

Consensus RDA across dissimilarity coefficients for canonical ordination of community composition data

F. GUILLAUME BLANCHET,^{1,2,3,5} PIERRE LEGENDRE,⁴ J. A. COLIN BERGERON,³ AND FANGLIANG HE³

¹*Department of Biology, Section of Ecology, University of Turku, Turku FIN-20014 Finland*

²*Mathematical Biology Group, Department of Biosciences, University of Helsinki, Helsinki FIN-00014 Finland*

³*Department of Renewable Resources, University of Alberta, 751 General Services Building, Edmonton, Alberta T6G 2H1 Canada*

⁴*Département de sciences biologiques, Université de Montréal, C.P. 6128, Succursale Centre-ville, Montréal, Québec H3C 3J7 Canada*

Abstract. Understanding how habitat structures species assemblages in a community is one of the main goals of community ecology. To relate community patterns to particular factors defining habitat conditions, ecologists often use canonical ordinations such as canonical redundancy analysis (RDA). It is a common practice to use dissimilarity coefficients to perform canonical ordinations through distance-based RDA (db-RDA) or transformation-based RDA (tb-RDA). Dissimilarity coefficients are measures of resemblance where the information about species communities is condensed into a symmetric square matrix of dissimilarities among sites. In this study, we compared 16 of the most commonly used dissimilarity coefficients to evaluate if the species-abundance distribution (SAD) of a community can be used to select an appropriate coefficient. Of these, 11 are designed to be used primarily with abundance data, although they can also be used with presence–absence data, whereas five can only be applied to presence–absence data. Using simulations, we compared the explained variance of RDAs differing only by their coefficients to evaluate how the abundance patterns of communities influence coefficient choice. We found that coefficients are largely equivalent, independently of the community SAD. In light of these findings, we propose the consensus RDA method, a new canonical ordination procedure that performs a consensus of RDAs across several coefficients. This new method focuses on the common relations found by independent RDAs differing only by their dissimilarity coefficients; this ensures the absence of a coefficient-related bias when interpreting the canonical ordination result. Also, because in our simulations the presence–absence data were directly derived from the abundance data, we were able to evaluate if the information in presence–absence data was equivalent to that in abundance data. We found that although some information was lost by converting abundance data into presence–absence, both data formats may be complementary. When applying consensus RDA to abundance and presence–absence data independently, a more complete understanding and interpretation of the ecological patterns is obtained. An ecological example illustrating consensus RDA and the conclusions of our simulations is presented, using Carabidae data collected at the Ecosystem Management Emulating Natural Disturbances (EMEND) project in northwestern Alberta, Canada.

Key words: abundance data; canonical redundancy analysis (RDA); Carabidae; dissimilarity coefficient; presence–absence data; species-abundance distribution (SAD); species-presence distribution (SPD).

INTRODUCTION

The species composition of an ecological community is heavily influenced by local variation in habitats. In theory, this intimate species–habitat relationship is due to evolutionary adaptations of the species to their environment; because of these adaptations, species have ecological niches (Hutchinson 1957), meaning that they are found at locations where they encounter appropriate living conditions. Whittaker (1967) illustrated this idea using the concept of environmental gradients, where different species use distinct sections of the same

gradient in a manner analogous to the dispersion of niches envisioned under Hutchinson's multivariate niche concept.

Numerous studies have shown that most communities that use a complex configuration of local habitats are composed of a few common species, plus a large proportion of less abundant or rare species. In contrast, species-poor communities with no dominant species are generally affected by only a few habitat gradients (Loreau 2010: Chapter 2). Thus, we may suggest that the complexity of species–habitat relationships influences the species-abundance structure of a community.

Variation in species abundance and the effects of multi-habitat gradients on this variation have been studied extensively. A common approach for depicting

Manuscript received 8 April 2013; revised 15 October 2013; accepted 5 November 2013. Corresponding Editor: H. H. Wagner.

⁵ E-mail: guillaume.blanchet@helsinki.fi

variation in species abundance is the species-abundance distribution (SAD), which ranks species in terms of the number of individuals of each species observed in sampling units from a community. SADs were mathematically described in the earlier ecological literature (Fisher et al. 1943 and Preston 1948). McGill et al. (2007) reviewed the various types of SADs and explained the utilities of SADs in describing and comparing communities.

At the community level, species-habitat relationships are often described using ordinations. Unconstrained ordinations such as principal component analysis (PCA; Pearson 1901), correspondence analysis (CA; Roux and Roux 1967), detrended CA (Hill and Gauch 1980), nonmetric multidimensional scaling (Shepard 1962), and principal coordinate analysis (PCoA; Gower 1966) have been widely used to study associations between species and habitat factors (Legendre and Legendre 2012: Chapter 9). More recently, constrained ordinations such as canonical redundancy analysis (RDA; Rao 1964) and canonical correspondence analysis (CCA; ter Braak 1986, 1987) have been used to more directly evaluate how specific habitat components affect species assemblages. It is well known that RDA is not well-suited to the analysis of species-abundance data collected along long gradients, which contain many zeros, because the Euclidean distance preserved in RDA does not have the property of being double-zero asymmetrical (Legendre and Legendre 2012: Subsection 7.2.2). Two variants of RDA have also been proposed to ecologists during the last 15 years: distance-based RDA (db-RDA; Legendre and Anderson 1999), which is a constrained version of PCoA, and transformation-based RDA (tb-RDA; Legendre and Gallagher 2001). Note that a PCA carried out on transformed data (tb-PCA) is the unconstrained version of tb-RDA. The transformations used in tb-PCA and tb-RDA make these ordination methods preserve one of the distances that is appropriate for the analysis of community composition data (Legendre and Legendre 2012: Sections 7.7, 9.1.10, and 11.1.5). In contrast with earlier approaches where the dissimilarity coefficient underlying the canonical ordination was fixed, db-RDA and tb-RDA make it possible to use an array of dissimilarity coefficients and data transformations to perform canonical ordinations, offering much more flexibility for the analysis of community data. A coefficient assesses the resemblance in species composition among sampled sites by condensing the community data into a symmetric square matrix of resemblance among sites. For example, the Euclidean distance (Table 1) computes Pythagoras' formula between all pairs of sites, which results in a symmetric square matrix where the species information is compared between two sites and condensed into a distance value.

Choosing a dissimilarity coefficient well-suited to study specific communities and particular ecological questions is a problem often faced by ecologists because of the overwhelming number of coefficients available in

the literature. As an example, Legendre and Legendre (2012: Chapter 7) describe 26 coefficients (distances and [dis]similarities) designed specifically for studying species assemblages. Although they propose theory-based guidelines and decision keys to help choose among coefficients (e.g., Legendre and Legendre 2012: Section 7.6), it often happens that more than one coefficient can be used to answer a particular ecological question. When such situations occur, Legendre and Gallagher (2001) suggest selecting the coefficient that yields the highest fraction of explained variance in canonical ordination; in other words, let the data determine which coefficient to use. Under this procedure, the abundance structure of a community can influence the selection of a coefficient used to describe it.

Although variation in SADs complicates coefficient selection, little is known about how variations in SADs affect the performance of coefficients. In this study, we compare the performance of dissimilarity coefficients commonly used in canonical-ordination and beta-diversity studies of community composition data and use simulations to evaluate the sensitivity of the coefficients to varying SADs. The comparisons are made for communities described either in terms of abundance or presence-absence data. The analysis meets two objectives. Firstly, by comparing the performance of coefficients within data type, we show that the choice of a coefficient based on the proportion of explained variance may influence the resulting interpretation of the species-habitat relationship. To solve this problem, we propose a new technique that computes a consensus among the canonical ordination results obtained from several coefficients. Secondly, by comparing coefficients between data types, we evaluate the extent to which information in abundance data is preserved after transformation to presence-absence data. We illustrate these effects using ground beetle (Carabidae) data from a boreal forest in northwestern Alberta, Canada.

DEFINING A COMMUNITY WITH A SAD

There are many ways to display a SAD. In this paper, we use a variation of Preston's (1948) graphs to describe species-abundance distributions where the abundance classes are arranged along the *x*-axis and increase according to a geometric progression, such that their lower bounds are 2^k , where *k* represents the successive integers from 0 and up. This approach was recommended by Gray et al. (2006) as the SAD construction that most accurately represents the species-abundance pattern of an ecological community. These graphs can be compared visually, making them effective tools to differentiate communities.

The 25 graphs shown in Fig. 1 present a range of possible SADs, most of which can be found in nature. All of them were employed to simulate site-by-species abundance matrices. For all SADs, the number of species was fixed at 20, but the total abundance varied from 261 (the sum of the abundance classes' lower limits

TABLE 1. List of dissimilarity coefficients compared in the study.

Dissimilarity coefficients	Equation	Reference	Comment
Binary symmetrical			
Simple-matching	$\sqrt{1 - \frac{a+d}{a+b+c+d}}$	Sokal and Michener (1958)	Directly related to Euclidean (see details in <i>RDA and dissimilarity coefficients</i>).
Binary probabilistic			
Raup-Crick	$1 - p(a_{hi})^\dagger$	Raup and Crick (1979) McCoy et al. (1986)	
Binary asymmetrical			
Jaccard	$\sqrt{1 - \frac{a}{a+b+c}}$	Jaccard (1901)	Binary equivalent of any variation of the modified Gower dissimilarity.
Sørensen	$\sqrt{1 - \frac{2a}{2a+b+c}}$	Sørensen (1948)	Binary equivalent of percentage difference.
Ochiai	$\sqrt{1 - \frac{a}{\sqrt{(a+b)(a+c)}}}$	Ochiai (1957)	
Distance between species profiles	$\sqrt{\frac{b+c}{(a+b)(a+c)}}$	Legendre and De Cáceres (2013)	Binary equivalent of the chord and Hellinger coefficient divided by $\sqrt{2}$.
Abundance symmetrical			
Euclidean	$\sqrt{\sum_{j=1}^p (y_{1j} - y_{2j})^2}$	Maor (2007)‡	Distance preserved in RDA.
Abundance asymmetrical			
Chord	$\sqrt{\sum_{j=1}^p \left(\frac{y_{1j}}{\sqrt{\sum_{j=1}^p y_{1j}^2}} - \frac{y_{2j}}{\sqrt{\sum_{j=1}^p y_{2j}^2}} \right)^2}$	Orlòci (1967)	On presence-absence data chord becomes $\sqrt{2(1 - \text{Ochiai})}$.
Hellinger	$\sqrt{\sum_{j=1}^p \left(\sqrt{\frac{y_{1j}}{y_{1+}}} - \sqrt{\frac{y_{2j}}{y_{2+}}} \right)^2}$	Rao (1995)	On presence-absence data Hellinger becomes $\sqrt{2(1 - \text{Ochiai})}$.
χ^2	$\sqrt{y_{++}} \sqrt{\sum_{j=1}^p \frac{1}{y_{+j}} \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2}$	Lebart and Fénelon (1971)	Dissimilarity preserved in CCA. Can also be used with presence-absence data.
Distance between species profiles	$\sqrt{\sum_{j=1}^p \left(\frac{y_{1j}}{y_{1+}} - \frac{y_{2j}}{y_{2+}} \right)^2}$	Legendre and Gallagher (2001)	Abundances standardized by division by row sums.
Percentage difference	$\frac{\sum_{j=1}^p y_{1j} - y_{2j} }{\sum_{j=1}^p (y_{1j} + y_{2j})}$	Odum (1950)	This dissimilarity is often wrongfully referred to as the Bray-Curtis index.§
$\sqrt{\text{Percentage difference}}$	$\frac{\sum_{j=1}^p \sqrt{y_{1j}} - \sqrt{y_{2j}} }{\sum_{j=1}^p (\sqrt{y_{1j}} + \sqrt{y_{2j}})}$	Clarke and Green (1988)	Taking the square root of the raw data prior to calculating percentage difference is often used when there is marked variation in abundance between species.

TABLE 1. Continued.

Dissimilarity coefficients	Equation	Reference	Comment
Abundance asymmetrical $\sqrt[p]{\text{Percentage difference}}$	$\frac{\sum_{j=1}^p \sqrt[p]{y_{1j}} - \sqrt[p]{y_{2j}} }{\sum_{j=1}^p (\sqrt[p]{y_{1j}} + \sqrt[p]{y_{2j}})}$	Clarke and Green (1988)	As in the previous row, but using a fourth root.
Modified Gower \log_2	$\frac{\sum_{j=1}^p w_j \lg_2(y_{1j}) - \lg_2(y_{2j}) }{\sum_{j=1}^p w_j}$	Anderson et al. (2006)	Different log bases are often used when there is marked variation in abundance between species. A high log base will generally reduce the emphasis of very abundant species more than a smaller one
Modified Gower \log_5	$\frac{\sum_{j=1}^p w_j \lg_5(y_{1j}) - \lg_5(y_{2j}) }{\sum_{j=1}^p w_j}$	Anderson et al. (2006)	See modified Gower \log_2 .
Modified Gower \log_{10}	$\frac{\sum_{j=1}^p w_j \lg_{10}(y_{1j}) - \lg_{10}(y_{2j}) }{\sum_{j=1}^p w_j}$	Anderson et al. (2006)	See modified Gower \log_2 .

