning, even though the assumptions of the model, *i.e.* equal variability for all traits and all species, is not realistic. Yet, it is better than the usual hypothesis of no variability at all. The diversity of eigenvalues [58] allows taking into account intra-specific variability rigorously but does not allow estimation-bias correction.

The alternative approach by Chiu and Chao appears to be more straightforward since it directly relies on the distance matrix. Yet, it is influenced by the way the distance matrix is built. If it is based on a generalized Gower's metric, the distribution of distances can be modified by transforming the numerical trait variables (*e.g.* take their logarithm) or redefining the (arbitrary) numerical values of ordered categorical ones. The parameter  $u$  does nothing else, but it does it explicitly. The fundamental difference between the two definitions of functional diversity is not on the role of  $u$  but on that of  $q$ , which allows to focus on ordinary or rare species, depending of the chosen measure of diversity.

We estimated functional  $\beta$  diversity between two tropical forest plots previously investigated by Marcon and Hérault [18] with different methods. The optimal value of  $u$  (implying the similarity of the most different species is about 8 %) results in quite flat profiles of diversity (Figure 3). This is not a numerical artifact (the average species ordinarieness is 0.26) but due to low variability (its standard deviation is 0.03) despite the efforts made to maximize the variability of similarity. The functional  $\beta$  diversity between these two quite similar forest communities is negligible (very close to 1). The main explanation is probably the high functional redundancy between these two forest plots. They are spaced by only a few hundred meters and are located on similar topography and soils. It is always strange to observe how tropical forests can locally be extremely rich in species but exhibit a very low turnover in space [59, 60]. This finding, well-documented for neutral diversity, seems here exacerbated from the perspective of the functional diversity. This supports the view that, in tropical forests, niche differentiation, if any, occur at very fine spatial scale through vertical distribution shaped by light access and micro-habitat complementarity [61]. Above this micro-habitat scale, functional assemblages are very similar to each other and governed by a few leading and continuous functional axes [29]. This is supported by the flat shape of diversity profiles: giving more importance to original species (lowering  $q$ ) has no effect because less ordinary species are

not frequent. In other words, according to our limited data and number of traits, functional diversity is brought by frequent species and quickly saturated: adding more species or more communities does not increase it. These conclusion clearly have to be completed by the alternative point of view of Chiu and Chao's functional diversity which can evaluate the importance of rare species rather than original ones.

The entropart package [62] for R allows diversity estimation presented in the paper.

## 5. Acknowledgments

This work has benefited from an "Investissement d'Avenir" grant managed by Agence Nationale de la Recherche (CEBA, ref. ANR-10-LABX-25-01). Funding came from the project Climfor (Fondation pour la Recherche sur la Biodiversité).