Notes: All binary dissimilarities are presented in the form: $D = \sqrt{1 - S}$, where S is a similarity. The variables a , b , c , and d are defined in Table 2. y_{++} is the total sum of table \mathbf{Y} , y_{+j} is the abundance of species j , and y_{i+} is the sum of all abundance of site i . w_j is used to exclude double zeroes by setting $w_j = 0$ whenever $y_{1j} = 0$ and $w_j = 1$ elsewhere. All coefficients are presented in a dissimilarity (distance) format.

† The variables h and i define two different sites.

‡ The Euclidean distance was first defined in Mesopotamia (~1800 BC), see Maor (2007) for details.

§ Bray and Curtis (1957) did not design this coefficient nor was it their purpose. They used a transformed version of Steinhaus coefficient (Motyka 1947) in their paper, which is equivalent to Odum's (1950) percentage difference. Their transformed coefficient is actually Whittaker's index described in Legendre and DeCáceres (2013).

¶ The function $\lg(y_{ij}) = \log(y_{ij}) + 1$ when $y_{ij} > 0$, otherwise $\lg(y_{ij}) = 0$.

for all species of the community depicted by Fig. 1a) to 20 460 (the sum of the abundance classes' upper limits for all species of the community depicted in Fig. 1j). Therefore, the SADs of Fig. 1 represent a huge variation of species-abundance distributions, as would typically be observed in real communities (see Dewdney [2000] for a comparison of 50 SADs constructed from many different species communities). SADs were selected to represent a broad range of species-abundance patterns found in natural communities.

Fig. 1a, b present communities with many rare species and no common species. Note that communities with a similar SAD structure but with a larger number of rare species are often found in nature; however, because the SADs in Fig. 1 were used to define the abundance of species in simulated communities, the SADs in Fig. 1a, b are the most extreme cases that would not generate empty sites in the site-by-species table.

Ecologists sometimes remove species with low abundances because the many zeros introduced by including these rare species can be troublesome for data analysis, especially with methods based on Euclidean distances, as explained by Legendre and Legendre (2012: Subsection

7.4.1). For example, in the classical oribatid mite study of Borcard et al. (1992), 14 poorly represented species which, together, summed to 50 individuals, were removed from the data matrix before analysis by CCA. Depending on the group of organisms studied, removing rare species can yield SADs similar to what is found in Fig. 1c–g, m–o, u–v.

In a recent paper, Gaston (2010) emphasized the importance of studying common instead of rare species. In light of that work, we included a few SADs (Fig. 1h–j, w–y) corresponding to communities composed mainly of common species. Other SADs have been found to well characterize certain groups of organisms. For example, boreal carabid communities often present bimodal SADs (Niemelä 1993) such as those in Fig. 1k, l. Finally, the SADs presented in Fig. 1p–t, are mainly theoretical and unlikely to be found in nature. We included them because analysis of such extreme cases may lead to a better understanding of dissimilarity coefficients.

RDA AND DISSIMILARITY COEFFICIENTS

In this study, we used the RDA framework to compare commonly used dissimilarity coefficients (Table 1), all of

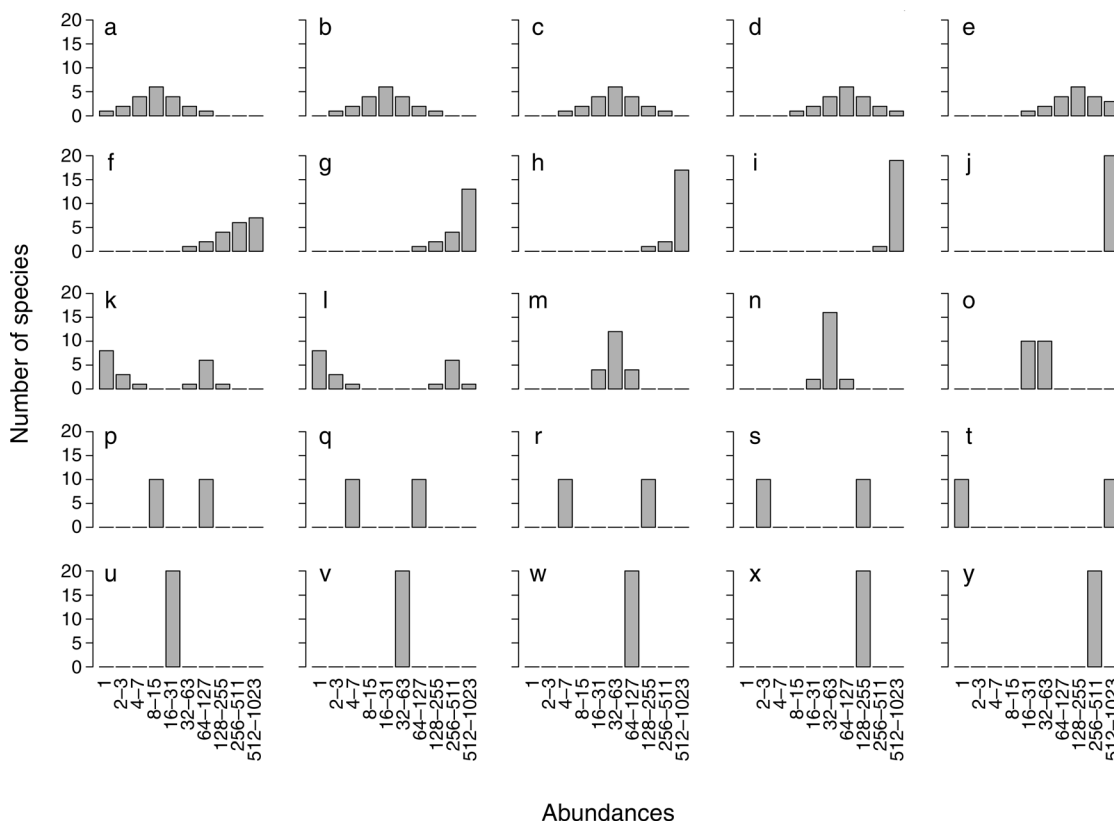


FIG. 1. Species-abundance distributions (SADs) used in the simulations. These SADs are presented using Preston (1948) graphs, where the abundance classes along the x-axis increase according to a geometric progression. The lower bound of the progression is made of the values 2^k , with k being the successive integers from 0 and up. The y-axis displays the number of species in each abundance class. These SADs were used as a basis for the simulations to generate a site-by-species data table. Each SAD represents a community of 20 species, and was constructed to encompass a wide range of variation in abundance patterns.

which can be used within db-RDA. Although most models were constructed through db-RDA, the chord, χ^2 , Hellinger, Ochiai, and distance-between-species-profiles coefficients were applied in tb-RDA because it is computationally more efficient. These five coefficients are mathematically equivalent in tb-RDA and db-RDA (Legendre and Legendre 2012: Section 7.7).

For presence-absence data, the Euclidean distance is equal to the square root of the complement of the simple-

matching coefficient (Table 1) multiplied by the number of species p : $\sqrt{p(1 - \text{simple-matching coefficient})}$; the formula reduces to $\sqrt{b + c}$ (see Table 2 for the meaning of b and c). This relationship was shown by Gower (1966) when he described PCA based on binary descriptors. A PCA based on binary data produces the same ordination as the principal coordinate analysis of a matrix of $\sqrt{1 - \text{simple-matching coefficient}}$; between the two ordinations, the coordinates are strictly proportion-

TABLE 2. Contingency table describing the similarity between two sites where species presence or absence were observed.

Site 1	Site 2	
	1 (species present)	0 (species absent)
1 (species present)	$a = \sum_{j=1}^p y_{1j}y_{2j}$	$b = \sum_{j=1}^p y_{1j} - \sum_{j=1}^p y_{1j}y_{2j}$
0 (species absent)	$c = \sum_{j=1}^p y_{2j} - \sum_{j=1}^p y_{1j}y_{2j}$	$d = p - a - b - c$

Notes: The variable a is the number of species present at sites 1 and 2, b is the number of species present at site 1 but absent at site 2, c is the number of species found at site 2 but not at site 1, and d is the number of species absent at both sites. The formulas for binary data in Table 1 describe how to combine the values a , b , c , and d to obtain the coefficients.

al and differ by a constant factor of \sqrt{p} . The same relationship holds when binary descriptors are used in an RDA, because it is the canonical extension of PCA. As a consequence, RDA based on binary data is equivalent to db-RDA of a matrix of $\sqrt{1 - \text{simple-matching coefficient}}$ and no data transformation is required.

By using the RDA framework for all coefficients, we were able to compare our simulation results directly. In particular, we used the χ^2 distance through the tb-RDA approach instead of calculating CCAs. In practice, tb-RDA with the χ^2 distance coefficient and CCA yield very similar, although not identical, ordination results (Legendre and Gallagher 2001).

An RDA is computed by regressing the community matrix \mathbf{Y} , composed of p species, on a matrix of m explanatory variables \mathbf{X} observed at the same n sites. This is carried out by a sum of squares minimization, leading to

$$\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

$$\hat{\mathbf{Y}} = \mathbf{XB} \quad (1)$$

where t indicates the transpose and -1 the inverse of a matrix. \mathbf{X} must either be centered by columns, or contain a column of 1's to estimate the regression intercepts. In Eq. 1, \mathbf{B} is the matrix of regression coefficients of all species in \mathbf{Y} on the explanatory variables \mathbf{X} . The residuals of the models are obtained through Eq. 2

$$\mathbf{Y}_{\text{res}} = \mathbf{Y} - \hat{\mathbf{Y}}. \quad (2)$$

By performing a PCA on $\hat{\mathbf{Y}}$, a matrix of eigenvectors \mathbf{U} defining the species scores and a diagonal matrix of eigenvalues $\mathbf{\Lambda}$ are obtained. The site scores can then be computed using \mathbf{X} (Eq. 3) or \mathbf{Y} (Eq. 4).

$$\mathbf{Z} = \mathbf{XB}\mathbf{U} = \hat{\mathbf{Y}}\mathbf{U} \quad (3)$$

$$\mathbf{F} = \mathbf{Y}\mathbf{U} \quad (4)$$

If required, the canonical coefficients can be calculated following Eq. 5:

$$\mathbf{C} = \mathbf{B}\mathbf{U}. \quad (5)$$

A more detailed description of the RDA algebra is available in Legendre and Legendre (2012: Section 11.1).

These calculations are exactly the same for tb-RDA, with the exception that the community matrix \mathbf{Y} is pre-transformed before calculating an RDA, using any of the transformations proposed by Legendre and Gallagher (2001). In db-RDA, a dissimilarity coefficient is applied to a community matrix, yielding a dissimilarity matrix. A PCoA is then calculated on this dissimilarity matrix and all the eigenvectors given by the PCoA are used as the \mathbf{Y} matrix in an RDA (Legendre and Anderson 1999). In db-RDA, the sites scores (Eqs. 3 and 4) and canonical coefficients (Eq. 5) are readily obtained. However, the species scores need to be

calculated a posteriori. We used the procedure proposed in the vegan package (Oksanen et al. 2013) to calculate the species scores

$$\mathbf{U}_{\text{post}} = \frac{\mathbf{Y}'\mathbf{Z}\mathbf{\Lambda}^{-1/2}}{\sqrt{n-1}}. \quad (6)$$

All binary similarity coefficients with the exception of the Raup-Crick coefficient were transformed into dissimilarities using $\sqrt{1 - \text{coefficient}}$ because Gower and Legendre (1986) have shown that this transformation makes them metric as well as Euclidean. This is important, because a PCoA of these transformed coefficients does not produce negative eigenvalues that would have to be corrected for before performing the RDA. Thus, this transformation facilitates the calculations. In contrast, the probabilistic nature of the Raup-Crick coefficient makes it special; on the one hand, its P value behaves like a dissimilarity, increasing as sites become more different in species composition; on the other hand, two sites with exactly the same species will not necessarily result in a dissimilarity of 0 for this coefficient; neither will two sites with completely different species automatically lead to a dissimilarity of 1. We decided to include it in our analyses because the probabilistic nature of the Raup-Crick coefficients may offer a solution to the double-zero problem.

The double-zero problem stems from the difficulty of relating two sites where a species has not been found (Legendre and Legendre 2012: Subsection 7.2.2). Double-zero asymmetrical dissimilarity coefficients are designed to ignore double zeroes altogether whereas double- x (where $x > 0$) reduces the dissimilarity; for binary dissimilarity coefficients, this amounts to ignoring the value d (Table 2) in the calculation of the coefficients. Conversely, double-zero symmetrical coefficients treat double zeroes as any other double- x value, which reduces the dissimilarity. For example, for the simple-matching coefficient, which is double-zero symmetrical, double zeroes (value d in Table 2) are considered as an indication of similarity in the same way as double presences (value a). Double-zero symmetrical coefficients should be used only when the goal of a study is to evaluate total changes in a community, for instance under the influence of pollution. Studies focusing on the impact of predation or disturbances may also find symmetrical coefficients interesting because the absence of a species at two sites is ecologically meaningful and should be considered (Anderson et al. 2011). For studies focusing on the variation in community composition among sites (i.e., beta diversity), double-zero asymmetrical coefficients should be preferred (Legendre and De Cáceres 2013).

In the present study, we performed simulations that reflected species variation in undisturbed communities where predation was not considered. The Euclidean and simple-matching coefficients are ill-adapted to these types of ecological problems because they are double-

zero symmetrical (Legendre and Legendre 2012: Subsection 7.4.1). We decided to include them when comparing coefficients within data types because both coefficients have been used, perhaps wrongfully, to study ecological communities through the use of RDA on abundance or presence-absence data (ter Braak and Verdonschot 1995).

The Jaccard, Sørensen, and simple-matching coefficients were computed with the *ade4* package (Dray and Dufour 2007). All other calculations were performed with the *vegan* package (Oksanen et al. 2013) with the exception of the Raup-Crick coefficient, which was programmed independently using the McCoy et al. (1986) permutation procedure. We used the McCoy et al. (1986) permutation approach because Legendre and Legendre (2012: Subsection 7.3.5) found that it was better at recognizing significant site associations, compared to the original permutation procedure of Raup and Crick (1979). All analyses were carried out using the R statistical language (R Development Core Team 2012).

SIMULATING COMMUNITIES WITH VARYING SPECIES ABUNDANCES

In our simulations, we constructed eight explanatory variables at 49 sites structured as a regular grid comprising 7×7 sites, using the *RsimSSDCOMPAS* package (M.-H. Ouellette, *personal communication*) within the R statistical language. The *RsimSSDCOMPAS* package is a wrapper for *SimSSD4*, a FORTRAN program used to simulate species, environment, and geographic coordinates. *SimSSD4* is available in ESA's Ecological Archives M075-017-S1, a supplement to the Legendre et al. (2005) paper. These explanatory variables (matrix \mathbf{X}) define linear gradients, waves, large patches, or random patterns. They are presented in Fig. A1 of Appendix A with a detailed description of how they were constructed. The same eight descriptors were used for all simulations.

In a simulated community, each of the 20 species had a different underlying structure constructed by combining pairs of the eight explanatory variables presented above. This structure remained constant for all simulated communities. The reference structure \mathbf{y}_{ref} of a species was constructed following Eq. 7, where ω is a weight, \mathbf{x}_i and \mathbf{x}_j are two of the eight explanatory variables, and $\boldsymbol{\varepsilon}$ is an error vector of standard normal deviates

$$\mathbf{y}_{\text{ref}} = \omega(\mathbf{x}_i + \mathbf{x}_j) + \boldsymbol{\varepsilon}. \quad (7)$$

The weight ω acts as a regression coefficient to influence the abundance of each species in the community, which is directly related to the size of the absolute value of ω (i.e., $|\omega|$). A value of ω was predefined for each species. A large $|\omega|$ generates species with larger abundances. Half of the species were constructed with positive weights and the other half with negative weights.

Ten species were characterized by strong links ($\omega = 2$ or -2) with the explanatory variables defining them. In

ecological terms, a large absolute weight represents a species that has a strong relationship with the measured environmental variables. The other 10 simulated species had smaller weights representing medium (two species with $\omega = 1$ or -1), weak (four species with $\omega = 0.5$ or -0.5), or very weak (four species with $\omega = 0.1$ or -0.1) relationships between a species and the descriptors controlling it.

As will be explained at the end of this section, additional sets of communities were simulated where the error $\boldsymbol{\varepsilon}$ was smaller, giving more importance to species with lower absolute weights. Note that Eq. 7 without the error term $\boldsymbol{\varepsilon}$ represents the true pattern defining a species. The reference structure of each species was determined following a predefined combination of ω , \mathbf{x}_i , and \mathbf{x}_j (Appendix A: Table A1). Also, the explanatory variables used to construct each species were carefully selected in such a way that each one was independently used to create five different species, making all explanatory variables equally important in the simulated community.

To construct a species, we transformed \mathbf{y}_{ref} for it to range from 0 to 1 in order to use the information it encompasses as a probability distribution. Eq. 8 was used if ω was positive and Eq. 9 if ω was negative. In these two equations, $|\mathbf{y}_{\text{ref}}|$ is the absolute value of \mathbf{y}_{ref} and \mathbf{y}_{prob} defines the probabilities of sampling a species at each of the 49 sites in the sampling area

$$\mathbf{y}_{\text{prob}} = \frac{|\mathbf{y}_{\text{ref}}|}{\sum |\mathbf{y}_{\text{ref}}|} \quad (8)$$

$$\mathbf{y}_{\text{prob}} = \frac{1}{|\mathbf{y}_{\text{ref}}|} \times \frac{1}{\sum |\mathbf{y}_{\text{ref}}|}. \quad (9)$$

Eq. 8 defines the probability of sampling a species directly related to the patterns in \mathbf{y}_{ref} , whereas Eq. 9 defines the probability of sampling a species inversely related to the patterns in \mathbf{y}_{ref} . If the probability of sampling a species is high for a site in proportion to the other sites, it is more likely for at least one individual of that species to be found at the site.

As explained in *Defining a community with a SAD*, the abundance pattern of each simulated community (defined as a group of species living in heterogeneous environment) followed one of the predefined SADs presented in Fig. 1. The SAD is a commonly used tool to rank species, based on the abundance of each species sampled from a community (McGill et al. 2007). The structures of these SADs were unaffected by the other steps of the simulation; SADs remained constant throughout the simulations. Each species was assigned to a bin of the SAD in order for the abundance distribution of the community to be reproduced when summing the number of individuals for each species in the site-by-species table. To define the exact abundance of a species in a simulated community, we randomly sampled the number of individuals of that species within its SAD bin boundary. Each number of individuals had

the same probability of being selected within a particular bin boundary. To allocate these individuals to specific sites, we sampled the sites (with replacement), using the species probability distribution y_{prob} . Because y_{prob} is constructed from y_{ref} , which has an underlying normal error, y_{prob} also follows a normal distribution.

By repeating this procedure for the 20 species, we obtained a site-by-species table representing one simulated community. We constructed 1000 communities for each of the 25 SADs in Fig. 1. Four other sets of 25 000 communities were also constructed where the error terms ϵ in Eq. 7 were standard normal deviates with standard deviations of 0.001, 0.250, 0.500, and 2.000. In all, we simulated 125 000 communities describing the abundance of species at each site.

To create site-by-species presence-absence tables, we transformed all abundances larger than 0 to 1s for all species-abundance community data generated above.

COMPARING DISSIMILARITY COEFFICIENTS USING EXPLAINED VARIANCE

The amount of explained variance in canonical ordinations was estimated with the coefficient of determination (R^2) and the comparison of dissimilarity coefficients was carried out following the procedure proposed by Legendre and Gallagher (2001). Coefficients of determination were calculated by dividing the total variance in $\hat{\mathbf{Y}}$ (which, incidentally, is also the sum of the canonical eigenvalues) by the total variance in \mathbf{Y} (which is also the sum of all eigenvalues, canonical and non-canonical). Note that $\hat{\mathbf{Y}}$ was constructed following Eq. 1, where \mathbf{X} is a matrix of explanatory variables, each of which are shown in Appendix A: Fig. A1. R^2 values range from 0 to 1. For example, if a model yields an R^2 of 0.2, it should be understood that this model explains 20% of the variance of the response.

In the present study, only the canonical eigenvalues associated with significant canonical axes ($P \leq 0.05$ after 999 random permutations) were considered in the calculation of R^2 . Fig. 2 compares the performance of RDAs for different dissimilarity coefficients for each of the 25 SADs presented in Fig. 1. The RDAs were carried out on the simulated species abundances constructed with the smallest error (normal distribution with a standard deviation of 0.001). Results of simulations with larger error are presented in Appendix B: Figs. B1–B4. All simulations yielded the same conclusions (see next paragraph) regardless of the error size. The only difference between the sets of simulations is that larger error when constructing species is associated with lower R^2 . The inverse relation between error term and variance explained, which is consistent for all coefficients compared, suggests that the amount of error does not favor (or disfavor) any coefficient. Note that if all canonical eigenvalues are used to calculate the R^2 instead of using only the significant eigenvalues, the conclusions are unchanged, because the fractions of the explained variance corresponding to the nonsignificant

canonical axes are too small to markedly affect the results. The variance explained by all the nonsignificant canonical axes considered together is above 0.1 only in extreme cases, and is usually around 0.06. The variance explained by a single nonsignificant canonical axis is usually less than 0.025.

In the simulation results presented in Fig. 2, the most striking feature is that the confidence intervals for all double-zero asymmetrical coefficients overlap considerably. Moreover, detailed inspection of the results shows that independent of the SAD structures, a community having a high R^2 for one coefficient generally also has high R^2 for other coefficients.

The R^2 values for the Euclidean distance differ most from the other coefficients, although its confidence intervals still overlap with the other coefficients (Fig. 2, top panel). This is because the Euclidean distance is a double-zero symmetrical coefficient. For the same reason, the confidence intervals are much wider for the Euclidean distance than for any other coefficient. At sites with the same environmental conditions, one should expect to find the same species, but species abundances usually vary. Although these variations in abundance may have important implications when species are rare, they should have only negligible effects on the results when species are common. In that instance, the Euclidean distance considers common and rare species similarly. The results associated with the Euclidean distance suggest that double-zero symmetrical coefficients should only be used to address ecological questions where double zeroes are ecologically meaningful, as suggested by Anderson et al. (2011).

Ecologists should also be careful in using the distance between species profiles, especially in the presence of many common species, because it seems to lose explanatory power in these circumstances (Fig. 1h–j, w–y). Legendre and De Cáceres (2013) have shown that the distance between species profiles lacks some of the important properties necessary for coefficients that are used to assess beta diversity. For this reason, the distance between species profiles suffers from the same problem as the Euclidean distance in the presence of common species, but to a lesser extent.

When comparing dissimilarity coefficients with simulated presence-absence data, the R^2 coefficients are very similar between coefficients across the different SADs (Fig. 3). Results for the Raup-Crick coefficient were the only exception, although its confidence intervals still overlap importantly with the others. It yields a somewhat lower R^2 when there are many common species (Fig. 3, central panel). Because a high R^2 for the Raup-Crick coefficient is generally associated with a high R^2 of the other coefficients, it may be that the Raup-Crick coefficient does not as effectively capture patterns as the other coefficients when many common species are sampled (Fig. 1h–j, y). These results are consistent with Legendre and Legendre (2012: Subsection 7.3.5), who showed that the statistical power of the