## References

- [1] Gallopín GC (2006) Linkages between vulnerability, resilience, and adaptive capacity. *Global Environmental Change* 16: 293–303.
- [2] Violle C, Navas ML, Vile D, Kazakou E, Fortunel C, et al. (2007) Let the concept of trait be functional! *Oikos* 116: 882–892.
- [3] Weiner J (1995) On the practice of ecology. *Journal of Ecology* 83: 153–158.
- [4] Shannon CE (1948) A Mathematical Theory of Communication. *The Bell System Technical Journal* 27: 379–423, 623–656.
- [5] MacArthur R (1955) Fluctuations of Animal Populations and a Measure of Community Stability. *Ecology* 36: 533–536.
- [6] Hurlbert SH (1971) The Nonconcept of Species Diversity: A Critique and Alternative Parameters. *Ecology* 52: 577–586.
- [7] Rényi A (1961) On Measures of Entropy and Information. In: Neyman J, editor, 4th Berkeley Symposium on Mathematical Statistics and Probability. Berkeley, USA: University of California Press, volume 1, pp. 547–561.
- [8] Patil GP, Taillie C (1982) Diversity as a concept and its measurement. *Journal of the American Statistical Association* 77: 548–561.
- [9] Dalton H (1920) The measurement of the inequality of incomes. *The Economic Journal* 30: 348–361.
- [10] Havrda J, Charvát F (1967) Quantification method of classification processes. Concept of structural  $\alpha$ -entropy. *Kybernetika* 3: 30–35.
- [11] Daróczy Z (1970) Generalized information functions. *Information and Control* 16: 36–51.
- [12] Tsallis C (1988) Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics* 52: 479–487.
- [13] MacArthur RH (1965) Patterns of species diversity. *Biological Reviews* 40: 510–533.
- [14] Adelman MA (1969) Comment on the "H" Concentration Measure as a Numbers-Equivalent. *The Review of Economics and Statistics* 51: 99–101.
- [15] Hill MO (1973) Diversity and Evenness: A Unifying Notation and Its Consequences. *Ecology* 54: 427–432.
- [16] Jost L (2007) Partitioning diversity into independent alpha and beta components. *Ecology* 88: 2427–2439.
- [17] Jost L (2006) Entropy and diversity. *Oikos* 113: 363–375.
- [18] Marcon E, Hérault B (2015) Decomposing phylodiversity. *Methods in Ecology and Evolution* in press.
- [19] Chiu CH, Jost L, Chao A (2014) Phylogenetic beta diversity, similarity, and differentiation measures based on Hill numbers. *Ecological Monographs* 84: 21–44.
- [20] Marcon E, Scotti I, Hérault B, Rossi V, Lang G (2014) Generalization of the partitioning of Shannon diversity. *Plos One* 9: e90289.
- [21] Grassberger P (1988) Finite sample corrections to entropy and dimension estimates. *Physics Letters A* 128: 369–373.
- [22] Chao A, Shen TJ (2003) Nonparametric estimation of Shannon's index of diversity when there are unseen species in sample. *Environmental and Ecological Statistics* 10: 429–443.
- [23] Haegeman B, Hamelin J, Moriarty J, Neal P, Dushoff J, et al. (2013) Robust estimation of microbial diversity in theory and in practice. *The ISME journal* 7: 1092–101.
- [24] Zhang Z, Grabchak M (2014) Entropic Representation and Estimation of Diversity Indices. *arXiv* 1403.3031: 1–12.
- [25] Leinster T, Cobbold C (2012) Measuring diversity: the importance of species similarity. *Ecology* 93: 477–489.
- [26] Chao A, Chiu CH, Jost L (2014) Unifying Species Diversity, Phylogenetic Diversity, Functional Diversity and Related Similarity/Differentiation Measures Through Hill Numbers. *Annual Review of Ecology, Evolution, and Systematics* 45: 297–324.
- [27] Hardy GH, Littlewood JE, Pólya G (1952) *Inequalities*. Cambridge University Press.
- [28] Podani J, Schmera D (2006) On dendrogram-based measures of functional diversity. *Oikos* 115: 179–185.
- [29] Baraloto C, Paine CETP, Patiño S, Bonal D, Hérault B, et al. (2010) Functional trait variation and sampling strategies in species rich plant communities. *Functional Ecology* 24: 208–216.
- [30] McGill BJ, Enquist BJ, Weiher E, Westoby M (2006) Rebuilding community ecology from functional traits. *Trends in ecology & evolution* 21: 178–85.

- [31] Eviner VT, Chapin FSI (2003) Functional Matrix: A Conceptual Framework for Predicting Multiple Plant Effects on Ecosystem Processes. *Annual Review of Ecology, Evolution, and Systematics* 34: 455–485.
- [32] Wolda H (1981) Similarity indices, sample size and diversity. *Oecologia* 50: 296–302.
- [33] Beck J, Holloway JD, Schwanghart W (2013) Undersampling and the measurement of beta diversity. *Methods in Ecology and Evolution* 4: 370–382.
- [34] Butturi-Gomes D, Junior MP, Giacomini HC, Junior PDM (2014) Computer intensive methods for controlling bias in a generalized species diversity index. *Ecological Indicators* 37, Part A: 90–98.
- [35] Ricotta C, Szeidl L (2006) Towards a unifying approach to diversity measures: Bridging the gap between the shannon entropy and rao’s quadratic index. *Theoretical Population Biology* 70: 237–243.
- [36] Chiu CH, Chao A (2014) Distance-based functional diversity measures and their decomposition: a framework based on hill numbers. *PloS one* 9: e100014.
- [37] Horvitz DG, Thompson DJ (1952) A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* 47: 663–685.
- [38] Tsallis C (1994) What are the numbers that experiments provide? *Química Nova* 17: 468–471.
- [39] Pavoine S, Vallet J, Dufour AB, Gachet S, Daniel H (2009) On the challenge of treating various types of variables: application for improving the measurement of functional diversity. *Oikos* 118: 391–402.
- [40] Leinster T (2013) The Magnitude of Metric Spaces. *Documenta Mathematica* 18: 857–905.
- [41] Leinster T (2009) A maximum entropy theorem with applications to the measurement of biodiversity. *ArXiv* 0910.0906: 1–26.
- [42] Dauby G, Hardy OJ (2012) Sampled-based estimation of diversity *sensu stricto* by transforming Hurlbert diversities into effective number of species. *Ecography* 35: 661–672.
- [43] Bonachela JA, Hinrichsen H, Muñoz MA (2008) Entropy estimates of small data sets. *Journal of Physics A: Mathematical and Theoretical* 41: 1–9.
- [44] Good IJ (1953) On the Population Frequency of Species and the Estimation of Population Parameters. *Biometrika* 40: 237–264.
- [45] Zhang Z, Huang H (2007) Turing’s formula revisited. *Journal of Quantitative Linguistics* 14: 222–241.
- [46] Routledge RD (1979) Diversity indices: Which ones are admissible? *Journal of Theoretical Biology* 76: 503–515.
- [47] Tuomisto H (2010) A diversity of beta diversities: straightening up a concept gone awry. Part 1. Defining beta diversity as a function of alpha and gamma diversity. *Ecography* 33: 2–22.
- [48] Marcon E, Puech F, Traissac S (2012) Characterizing the relative spatial structure of point patterns. *International Journal of Ecology* 2012: 11.
- [49] Rao C, Nayak T (1985) Cross entropy, dissimilarity measures, and characterizations of quadratic entropy. *Information Theory, IEEE Transactions on* 31: 589–593.
- [50] Condit R (1998) *Tropical Forest Census Plots*. Berlin, Germany, and Georgetown, Texas: Springer-Verlag and R. G. Landes Company, 1–224 pp.
- [51] Hubbell SP (1999) Light-Gap Disturbances, Recruitment Limitation, and Tree Diversity in a Neotropical Forest. *Science* 283: 554–557.
- [52] Hubbell SP, Condit R, Foster RB (2005). *Barro Colorado Forest Census Plot Data*.
- [53] Hérault B, Bachelot B, Poorter L, Rossi V, Bongers F, et al. (2011) Functional traits shape ontogenetic growth trajectories of rain forest tree species. *Journal of Ecology* 99: 1431–1440.
- [54] Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K (2014). *cluster: Cluster Analysis Basics and Extensions*.
- [55] R Development Core Team (2014). *R: A Language and Environment for Statistical Computing*. URL <http://www.r-project.org>.
- [56] Messier J, McGill BJ, Lechowicz MJ (2010) How do traits vary across ecological scales? a case for trait-based ecology. *Ecology letters* 13: 838–48.
- [57] Lepš J, De Bello F, Lavorel S, Berman S (2006) Quantifying and interpreting functional diversity of natural communities: practical considerations matter. *Preslia* 78: 481–501.
- [58] Pavoine S, Izsák J (2014) New biodiversity measure that includes consistent interspecific and intraspecific components. *Methods in Ecology and Evolution* 5: 165–172.
- [59] Condit R, Pitman N, Leigh EG, Chave J, Terborgh J, et al. (2002) Beta-diversity in tropical forest trees. *Science* 295: 666–669.
- [60] Novotny V, Miller SE, Hulcr J, Drew RaI, Basset Y, et al. (2007) Low beta diversity of herbivorous insects in tropical forests. *Nature* 448: 692–5.
- [61] Kraft NJB, Valencia R, Ackerly DD (2008) Assembly in an Amazonian Forest Tree Community Functional Traits and Niche-Based. *Science* 322: 580–582.
- [62] Marcon E, Hérault B (2014) *entropart*, an R package to partition diversity. *Journal of Statistical Software* in press.

# Appendix S1: Estimation of $V$ and $W$

---

We first estimate  $V$ .

Let  $\Delta = [\bar{z}(1 - p_s) + p_s]^{q-1}$  and  $\Delta', \Delta'', \dots, \Delta^{(v)}$  be the  $v^{\text{th}}$  partial derivative of  $\Delta$  with respect to  $p_s$ . We have:

$$\begin{aligned}\Delta &= [\bar{z} + (1 - \bar{z})p_s]^{q-1} \\ \Delta' &= (q-1)(1 - \bar{z})[\bar{z} + (1 - \bar{z})p_s]^{q-2} \\ \Delta'' &= (q-1)(q-2)(1 - \bar{z})^2[\bar{z} + (1 - \bar{z})p_s]^{q-3} \\ &\dots \\ \Delta^{(v)} &= \left[ \prod_{i=1}^v (q-i) \right] (1 - \bar{z})^v [\bar{z} + (1 - \bar{z})p_s]^{q-v-1}\end{aligned}$$