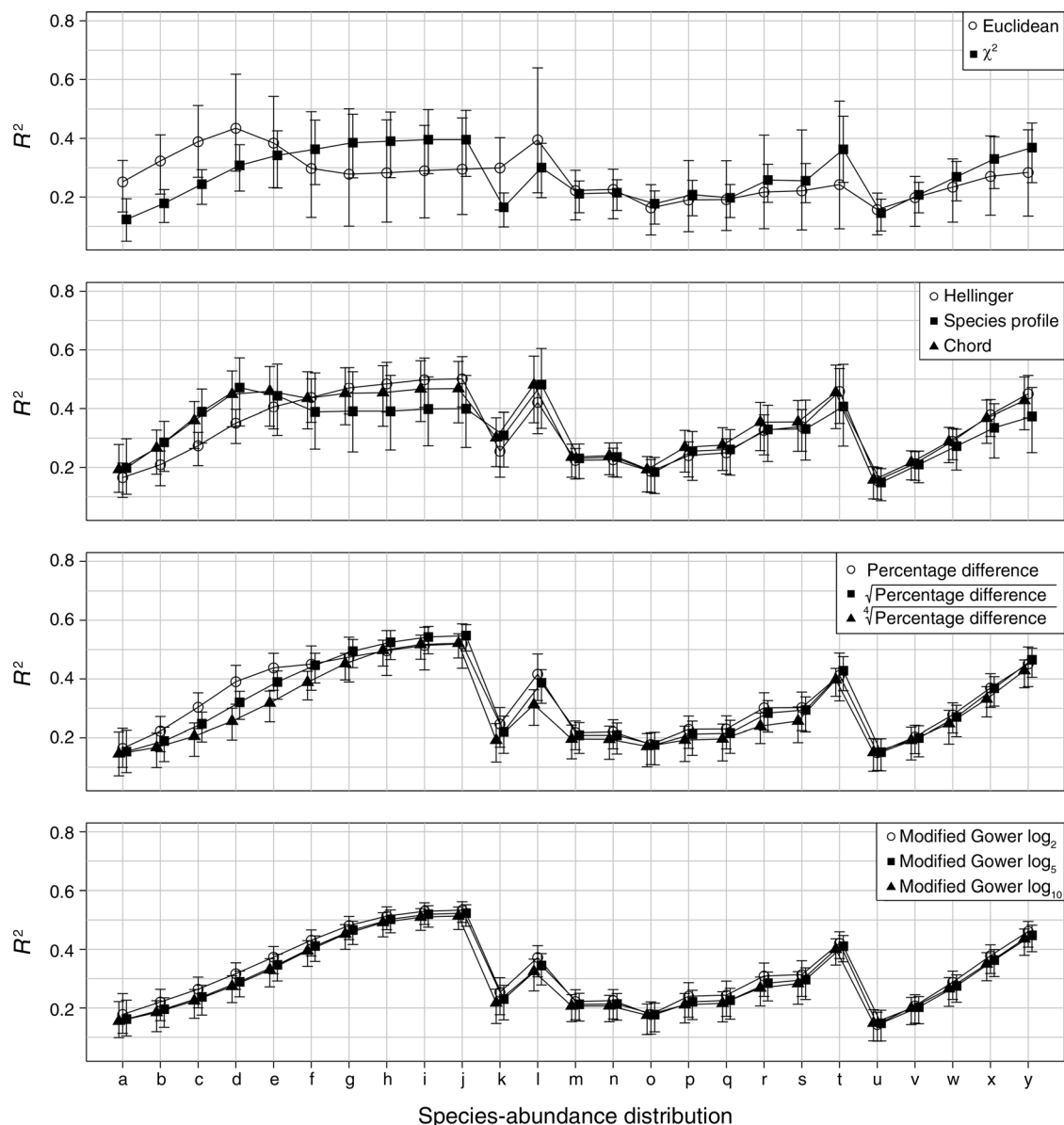


FIG. 2. Comparison of explained variance (R^2) for 11 dissimilarity coefficients calculated from simulated communities following different SADs using abundance data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations, and error bars represent 95% confidence intervals. Coefficients are presented in different panels strictly for visual clarity; scale remains consistent throughout. SADs under comparison are the same as those presented in Fig. 1, and share the same letter identification. A line was drawn along the R^2 results of each coefficient to facilitate comparisons between coefficients. Results are based on species simulated with an error term sampled from a normal distribution (mean = 0, standard deviation = 0.001). A thousand simulations were run for each SAD.

Raup-Crick coefficient to detect significant association between pairs of sites is low even when McCoy et al.'s (1986) permutation procedure is used.

We were surprised that the simple-matching coefficient produced results equivalent to other coefficients (Fig. 3, central panel). We expected it to be burdened by the same problems as the Euclidean distance because the simple-matching coefficient is the presence-absence equivalent of the Euclidean distance, making it a

double-zero symmetrical coefficient. However, it seems that when abundances are considered, the importance of double zeroes increases. If a single species is sampled in large abundances at two sites, the Euclidean distance between these sites for that particular species is likely to be somewhat far from 0, even though it is clear that these sites are quite similar. The Euclidean distance thus overemphasizes the differences between two sites where a species is found in large but unequal abundances, a

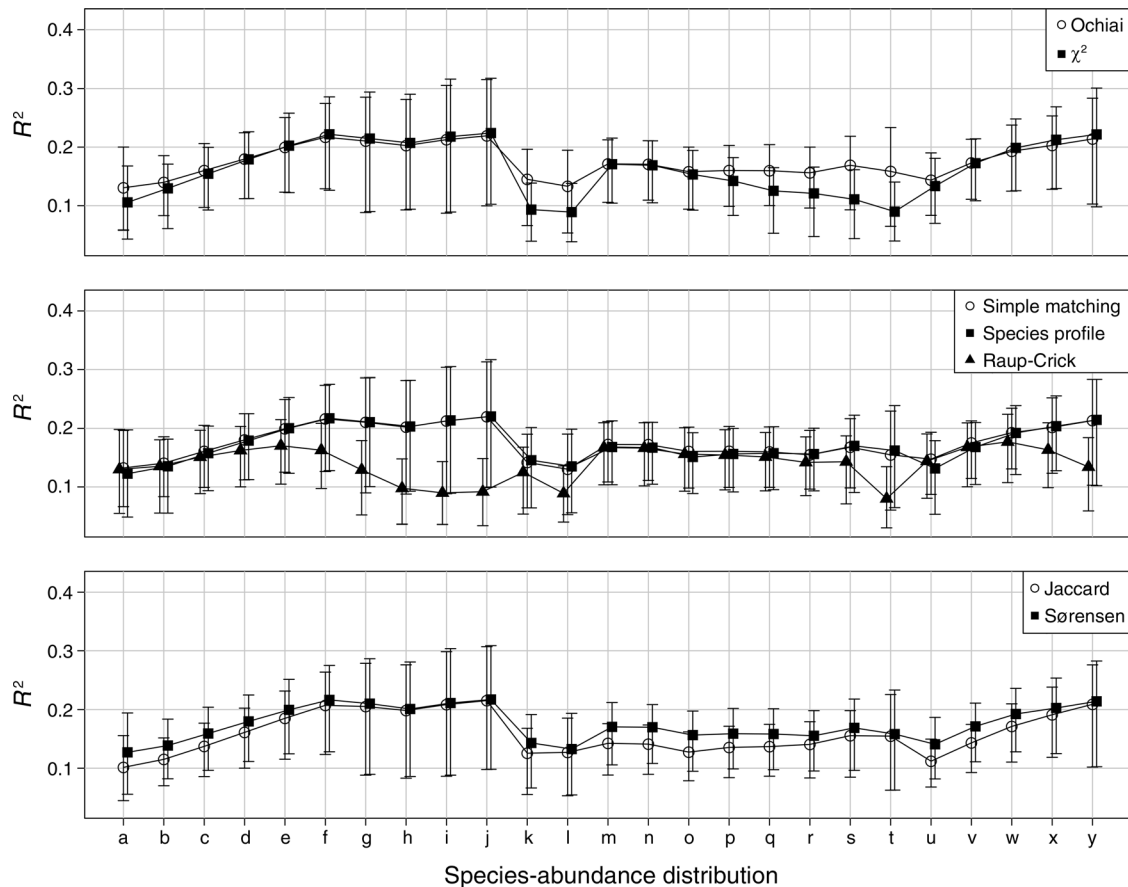


FIG. 3. Comparison of explained variance (R^2) for seven dissimilarity coefficients calculated from simulated communities following different SADs using presence-absence data. Details are as in Fig. 2.

problem that does not exist for the simple-matching coefficient because the species will be recorded as present (or 1) for both sites, yielding a distance of exactly 0.

Another aspect of our simulations is the increase in explained variance with the number of common species (progression of R^2 from SAD a to j in Fig. 1). This trend is consistent for all coefficients compared (with the exception of the Euclidean, species profile, and Raup-Crick coefficients, discussed earlier in this section), in abundance and presence-absence data alike, although it is weaker for presence-absence data (Appendix B: Figs. B5–B8). Similar conclusions were found with communities simulated with larger error (Appendix B: Figs. B1–B8).

A NEW WAY TO PERFORM CANONICAL ORDINATIONS

The previous simulations have shown that within data types, double-zero asymmetrical coefficients yield similar values of R^2 , for each SAD compared (Figs. 2 and 3 and Appendix B: Figs. B1–B8). This is shown by the substantial overlap between confidence intervals of all coefficients calculated for any particular SAD. Each coefficient has particularities making it more appropriate for specific ecological situations or research questions, and less so for others. With the wealth of

coefficients available in the ecological literature, it is common for more than one coefficient to be appropriate for a particular ecological study. In that context, the question “which dissimilarity coefficient should be used?” remains incompletely answered.

Here, we propose a three-step procedure to handle this problem. Even though most of the information highlighted by the different coefficients is often quite similar, the mathematical properties of each coefficient emphasize certain characteristics in the data that other coefficients do not, and vice versa.

In that respect, the first step is to compare coefficients and evaluate how much the information they explain diverges. This is accomplished by comparing all aspects of the canonical ordination models (i.e., the sites, the species, and the canonical coefficients), not only the variance explained. Secondly, a selection may be carried out among the coefficients, if necessary. The RDA models constructed using coefficients that differ markedly from the others should be considered separately, or their use should be reevaluated. The differences between RDA models can be in the ordination of the sites, the site–species relationships, and/or the relationships between canonical coefficients and the sites and species. In

a nutshell, potential differences among RDA models should be sought in all aspects of the models. Comparison and selection of coefficients is recommended, because if a coefficient is markedly different from the others, its inclusion in the following consensus step may blur ecological relationships that could appear if this coefficient was removed. Thirdly, the information common to RDA models that only differ by their dissimilarity coefficients should be synthesized. It is important to focus only on the information shared by the different RDA models to ensure that no misguided ecological interpretations are made. Because it is difficult to extract common information by an examination of independent canonical ordination triplots, we propose a new method that computes a consensus among canonical ordinations that differ only by the coefficients used to construct them. The consensus focuses on the patterns found by all RDA models, leaving out the information extracted by only one or a few coefficients. We call this new approach “consensus RDA.” A detailed explanation of how these three steps are carried out is presented in the next sections.

Comparison of RDA models.—To compare RDA models that only differ by their coefficients, the first step is to isolate the significant components found in each \mathbf{Z} matrix (site scores calculated using the explanatory variables, Eq. 3), e.g., the axes with a P value ≤ 0.05 . Model comparisons rely on the \mathbf{Z} matrices, which contain the canonical ordination coordinates of the sites; the variance of each canonical axis in \mathbf{Z} is equal to its associated eigenvalue when the distances among sites are preserved in the ordination results (RDA scaling 1). In the RDA framework, the canonical eigenvalues measure the variance explained by the canonical axes.

We correlated the significant canonical axes of the \mathbf{Z} matrix obtained for each dissimilarity index to those obtained with the other indices using RV coefficients (Escoufier 1973, Robert and Escoufier 1976). The RV coefficient is a multivariate generalization of the squared Pearson correlation that correlates two matrices with corresponding rows (sites). It produces values that range between 0 (no correlation) and 1 (perfect correlation). The RV coefficients for all pairs of dissimilarity indices were assembled in a matrix of pairwise RV coefficients. Using this matrix, we drew a minimum spanning tree (MST; Legendre and Legendre 2012: Section 8.2) to compare dissimilarity indices. This required the matrix of RV coefficients to be transformed into a dissimilarity matrix. We used $(1 - \text{RV})$ to perform the transformation because it ensured that the correlation information brought by the RV coefficients was conserved. These dissimilarities ranged from 0 to 1.

Selection of RDA models.—After examination of the MST, a selection of concordant dissimilarity indices can be made. We leave it at the discretion of users to decide how dissimilarity indices should be selected. For example, the dissimilarity indices linked by the longest MST branches can be removed if these branches are

much longer than the average branch. If the longest branch in the MST links two groups of dissimilarity indices, it may be interesting to calculate two consensus RDAs, one for each group of indices.

Consensus RDA.—To calculate a consensus RDA, the significant components of the \mathbf{Z} matrices selected to compare RDA models are used again (Fig. 4b). Of course, only the \mathbf{Z} matrices from coefficients that have been selected in the previous step should be considered. In consensus RDA, all significant components are grouped in a large matrix (Fig. 4c). Using this large matrix as a response in an RDA where the matrix of explanatory variables is \mathbf{X} (Fig. 4c), compute the consensus RDA site scores \mathbf{Z}^* and the consensus RDA canonical coefficients \mathbf{C}^* . This RDA also yields eigenvalues (Λ^*), which express the amount of variance represented by each \mathbf{Z}^* component, and more generally by each axis of the consensus RDA. These eigenvalues can be used to measure the strength of the consensus. The consensus RDA species scores \mathbf{U}^* need to be calculated following Eq. 6. In other words, \mathbf{U}^* is obtained following the same procedure as in db-RDA.

When performing an RDA, the results can be presented in either a distance (scaling 1) or a correlation (scaling 2) triplot. Scaling can also be used in consensus RDA. All the calculations presented above are carried out using the scaling 1 matrices \mathbf{Z} because, as explained in *Comparison of RDA models*, the consensus method relies on a property of \mathbf{Z} that is only present in scaling 1. To obtain a consensus result in scaling 2, the consensus site scores (matrices \mathbf{Z}^*) need to be rescaled following $\mathbf{Z}^*(\Lambda^*)^{-1/2}$. A similar procedure is used to apply scaling 2 on the species scores consensus ($\mathbf{U}^*(\Lambda^*)^{-1/2}$).

An interesting aspect of this new method is that as long as the dissimilarity indices represent the only aspect that differs between the different RDAs, a consensus RDA can be computed. This also includes partial RDAs.

All calculations necessary to obtain a consensus RDA rely on the \mathbf{Z} matrices of the different RDA models constructed with a group of relevant dissimilarity indices. The components in these \mathbf{Z} matrices contain the fitted site scores for the RDAs; they do not include the residuals components of \mathbf{Y} , which are part of the \mathbf{F} matrices (ter Braak 1994, Legendre and Legendre 2012: Subsection 11.1.3). Because it is often more interesting to study an RDA triplot in a projection where residuals are not included, a consensus of \mathbf{F} matrices was not incorporated in our description of consensus RDA.

The explanations of how to perform consensus RDA indicate that any number of axes can be used for any of the RDAs that are considered in the calculation of the consensus. However, it is not clear whether all, or only the significant canonical axes, should be used in a consensus RDA to obtain the model that best explains the community data. To evaluate which approach should be used, the simulated site-by-species tables presented in *Simulating communities with varying species abundances* were used. Each site-by-species table was

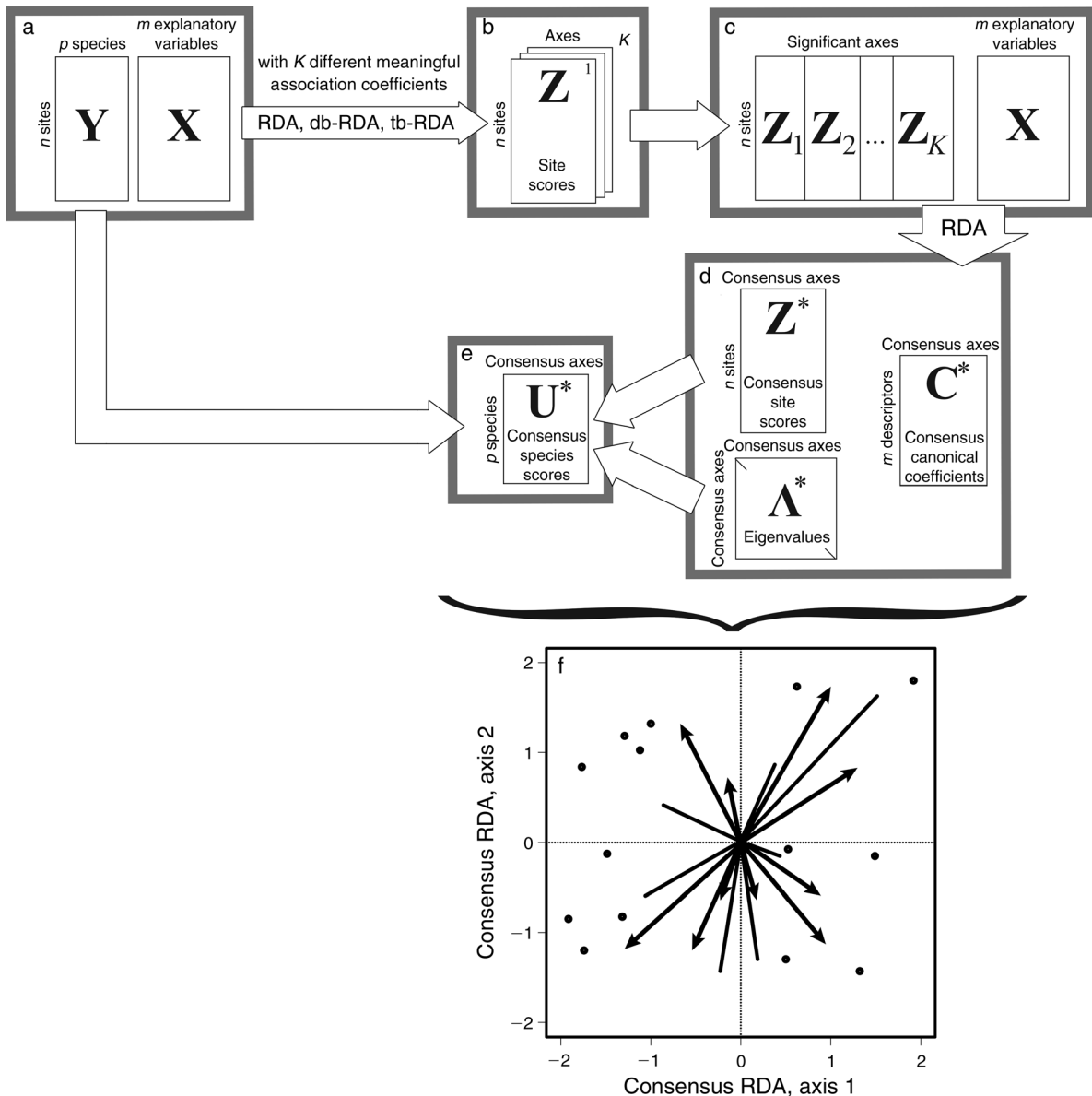


FIG. 4. Schematic representation of consensus redundancy analysis (RDA). (a) The first step of the procedure is to perform a series of RDAs (transformation-based RDA [tb-RDA] or distance-based RDA [db-RDA]) to model the community data Y using explanatory variables X . Each RDA is computed with a different dissimilarity coefficient using scaling type 1 (distance triplot, Z matrices). In the figure, K different dissimilarity coefficients are used. (b) For each of the K dissimilarity coefficients, the significant axes within each Z matrix are grouped in a large matrix. (c) An RDA is then performed on this large matrix using X as the explanatory variables. (d) This RDA yields the site scores consensus matrix Z^* , a diagonal matrix of eigenvalues Λ^* , and the consensus canonical coefficients C^* . (e) Eq. 6 is then used to obtain the consensus species scores U^* . (f) Z^* , U^* , and C^* can be used to draw a consensus RDA triplot; the eigenvalues in Λ^* show the importance of each axis in the consensus triplot.

correlated with Z^* (consensus site scores), which was calculated using all canonical axes, using the RV coefficient. We then compared these RV coefficients with RV coefficients correlating the site-by-species tables with the consensus site scores calculated using only the significant axes. The comparisons were carried out using both abundance and presence-absence simulated data. All dissimilarity coefficients discussed were used in the construction of the consensus site scores.

The results in Fig. 5 were obtained using abundance data where the error was the largest (ϵ in Eq. 7 followed a normal distribution with mean = 0 and standard deviation = 2), which yielded the largest variation in the comparisons made. In Fig. 5 (note the narrow range of the y -axis), the differences between the RV coefficients calculated using all canonical axes and the RV coefficients computed using only significant axes ranges almost always between 0.05 and -0.02 . Although, for

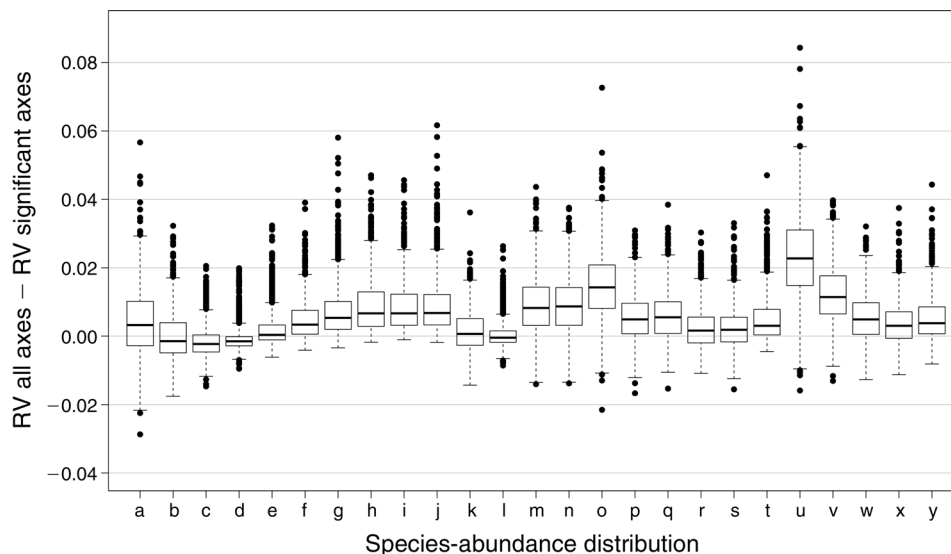


FIG. 5. Comparison of consensus RDAs across SADs (same as in Fig. 1), constructed using all canonical axes with consensus RDAs using only the significant canonical axes. The Z^* matrices calculated from the abundance data were used in the comparison. The y-axis presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes; all difference values were in the interval $[-0.03, 0.09]$. The RV coefficient is a multivariate generalization of the squared Pearson correlation that correlates two matrices with corresponding rows (sites). The results are presented as boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box is the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a normal distribution (mean = 0, standard deviation = 2). A thousand simulations were run for each SAD.

certain extreme cases, slightly more information can be obtained using all canonical axes, in the majority of situations very little information is gained (or sometimes lost) from using all canonical axes instead of only the significant ones. Results from the simulations where communities were generated with smaller error terms are presented in Appendix C. In these simulations, presence-absence and abundance data were considered. For abundance data, the results yield the same conclusions as the one presented in Fig. 5. For presence-absence data, it is slightly better to use all canonical axes; however, the information gain is minimal. In doubtful cases, the best solution is found by comparing a consensus RDA obtained using all canonical axes with a consensus RDA constructed with only the significant axes and choosing the solution that yields the largest RV coefficient. This approach ensures that the result of the consensus RDA always represents the largest amount of information from the community data.

A comparison of dissimilarity indices and a consensus RDA are presented in *Ecological illustration*, for abundance and presence-absence data.

SHOULD WE USE PRESENCE-ABSENCE DATA?

Modeling presence-absence data is more challenging than abundance data because information on species abundances is missing. The results of our simulations confirm this statement; the R^2 values are consistently higher for abundance (Fig. 2, Appendix B: Figs. B1–B4)

than for presence-absence data (Fig. 3, Appendix B: Figs. B5–B8). This result is not surprising, because one would expect to obtain better species-environment linear models when using more informative data. This finding remains the same irrespective of the error level in the data (Fig. 2, Appendix B: Figs. B1–B8). However, comparison between presence-absence and abundance data using R^2 does not reflect how well the true species structure is modeled. To compare the canonical ordination results of abundance and presence-absence data, we first need to measure how much information from the true species (Eq. 7 without the error term) structure is extracted by the canonical analyses. As explained at the end of *RDA and dissimilarity coefficients*, the Euclidean and simple-matching coefficients are double-zero symmetrical; they are designed to answer ecological questions where double zeroes are ecologically meaningful. In our simulations, double zeroes do not necessarily reflect a strong similarity between sites. For this reason, double-zero symmetrical coefficients were not included in the comparison between abundance and presence-absence data. For both data types, we calculated RV coefficients between the true species structure (Eq. 7 without the error term) and the significant canonical axes.

We regrouped all RV coefficient results within data type and compared the grouped abundance to the grouped presence-absence results (Fig. 6). According to the results obtained by comparing dissimilarity coefficients within data type (Figs. 2, 3, Appendix B: Figs.

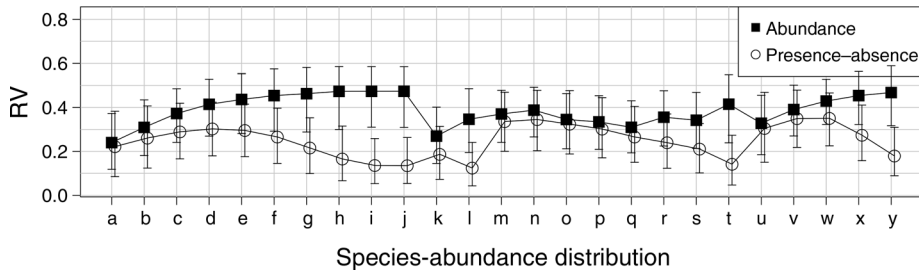


FIG. 6. Comparison between abundance and presence-absence data by SADs (as in Fig. 1), showing how much of the true species structure (Eq. 7 without the error term) is modeled by the canonical ordination models. For each data type (abundance and presence-absence), the significant canonical axes computed using all dissimilarity coefficients (excluding the double-zero symmetrical coefficients) were grouped. RV coefficients were then used to correlate the true species structure with the grouped significant canonical axes. Error bars represent 95% confidence intervals. A line was drawn along the R^2 results of each dissimilarity coefficient to facilitate comparison between the two data types. Results are based on species simulated with an error term sampled from a normal distribution (mean = 0, standard deviation = 0.001). A thousand simulations were run for each SAD.

B1–B8), it is valid to group dissimilarity coefficients used on the same data type because no dissimilarity coefficient dominates over the others for any SAD. Fig. 6 illustrates the grouped results for simulations where the error is the smallest (standard deviation = 0.001). What is striking about these results is that when there are many common species (Fig. 1i, j, y), the amount of information extracted by canonical ordinations is much less for presence-absence than for abundance data. These conclusions can be extended to situations where there are at least as many common as there are rare species (Fig. 6, SADs g, h, l, and t) because the overlap between confidence intervals is small in these situations. This suggests that for communities with at least as many common as rare species, the information lost in canonical ordinations on occurrences should not be interpreted in the same way as results obtained from canonical ordinations on abundance data. Similar results were obtained for data simulated with larger errors (Appendix D: Figs. D1–D4). We will show in *Ecological illustration* how these findings apply to real ecological data.