At  $p_s = 1$ , we evaluate these derivatives to be:

$$\begin{aligned}\Delta' &= (q-1)(1 - \bar{z}) \\ \Delta'' &= (q-1)(q-2)(1 - \bar{z})^2 \\ &\dots \\ \Delta^{(v)} &= \left[ \prod_{i=1}^v (q-i) \right] (1 - \bar{z})^v = \left[ \prod_{i=1}^v (i-q) \right] (1 - \bar{z})^v (-1)^v\end{aligned}$$

Assuming  $\sum_{s \geq 1} p_s [\bar{z}(1 - p_s) + p_s]^{q-1} < \infty$  and  $\bar{z}$  is a small positive number, we can introduce the Taylor expansion of  $\Delta$  at  $p = 1$  in  $V$ . We have:

$$\begin{aligned}V &= \sum_{s \geq 1} p_s [\bar{z}(1 - p_s) + p_s]^{q-1} = \sum_{s \geq 1} p_s [\bar{z} + (1 - \bar{z})p_s]^{q-1} \\ &= \sum_{s \geq 1} p_s \left\{ 1 + \sum_{v=1}^{\infty} \frac{1}{v!} \left[ \prod_{i=1}^v (i-q) \right] (1 - \bar{z})^v (-1)^v (p_s - 1)^v \right\} \\ &= \sum_{s \geq 1} p_s \left\{ 1 + \sum_{v=1}^{\infty} \frac{1}{v!} \left[ \prod_{i=1}^v (i-q) \right] (1 - \bar{z})^v (1 - p_s)^v \right\} \\ &= 1 + \sum_{s \geq 1} p_s \left\{ 1 + \sum_{v=1}^{\infty} \frac{(1 - \bar{z})^v}{v!} \left[ \prod_{i=1}^v (i-q) \right] (1 - p_s)^v \right\}\end{aligned}$$

$$= 1 + \sum_{v=1}^{\infty} \frac{(1-\bar{z})^v}{v!} \left[ \prod_{i=1}^v (i-q) \right] \sum_{s \geq 1} p_s (1-p_s)^v$$

Denote  $\zeta_{1,v} = \sum_{s \geq 1} p_s (1-p_s)^v$ :

$$\begin{aligned} V &= \zeta_{1,0} + \sum_{v=1}^{\infty} \frac{(1-\bar{z})^v}{v!} \left[ \prod_{i=1}^v (i-q) \right] \zeta_{1,v} \\ &= \zeta_{1,0} + \sum_{v=1}^{\infty} (1-\bar{z})^v \left[ \prod_{i=1}^v \left( \frac{i-q}{i} \right) \right] \zeta_{1,v} \end{aligned}$$

$V$  is thus a linear diversity index. It can be estimated according to Zhang and Grabchak (2014) by:

$$\hat{V} = 1 + \sum_{s \leq T} \frac{n_s}{n} \sum_{v=1}^{n-n_s} (1-\bar{z})^v \left[ \prod_{i=1}^v \frac{i-q}{i} \right] \left[ \prod_{j=1}^v \left( 1 - \frac{n_s-1}{n-j} \right) \right]$$

$W$  is calculated the same way.

Let  $\Gamma = -\ln[\bar{z}(1-p_s) + p_s]$ . We have:

$$\begin{aligned} \Gamma &= -\ln[\bar{z} + (1-\bar{z})p_s] \\ \Gamma' &= -\frac{1-\bar{z}}{\bar{z} + (1-\bar{z})p_s} \\ \Gamma'' &= \frac{(1-\bar{z})^2}{[\bar{z} + (1-\bar{z})p_s]^2} \\ \Gamma^{(3)} &= \frac{-2(1-\bar{z})^3}{[\bar{z} + (1-\bar{z})p_s]^3} \\ &\dots \\ \Gamma^{(v)} &= \frac{(v-1)!(1-\bar{z})^v(-1)^v}{[\bar{z} + (1-\bar{z})p_s]^v} \end{aligned}$$

At  $p_s = 1$ :

$$\Gamma^{(v)} = (v-1)!(1-\bar{z})^v(-1)^v$$



We introduce the Taylor expansion of  $\Gamma$  at  $p = 1$  in  $W$ :

$$\begin{aligned}
W &= - \sum_{s \geq 1} p_s \ln[\bar{z} + (1 - \bar{z})p_s] \\
&= \sum_{s \geq 1} p_s \left\{ \sum_{v=1}^{\infty} \frac{1}{v!} (v-1)! (1 - \bar{z})^v (-1)^v (p_s - 1)^v \right\} \\
&= \sum_{s \geq 1} p_s \left\{ \sum_{v=1}^{\infty} \frac{(1 - \bar{z})^v}{v} (1 - p_s)^v \right\} \\
&= \sum_{v=1}^{\infty} \frac{(1 - \bar{z})^v}{v} \sum_{s \geq 1} p_s (1 - p_s)^v \\
&= \sum_{v=1}^{\infty} \frac{(1 - \bar{z})^v}{v} \zeta_{1,v}
\end{aligned}$$