ECOLOGICAL ILLUSTRATION: CARABIDAE OF NORTHWESTERN ALBERTA

To show how the previous findings may be applied in real ecological situations, we extend the analysis to a data set about ground beetles (Carabidae) sampled at 192 sites in a never-harvested mature boreal mixedwood forest (see Bergeron et al. 2011, Blanchet et al. 2013). In this illustration, we aim at finding how trees influence the ground beetle community in the boreal forest. This question has already been approached with the same data by Bergeron et al. (2011). The difference here is that we used consensus RDA based on several dissimilarity indices detailed later in this section, for abundance and presence-absence data. Bergeron et al. (2011) performed all their analyses using a single dissimilarity calculated on abundance data.

The sites, which formed a near-regular grid in an area of 70 km², were located in the Ecosystem Management

Emulating Natural Disturbances (EMEND) experimental area in northwestern Alberta, Canada. Each site contained three pitfall traps (Spence and Niemelä 1994) located on the perimeter of a 15 m radius circle. From the center of the circle, a trap points due north while the other two are separated by 120 degrees. The community data are composed of 37 ground beetle species sampled throughout the summer of 2003. Beetle abundances were divided by the number of days each trap was active to remove the effect of trap disturbance and of non-demonic intrusions (Hurlbert 1984). Presence-absence data for each site were obtained by transforming all abundances larger than 0 to 1.

As explanatory variables, the relative basal areas of the 25 trees closest to the center of each site were used. Eight tree species were present in the experimental area and the relative basal area of each species was used as an explanatory variable. Further analysis of this data set may be found in Blanchet et al. (2013) and in Bergeron et al. (2011, 2012). The Hellinger distance was used by Blanchet et al. (2013) and Bergeron et al. (2012), and the percentage difference distance was employed by Bergeron et al. (2011). Note that Bergeron et al. (2011) used nonmetric multidimensional scaling (Legendre and Legendre 2012: Section 9.4) to study carabids, unlike Blanchet et al. (2013) and Bergeron et al. (2012), who used tb-RDAs.

In this ecological illustration, we compare canonical ordinations calculated on abundance and presence-absence data, considering results from all dissimilarity coefficients used in our simulations, with the exception of the double-zero symmetrical coefficients. We did not use double-zero symmetrical coefficients because they consider double zeroes (the absence of a species at two sites) as informative, which may lead to wrongful interpretations. The carabid data set used in this illustration was sampled to study how habitat variation influenced the ground beetle community. Blanchet et al. (2013) have shown that this community is mostly unaffected by anthropogenic disturbances. In this context, Anderson et al. (2011) explained that double zeroes are not necessarily

ecologically meaningful, making the use of double-zero symmetrical coefficients inappropriate for studying this particular carabid community.

An RV comparison of the RDA models constructed with different dissimilarity coefficients is presented using MSTs in Fig. 7b for abundance data and Fig. 7e for presence-absence data. Each MST was constructed from a dissimilarity matrix of RV coefficients correlating all pairs of RDA models obtained from the different coefficients following the procedure presented in *A new way to perform canonical ordinations*. As a reference, we included in Table 3 the amount of variance explained (R^2) by the full db-RDA or tb-RDA models based upon the different dissimilarity coefficients. We used the full RDA models because the final consensus RDA results were more informative than when only the significant axes were used. This was true for abundance and presence-absence data.

We found that for both abundance and presence-absence data, the RDA model constructed using the χ^2 distance was the most different from the others (Fig. 7b, e). This is likely due to *Notiophilus directus*, a rare species found with low abundance at three sites where *Pterostichus punctatissimus* and *Miscodera arctica* (one site), or only *Pterostichus punctatissimus* (two sites), were encountered. Legendre and Legendre (2012: Subsection 7.4.1) and Greenacre (2013) explained that the χ^2 distance gives higher weights to species represented by only a few individuals at sites where only a few other species are found. Legendre and De Cáceres (2013) also found that the χ^2 distance lacked an important property for analysis of community composition data. Because we did not want to give undue importance to rare species, we did not further consider the χ^2 distance in the analyses of this carabid community.

Using the remaining coefficients, we constructed a consensus RDA. We plotted as many species as we could in the consensus RDA triplots without losing overall interpretability. The species not presented on the diagrams were consistently near the center of the consensus triplots, which made it impossible to interpret the ecological relationships of these species with respect to the tree basal areas. The first two axes of the consensus RDA represent 88.4% of the variance for abundance data and 85.2% for presence-absence data, and thus represent well the information in the different RDA models. The third consensus axis explained less than 7% of the variance, for abundance as well as presence-absence data, making the information presented by any subsequent axes too small to justify their use. Note that the R^2 in consensus RDA represents the strength of the consensus, i.e., how much of the variation is joint (or consensual), not the strength of the individual model, as it is the case for traditional RDA.

Although the amount of information explained by the first two axes of the consensus RDAs based on abundance (Fig. 7c) and presence-absence data (Fig. 7f) is similar, the underlying information is different.

For example, segregation of beetle species niches between coniferous (*Abies balsamea* [Ab], *Larix laricina* [Ll], *Picea mariana* [Pm], and *Picea glauca* [Pg]) and deciduous forest (*Betula papyrifera* [Bp], *Populus tremuloides* [Pt], and *Populus balsamifera* [Pb]) on the positive side of the x-axis is better achieved using the consensus RDA based on abundance data. Because these beetle and tree species are all characteristic of upland mixedwood forest (Bergeron et al. 2011), the comparison between abundance (Fig. 7c) and presence-absence (Fig. 7f) consensus triplots suggests that beetle species occur all along the deciduous-coniferous forest gradient, but it is their abundance that varies according to habitat. Also, relationships between beetle and tree species were not always consistent between the two consensus ordinations. For example, *Calathus advena* (Calaadve) is more closely related to *P. glauca* in the abundance ordination (Fig. 7f) than it is in the presence-absence ordination (Fig. 7c).

The results about *Agonum gratiosum* (Agongrat), *Carabus chamissonis* (Caracham), *Platynus mannerheimii* (Platmann), and *Pterostichus brevicornis* (Pterbrev) are impossible to interpret in the consensus RDA computed using species abundance because they are too close to the triplot center. However, in the presence-absence consensus triplot, these species are interpretable. This may relate to the fact that presence-absence data give more weight to less abundant species than relative abundance data (Anderson et al. 2011, 2006). In this beetle assemblage, *P. mannerheimii* is of special ecological interest even if it is not a common species, because it has a narrow habitat requirement that is locally restricted to old wet and productive forest (Bergeron et al. 2011). This species, as well as *A. gratiosum*, is only found at sites dominated by *P. mariana* and *L. laricina*, even if their abundance at these sites is not as high as that of the more common species. In that instance, it makes sense that the consensus RDA on presence-absence data makes these two species stand out. Such results show that performing community analyses on both abundance and presence-absence data concurrently makes it possible to extract interesting information out of the data.

The SAD (Fig. 7a) of this beetle community depicts many rare and many common species, which is typical for carabid communities (Niemelä 1993). The species-presence distribution that describes species occurrence for these data highlights more sharply the two groups of species in the carabid data (Fig. 7d). Our simulations suggest that a community composed of many rare and many abundant species (Fig. 11, t) does not preserve well community patterns after having been transformed into presence-absence (Fig. 6, SADs l and t). This is reflected in the ecological analysis, where abundance-based ordination achieves a better segregation of the beetle ecological niches. Although this may suggest that presence-absence ordinations are not useful on their own, differences between the abundance and the

presence-absence ordinations may have an ecological foundation. It may be that differences between ordinations based on abundance and presence-absence data reflect the spatial aggregation of carabid species. It is also possible that the consensus RDA calculated on abundance data brings complementary information to the consensus RDA results obtained from presence-absence data. To know if the differences between the two consensus RDAs are determined by ecological processes, a detailed study of this carabid community needs to be carried out contrasting presence-absence and abundance data at multiple scales using other variables characterizing the habitat of Carabidae, in addition to the tree basal areas.

In this carabid example, consensus RDA gives strong confidence in the ecological associations discovered between ground beetles and trees. Here are a few examples of discoveries made by Bergeron et al. (2012) that hold true in consensus RDA (Fig. 7c). *Agonum retractum* (Agonretr) and *Platynus decentis* (Platdece) prefer forest containing *P. balsamifera*, *P. tremuloides*, and to a lesser extent *B. papyrifera*, which are all upland deciduous trees. *Stereocerus haematopus* (Sterhaem) and *Calathus advena* (Calaadve) were commonly found in coniferous forest dominated by *P. glauca* and *A. balsamea*. *Calathus ingratus* (Calaingr) and *Pterostichus adstrictus* (Pteadst) typically occurred in both deciduous and coniferous upland forest where *P. mariana* and *L. laricina* are usually absent. *Pterostichus punctatissimus* (Pterpunc) and *P. mariana* present a strong ecological association. Because consensus RDA is based on many dissimilarity coefficients, the ecological associations discovered between trees and ground beetle species of the mixedwood boreal forest should be ecologically more meaningful and reliable, and not biased by the dissimilarity coefficient used in the calculation of the consensus ordination.

It is not the goal of this study to present a detailed ecological study of northwestern Alberta boreal carabids. However, by comparing the consensus RDA calculated on the carabid abundance data (Fig. 7c) with the canonical ordination results from Bergeron et al. (2012, Fig. 4) who also studied the relationship between carabids and tree relative basal area with the same data using RDA, differences can be found that are solely attributed to the dissimilarity coefficient used. For example, in our results, *S. haematopus* is more closely related to *P. glauca* than it is in Bergeron et al. (2012). These authors based the canonical ordination on the Hellinger transformation of the beetle abundance data, which emphasizes the composition nature of the data rather than raw abundance (Anderson et al. 2011). The consensus RDA of Fig. 7c, which uses a variety of coefficients along the composition-abundance gradient, indicates that the abundance pattern of *S. haematopus* is more closely associated with *P. glauca* than previously discovered by Bergeron et al. (2012). To prevent a biased interpretation resulting from the use of a specific

TABLE 3. Variance explained (R^2) by RDA models constructed independently with each dissimilarity coefficient using data from the ecological illustration, where the tree relative basal area was used to model a ground beetle (Carabidae) assemblage.

Dissimilarity coefficient	R^2
Abundance data	
Species profiles	0.303
Chord	0.321
Hellinger	0.340
χ^2	0.094
Percentage difference	0.203
$\sqrt{\text{Percentage difference}}$	0.238
$\sqrt{\text{Percentage difference}}$	0.249
Modified Gower log ₂	0.297
Modified Gower log ₅	0.304
Modified Gower log ₁₀	0.302
Presence-absence	
Species profiles	0.225
Ochiai	0.244
Raup-Crick	0.190
χ^2	0.048
Jaccard	0.188
Sørensen	0.244

Notes: The abundance data are the abundances of carabids divided by the number of days the traps were active at each site, while the presence-absence data are the occurrence of species at each site. Results are given for all but the double-zero symmetrical coefficients. The coefficients are described in Table 1.

dissimilarity coefficient, as was the case for *S. haematopus*, consensus RDA is a better option.

DISCUSSION

This study presents a new approach to performing canonical ordination using a group of dissimilarity coefficients and proposes a new framework to analyze species communities using abundance and presence-absence data together.

A surprising result of this study is that the SAD of a community is not an important criterion for choosing a coefficient (Figs. 2, 3, and 6; Appendix B: Figs. B1–B8 and Appendix D: Figs. D1–D4) in canonical ordinations. This is what prompted us to develop consensus RDA. These results may also bring insight into the comparison of SADs, an important line of research (McGill et al. 2007). Using the results in Fig. 6 (and also Appendix D: Figs. D1–D4) obtained from abundance data, we can compare SADs, because the communities simulated with different SADs were correlated with the same true underlying structure of the data (Eq. 7 without the error term). The true underlying structure of the data serves as a reference to know how well a SAD determines the raw community structure because it is the basic information from which all species are constructed in the simulation. From the discussion in McGill et al. (2007) on SAD comparison, it can be expected that SADs defining notably different abundance patterns (e.g., Fig. 1b, l, m, q, u) would correlate differently with the true underlying structure of the data.

However, in Fig. 6 they all correlate equally well with the true underlying structure of the simulated communities. Moreover, the fairly broad range of the RV coefficient 95% confidence intervals for any one of the 25 SADs indicates that the variations in the raw multivariate community data can be surprisingly important, even if species have the same abundance structure. Such results may suggest that the SAD of a community may present only a small fraction of the information that characterizes a community matrix. However, further research is still needed to confirm the findings we made that the information lost when constructing SADs may make it difficult to develop a valuable approach to compare communities using SADs.

Our study shows that the choice of dissimilarity coefficients in canonical ordinations should primarily be based on the ecological knowledge available for the community under study. The ecological questions and the data type should guide the choice of one or a group of coefficients. Legendre and Legendre (2012, Table 7.4) offer a decision key designed to help ecologists select coefficients for community composition data based on data types (presence-absence or abundance) and type of information to be extracted. If a canonical analysis is performed using only one coefficient when more than one can potentially be used, Legendre and Gallagher (2001) would select the coefficient that explains the largest amount of variance. However, the properties of the selected coefficient may influence the interpretation.

If more than one coefficient is chosen, it is important to compare them using an MST based on dissimilarities of pairwise RV coefficients to determine if any of them present results markedly different from the others. This comparison can be seen as a selection procedure for coefficients. It evaluates the similarities between different RDA models where coefficients are the only element differentiating the models and finds which model(s) differ(s) notably from the others. This comparison will help decide if any coefficient(s) should be discarded. Comparing RDA models constructed using different coefficients through an MST is a first step in the analysis of a community through more than one coefficient.

When more than one coefficient presents similar information, a consensus RDA allows one to extract the most information out of the data because it focuses on the common information brought out by different coefficients. Using only one coefficient may put too much emphasis on a particular aspect of the data because each coefficient was designed to highlight different particularities of a community matrix. This may lead to a suboptimal ecological interpretation. Consensus RDA prevents this problem from occurring by extracting only the common information generated by a group of coefficients. In that respect, consensus RDA indirectly solves the technical problem of choosing a coefficient by using all the ones that can be suitable to analyze the data. Also, because it diminishes the importance of the information highlighted by one or a

few coefficients, it produces a result less influenced by the mathematical properties of a coefficient. For this reason, consensus RDA gives a more accurate representation of a community and will help researchers better understand the factors structuring the species in the community they study.

Conceptually, the new canonical ordination procedure proposed in this study has similarities with model averaging (Burnham and Anderson 2004, Anderson 2008). In model averaging, the best models are given more weight than the poor ones. This can be related to the selection procedure we propose, where coefficients are considered independently, discarded, or used to construct a consensus model. In consensus RDA, the different RDA models are weighted by the sum of the canonical eigenvalues of their components included in the construction of the consensus. In that respect, models constructed with different dissimilarity coefficients have different weights, which will influence the consensus; this is another similarity with model averaging. However, model averaging is more flexible because the choice of variables may vary between models, whereas only dissimilarity coefficients vary in consensus RDA.

As Økland (1996) pointed out, unconstrained ordination is a useful method to generate hypotheses when no explanation of the community variation has been proposed, whereas the main purpose of constrained ordination is hypothesis testing. With the development of constrained ordination methods, more complex analyses have been proposed and used by ecologists. For example, one may test a hypothesis by RDA using a set of explanatory variables or experimental factors, then examine the PCA ordination of the non-canonical variation to generate new hypotheses about the origin of the residual variation not explained by the explanatory variables.

Using a different approach, Borcard and Legendre (1994) showed that spatially constrained ordination of community composition data can help ecologists generate hypotheses about the processes that produced the spatial variation of the community. In their 1994 paper, they used a polynomial of the geographic coordinates as a constraining factor in CCA. In subsequent papers, they developed spatial eigenfunction analysis based on Moran's eigenvector maps (MEM, originally called PCNM; Borcard and Legendre 2002, Borcard et al. 2004, Dray et al. 2006); this is a much more powerful method for modeling fine-scaled spatial variation. Because ecological data are often spatially correlated (Legendre 1993), inclusion of spatial variables such as MEMs in canonical ordination is important to understand and test the significance of species-environment relationships. Dray et al. (2012) reviewed different ways of considering space in community ecology and including it in canonical ordinations.

Another aspect that researchers need to consider when performing canonical ordinations such as CCA, RDA, or consensus RDA is that these methods compute

a linear model of the explanatory variables for each species in the community. In the case of CCA, the species data are chi-square transformed before computing multiple regressions, and the regression involves the total abundances of the sites as weights (Legendre and Legendre 2012: Section 11.2). Qualitative explanatory variables (factors) can be included in these models, as is also the case in multiple regression. Because the species–environment relationships in nature are not necessarily linear, it has been proposed to include polynomials of the explanatory variables in the explanatory matrix, instead of the explanatory variables only, to make it possible to model nonlinear relationships between the species and the explanatory variables (e.g., Legendre and Legendre 2012: Ecological Application 14.1b). Note that dissimilarity coefficients and data transformations do not account by themselves for the nonlinearity of the species–environment relationship. They were designed to give more (or less) weight to common (or rare) species and to account for the double-zero problem. This approach can be applied to all canonical ordination methods, including consensus RDA.

A problem that we have not approached but warrants further investigation is selection of explanatory variables in consensus RDA. Methods such as forward selection (e.g., Blanchet et al. 2008) assume that an RDA is performed using only one dissimilarity coefficient. Consensus RDA requires all explanatory variables to be the same and that only the coefficient differs between RDAs. If an automatic variable selection procedure is used independently for each RDA, it is likely that different sets of variables will be selected. In this situation, we propose three variable selection approaches. (1) A consensus analysis should employ the union of all explanatory variables selected for the various coefficients. That is, if for a coefficient, explanatory variables A and B are selected, and with another coefficient it is explanatory variables A and C that are chosen, the union of the explanatory variables for the consensus RDA would be variables A, B, and C. Using this approach, one can at least eliminate the explanatory variables that are totally useless. This idea of using the union of the selected variables is inspired by the selection method of Peres-Neto and Legendre (2010) for Moran's eigenvector maps eigenfunctions. (2) The variable selection is carried out on the consensus RDA result without any variable selection carried out on individual RDAs. (3) Use the union of the selected variables on individual RDAs, as explained in (1), and then carry out a further selection for the consensus RDA. Further studies will need to be carried out to evaluate which of these three approaches yields the models that best define a species community.

Species-abundance data contain more information than presence–absence data and often lead to a better understanding of community variations through RDA, although community ecologists generally consider that the single most important information about a species is

its presence. However, for certain organisms, abundance data are not reliable. In palynology, for example, presence–absence data are often favored because abundance data are subject to large bias (Davis 2000). Presence–absence data are also more suitable when studying ant communities using pitfall traps because the ants' social behavior and propensity at creating foraging trails has an enormous influence on abundance data (Higgins and Lindgren 2012). Similarly, in studies of fish biodiversity, variation in size of fish species living in the same area demands that different instruments be used to catch them, and thus the abundance data are not comparable. The only way to consider all species of fish together in a consistent analysis is by using presence–absence data (biomass data can also be used for fish of all sizes caught by electrofishing or recorded during underwater visual census). This is likely to be true for any communities where variations in size between species require that different trapping methods be used to catch enough species to have a representative fraction of the studied species community.

When working with presence–absence data, we suggest that one should first draw a species presence distribution, as we did in Fig. 7d. The ratio between common and rare species should serve as a general guideline when drawing ecological conclusions. Although it is possible that canonical ordinations performed on presence–absence data show biased results, it is more likely that such ordinations can be complementary to those computed for abundance data. Certain environmental factors may be necessary for a species to occur in an area (e.g., certain plant species are found only in the presence of certain geological formations) while other factors may make species abundances vary (e.g., precipitation). Variation in abundance is efficient in describing how a species is related to a gradient (environmental, physical, or other). However, species abundances may conceal the strict relationship a species has with its habitat. This strict relationship is what makes a species occur or not occur at a site. In that respect, considering both abundance and presence–absence data may be ecologically valuable to better understand the factors structuring a community. The idea to use both abundance and presence–absence data to better understand an ecological system has been proposed before (see e.g., Van Buskirk 2005). As explained in the previous paragraphs, abundance data may be unreliable when sampling certain groups of organisms. However, for all communities where species abundances can be sampled without diminishing the value of the data, presence–absence data can be easily obtained by transforming all abundances larger than 0 to 1s, allowing ecologists to get a more complete understanding of the data they collected.

In this study we presented a new approach to perform canonical ordination in community ecology research. This approach has the potential to be used in other fields of research where the structure of the data is similar to

that of community ecology. Population and landscape genetics are examples of research areas where consensus RDA could potentially be useful.

ACKNOWLEDGMENTS

We are grateful to Xianli Wang, John R. Spence, and Dave W. Roberts for insightful comments on an early draft of the manuscript. This research was supported by GEOIDE Canada and an NSERC to F. He and NSERC grant number 7738 to P. Legendre.

LITERATURE CITED

- Anderson, D. R. 2008. Model based inference in the life sciences—a primer on evidence. Springer, New York, New York, USA.
- Anderson, M. J., K. E. Ellingsen, and B. H. McCordle. 2006. Multivariate dispersion as a measure of beta diversity. *Ecology Letters* 9:683–693.
- Anderson, M. J., et al. 2011. Navigating the multiple meanings of beta diversity: a roadmap for the practicing ecologist. *Ecology Letters* 14:19–28.
- Bergeron, J. A. C., F. G. Blanchet, J. R. Spence, and W. J. A. Volney. 2012. Ecosystem classification and inventory maps as surrogates for ground beetle assemblages in boreal forest. *Journal of Plant Ecology* 5:97–108.
- Bergeron, J. A. C., J. R. Spence, and W. J. A. Volney. 2011. Landscape patterns of species-level associations between ground-beetles (Coleoptera: Carabidae) and overstory trees in boreal forests of western Canada (Coleoptera: Carabidae). *ZooKeys* 147:577–600.
- Blanchet, F. G., J. A. C. Bergeron, J. R. Spence, and F. He. 2013. Landscape effects of disturbance, habitat heterogeneity and spatial autocorrelation for a ground beetle (Carabidae) assemblage in mature boreal forest. *Ecography* 36:636–647.
- Blanchet, F. G., P. Legendre, and D. Borcard. 2008. Forward selection of explanatory variables. *Ecology* 89:2623–2632.
- Borcard, D., and P. Legendre. 1994. Environmental control and spatial structure in ecological communities: an example using oribatid mites (*Acari, Oribatei*). *Environmental and Ecological Statistics* 1:37–61.
- Borcard, D., and P. Legendre. 2002. All-scale spatial analysis of ecological data by means of principal coordinates of neighbour matrices. *Ecological Modelling* 153:51–68.
- Borcard, D., P. Legendre, C. Avois-Jacquet, and H. Tuomisto. 2004. Dissecting the spatial structure of ecological data at multiple scales. *Ecology* 85:1826–1832.
- Borcard, D., P. Legendre, and P. Drapeau. 1992. Partialling out the spatial component of ecological variation. *Ecology* 73:1045–1055.
- Bray, J. R., and J. T. Curtis. 1957. An ordination of the upland forest communities of southern Wisconsin. *Ecological Monographs* 27:325–349.
- Burnham, K. P., and D. R. Anderson. 2004. Multimodel inference—understanding AIC and BIC in model selection. *Sociological Methods and Research* 33:261–304.
- Clarke, K. R., and R. H. Green. 1988. Statistical design and analysis for a biological effects study. *Marine Ecology Progress Series* 46:213–226.
- Davis, M. B. 2000. Palynology after Y2K—understanding the source area of pollen in sediments. *Annual Review of Earth and Planetary Sciences* 28:1–18.
- Dewdney, A. K. 2000. A dynamical model of communities and a new species-abundance distribution. *Biological Bulletin* 198:152–165.
- Dray, S., and A. B. Dufour. 2007. The ade4 package: implementing the duality diagram for ecologists. *Journal of Statistical Software* 22:1–20.
- Dray, S., et al. 2012. Community ecology in the age of multivariate multiscale spatial analysis. *Ecological Monographs* 82:257–275.
- Dray, S., P. Legendre, and P. Peres-Neto. 2006. Spatial modelling: a comprehensive framework for principal coordinate analysis of neighbor matrices (PCNM). *Ecological Modelling* 196:483–493.
- Escoufier, Y. 1973. Le traitement des variables vectorielles. *Biometrics* 29:751–760.
- Fisher, R. A., A. S. Corbet, and C. B. William. 1943. The relation between the number of species and the number of individuals in a random sample of an animal population. *Journal of Animal Ecology* 12:42–58.
- Gaston, K. J. 2010. Valuing common species. *Science* 327:154–155.
- Gower, J. C. 1966. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53:325–338.
- Gower, J. C., and P. Legendre. 1986. Metric and Euclidean properties of dissimilarity coefficients. *Journal of Classification* 3:5–48.
- Gray, J. S., A. Bjørgesaeter, and K. I. Ugland. 2006. On plotting species abundance distributions. *Journal of Animal Ecology* 75:752–756.
- Greenacre, M. 2013. The contributions of rare objects in correspondence analysis. *Ecology* 94:241–249.
- Higgins, R. J., and B. S. Lindgren. 2012. An evaluation of methods for sampling ants (Hymenoptera: Formicidae) in British Columbia, Canada. *Canadian Entomologist* 144:491–507.
- Hill, M. O., and H. G. Gauch. 1980. Detrended correspondence analysis, an improved ordination technique. *Vegetation* 42:47–58.
- Hurlbert, S. H. 1984. Pseudoreplication and the design of ecological field experiments. *Ecological Monographs* 54:187–211.
- Hutchinson, G. 1957. Concluding remarks. *Cold Spring Harbor Symposium on Quantitative Biology* 22:415–427.
- Jaccard, P. 1901. Étude comparative de la distribution florale dans une portion des Alpes et du Jura. *Bulletin de la Société Vaudoise des Sciences naturelles* 37:547–579.
- Lebart, L., and J. P. Fénelon. 1971. *Statistique et Informatique Appliquées*. Dunod, Paris, France.
- Legendre, P. 1993. Spatial autocorrelation: trouble or new paradigm? *Ecology* 74:1659–1673.
- Legendre, P., and M. J. Anderson. 1999. Distance-based redundancy analysis: testing multispecies responses in multifactorial ecological experiments. *Ecological Monographs* 69:1–24.
- Legendre, P., D. Borcard, and P. R. Peres-Neto. 2005. Analyzing beta diversity: partitioning the spatial variation of community composition data. *Ecological Monographs* 75:435–450.
- Legendre, P., and M. De Cáceres. 2013. Beta diversity as the variance of community data: dissimilarity coefficients and partitioning. *Ecology Letters* 16:951–963.
- Legendre, P., and E. Gallagher. 2001. Ecologically meaningful transformations for ordination of species data. *Oecologia* 129:271–280.
- Legendre, P., and L. Legendre. 2012. *Numerical ecology*. Third English edition. Elsevier, Amsterdam, Netherlands.
- Loreau, M. 2010. From populations to ecosystems: theoretical foundations for a new ecological synthesis. Princeton University Press, Princeton, New Jersey, USA.
- Maor, E. 2007. The Pythagorean theorem: a 4,000-year history. Princeton University Press, Princeton, New Jersey, USA.
- McCoy, E. D., S. S. Bell, and K. Walters. 1986. Identifying biotic boundaries along environmental gradients. *Ecology* 67:749–759.
- McGill, B. J., et al. 2007. Species abundance distributions: moving beyond single prediction theories to integration within an ecological framework. *Ecology Letters* 10:995–1015.
- Motyka, J. 1947. O zadaniach i metodach badan geobotanicznych. Sur les buts et les méthodes des recherches géo-