Its estimator is:

$$\hat{W} = \sum_{s \leq T} \frac{n_s}{n} \sum_{v=1}^{n-n_s} \frac{(1 - \bar{z})^v}{v} \left[ \prod_{j=1}^v \left( 1 - \frac{n_s - 1}{n - j} \right) \right]$$

## Appendix S2: decomposition of entropy

---

We start from the multiplicative decomposition of diversity:

$${}^q D_\gamma^Z = {}^q D_\alpha^Z {}^q D_\beta^Z$$

We write the deformed logarithm of the equality:

$$\Leftrightarrow \ln_q {}^q D_\gamma^Z = \ln_q {}^q D_\alpha^Z + \ln_q {}^q D_\beta^Z - (q-1)(\ln_q {}^q D_\alpha^Z)(\ln_q {}^q D_\beta^Z)$$

We replace  $\ln_q {}^q D_\alpha^Z$  by  ${}^q H_\alpha^Z$  and we factorize the last two terms:

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = (\ln_q {}^q D_\beta^Z)[1 - (q-1) {}^q H_\alpha^Z]$$

We replace  ${}^q H_\alpha^Z$  by its value:

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = (\ln_q {}^q D_\beta^Z) \left[ 1 - (q-1) \frac{1 - \sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1}}{q-1} \right]$$

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = (\ln_q {}^q D_\beta^Z) \left( \sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1} \right)$$

We replace  $\ln_q {}^q D_\beta^Z$  by  $\ln_q ({}^q D_\gamma^Z / {}^q D_\alpha^Z)$  and introduce probabilities:

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = \frac{1 - \frac{\sum_s p_s(\mathbf{Zp})_s^{q-1}}{\sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1}}}{q-1} \left( \sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1} \right)$$

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = \frac{\sum_i w_i \sum_s p_{si}(\mathbf{Zp})_{si}^{q-1} - \sum_s p_s(\mathbf{Zp})_s^{q-1}}{q-1}$$

We replace  $p_s$  by  $\sum_i w_i p_{si}$  and factorize:

$$\Leftrightarrow {}^q H_\gamma^Z - {}^q H_\alpha^Z = \sum_i w_i p_{si} \frac{(\mathbf{Zp})_{si}^{q-1} - (\mathbf{Zp})_s^{q-1}}{q-1}$$

This result can be written:

$${}^q H_\beta^Z = \sum_i w_i \sum_s p_{si} \left( \ln_q \frac{1}{(\mathbf{Zp})_s} - \ln_q \frac{1}{(\mathbf{Zp})_{si}} \right)$$

## Appendix S3: decomposition of diversity

---

Recall the definitions of  $\gamma$  and  $\alpha$  diversity:

$${}^q D_\gamma^Z = \left( \sum_s p_s(\mathbf{z}\mathbf{p})_s^{q-1} \right)^{\frac{1}{1-q}}$$

$${}^q D_\alpha^Z = \left( \sum_i w_i \sum_s p_{si}(\mathbf{z}\mathbf{p})_{si}^{q-1} \right)^{\frac{1}{1-q}}$$

$\beta$  diversity is  $\gamma$  divided by  $\alpha$ :

$$\begin{aligned} {}^q D_\beta^Z &= \frac{{}^q D_\gamma^Z}{{}^q D_\alpha^Z} = \left( \frac{\sum_s p_s(\mathbf{z}\mathbf{p})_s^{q-1}}{\sum_i w_i \sum_s p_{si}(\mathbf{z}\mathbf{p})_{si}^{q-1}} \right)^{\frac{1}{1-q}} = \left( \sum_i w_i \frac{\sum_s p_{si}(\mathbf{z}\mathbf{p})_{si}^{q-1}}{\sum_s p_s(\mathbf{z}\mathbf{p})_s^{q-1}} \right)^{\frac{1}{q-1}} \\ &= \left( \sum_i w_i \left( \frac{{}^q D_\alpha^Z}{{}^q D_\gamma^Z} \right)^{1-q} \right)^{\frac{1}{q-1}} \end{aligned}$$

This can be rewritten:

$${}^q D_\beta^Z = \left( \sum_i w_i \left( \frac{1/{}^q D_\alpha^Z}{1/{}^q D_\gamma^Z} \right)^{q-1} \right)^{\frac{1}{q-1}}$$