- botaniques. Pages viii+168 in *Annales Universitatis Mariae Curie-Skłodowska* (Lublin, Polonia), Sectio C, Supplementum I.
- Niemelä, J. 1993. Mystery of the missing species: species-abundance distribution of boreal ground-beetles. *Annales Zoologici Fennici* 30:169–172.
- Ochiai, A. 1957. Zoogeographic studies on the soleoid fishes found in Japan and its neighbouring regions. *Bulletin of the Japanese Society of Scientific Fisheries* 22:526–530.
- Odum, E. P. 1950. Bird populations of the highlands (North Carolina) plateau in relation to plant succession and avian invasion. *Ecology* 31:587–605.
- Økland, R. H. 1996. Are ordination and constrained ordination alternative or complementary strategies in general ecological studies? *Journal of Vegetation Science* 7:289–292.
- Oksanen, J., F. G. Blanchet, R. Kindt, P. Legendre, P. R. Minchin, R. B. O'Hara, G. L. Simpson, P. Sölymos, M. H. H. Stevens, and H. Wagner. 2013. *vegan: community ecology package*. <http://CRAN.R-project.org/package=vegan>
- Orlói, L. 1967. An agglomerative method for classification of plant communities. *Journal of Ecology* 55:193–206.
- Pearson, K. 1901. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine* 2:559–572.
- Peres-Neto, P. R., and P. Legendre. 2010. Estimating and controlling for spatial structure in the study of ecological communities. *Global Ecology and Biogeography* 19:174–184.
- Preston, F. W. 1948. The commonness, and rarity, of species. *Ecology* 29:254–283.
- R Development Core Team. 2012. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org/>
- Rao, C. R. 1964. The use and interpretation of principal component analysis in applied research. *Sankhya, the Indian Journal of Statistics* 26:329–358.
- Rao, C. R. 1995. A review of canonical coordinates and an alternative to correspondence analysis using Hellinger distance. *Quaestio* 19:23–63.
- Raup, D. M., and R. E. Crick. 1979. Measurement of faunal similarity in paleontology. *Journal of Paleontology* 53:1213–1227.
- Robert, P., and Y. Escoufier. 1976. A unifying tool for linear multivariate statistical methods: the RV-coefficient. *Applied Statistics* 25:257–265.
- Roux, G., and M. Roux. 1967. À propos de quelques méthodes de classification en phytosociologie. *Revue de Statistique Appliquée* 15:59–72.
- Shepard, R. N. 1962. The analysis of proximities: multidimensional scaling with an unknown distance. I. *Psychometrika* 27:125–140.
- Sokal, R., and C. Michener. 1958. A statistical method for evaluating systematic relationships. *University of Kansas Scientific Bulletin* 38:1409–1438.
- Sørensen, T. 1948. A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analysis of vegetation on Danish commons. *Biologiske skrifter* 5:1–34.
- Spence, J. R., and J. K. Niemelä. 1994. Sampling carabid assemblages with pitfall traps: the madness and the method. *Canadian Entomologist* 126:881–894.
- ter Braak, C. J. F. 1986. Canonical correspondence-analysis—a new eigenvector technique for multivariate direct gradient analysis. *Ecology* 67:1167–1179.
- ter Braak, C. J. F. 1987. The analysis of vegetation–environment relationships by canonical correspondence analysis. *Vegetatio* 69:69–77.
- ter Braak, C. J. F. 1994. Canonical community ordination. Part I: basic theory and linear methods. *Écoscience* 1:127–140.
- ter Braak, C. J. F., and P. F. M. Verdonschot. 1995. Canonical correspondence analysis and related multivariate methods in aquatic ecology. *Aquatic Sciences Research Across Boundaries* 57:255–289.
- Van Buskirk, J. 2005. Local and landscape influence on amphibian occurrence and abundance. *Ecology* 86:1936–1947.
- Whittaker, R. H. 1967. Gradient analysis of vegetation. *Biological Review* 49:207–264.

SUPPLEMENTAL MATERIAL

Appendix A

Explanation of the construction of the explanatory variables and how they were combined for the simulation study (*Ecological Archives* M084-018-A1).

Appendix B

Comparison of association coefficients using a coefficient of determination (R^2) (*Ecological Archives* M084-018-A2).

Appendix C

Comparison of consensus RDA constructed using only the significant canonical axes with consensus RDA constructed with all canonical axes (*Ecological Archives* M084-018-A3).

Appendix D

Comparison of canonical ordination models for abundance and presence–absence data using simulations (*Ecological Archives* M084-018-A4).

Appendix E

Species code and names for Carabidae and tree species (*Ecological Archives* M084-018-A5).

Supplement

The R package **ordiconsensus** compiled for all platforms (*Ecological Archives* M084-018-S1).

Data Availability

Data associated with this paper have been deposited in Dryad: <http://dx.doi.org/10.5061/dryad.8gs3n>

APPENDIX A

Ecological Archives EXXX-XXX-A1

EXPLANATION OF THE CONSTRUCTION OF THE EXPLANATORY VARIABLES AND HOW THEY WERE COMBINED FOR THE SIMULATION. ONE FIGURE (FIG. A1), ONE TABLE (TABLE A1) AND R CODE.

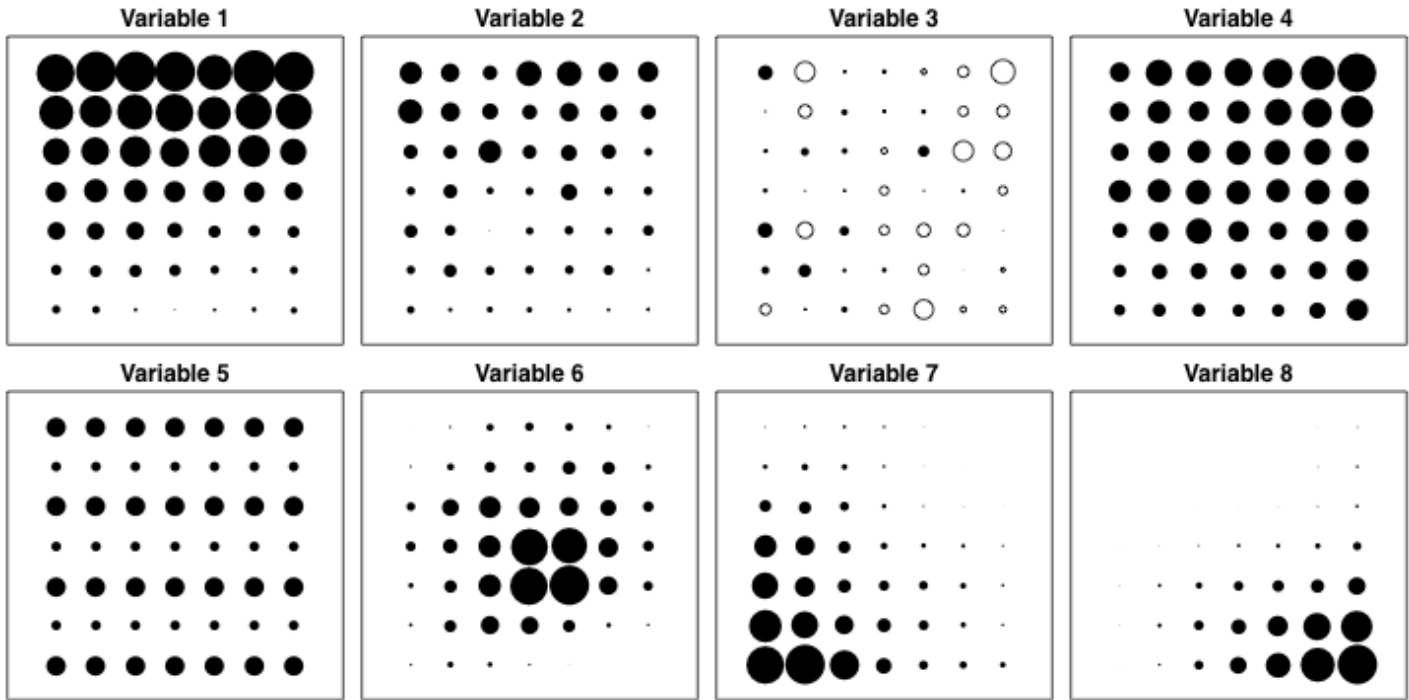


Fig. A1. Eight variables used in the construction of the simulated species

These variables were constructed using the RsimSSDCOMPAS package through the R statistical language using the following R code:

```
variable1<-SimSSDR(7,7,1,10,range11=5,range12=5,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=TRUE,
  SAR=FALSE)$E
```

```
variable2<-SimSSDR(7,7,1,5,range11=1,range12=1,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=TRUE,
  SAR=FALSE)$E
```

```
variable3<-SimSSDR(7,7,0,range11=5,range12=5)$E
```

```
variable4<-SimSSDR(7,7,2,10,range11=5,range12=5,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=TRUE,
  SAR=FALSE)$E
```

```

variable5<-SimSSDR(7,7,4,5,range11=2,range12=2,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=FALSE,
  SAR=FALSE)$E

variable6<-SimSSDR(7,7,3,10,range11=10,range12=10,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=FALSE,
  SAR=FALSE,centroide=list(c(0,0)))$E

variable7<-SimSSDR(7,7,3,10,range11=10,range12=10,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=FALSE,
  SAR=FALSE,centroide=list(c(1,1)))$E

variable8<-SimSSDR(7,7,3,10,range11=10,range12=10,range21=1,range22=1,
  nsp1=1,nsp2=1,varnor=list(rep(0,3)),SAE=FALSE,
  SAR=FALSE,centroide=list(c(10,0)))$E

```

Table A1: Combinations of explanatory variables and weight (regression coefficient) used to construct each species. The number associated to each species is the order given in the site-by-species table

Species	Explanatory variables combined	Weight given to (regression coefficient of) each species
1	1 and 4	2
2	1 and 5	0.1
3	1 and 6	-2
4	1 and 7	-0.1
5	1 and 8	2
6	2 and 3	0.5
7	2 and 5	-2
8	2 and 6	-0.5
9	2 and 7	2
10	2 and 8	1
11	3 and 5	-2
12	3 and 6	-1
13	3 and 7	2
14	3 and 8	0.5
15	4 and 5	-2
16	4 and 6	-0.5
17	4 and 7	2
18	4 and 8	0.1
19	5 and 8	-2
20	6 and 7	-0.1

APPENDIX B

Ecological Archives EXXX-XXX-A2

COMPARISON OF ASSOCIATION COEFFICIENTS USING A COEFFICIENT OF DETERMINATION (R^2). EIGHT FIGURES (FIGS. B1, B2, B3, B4, B5, B6, B7, AND B8)

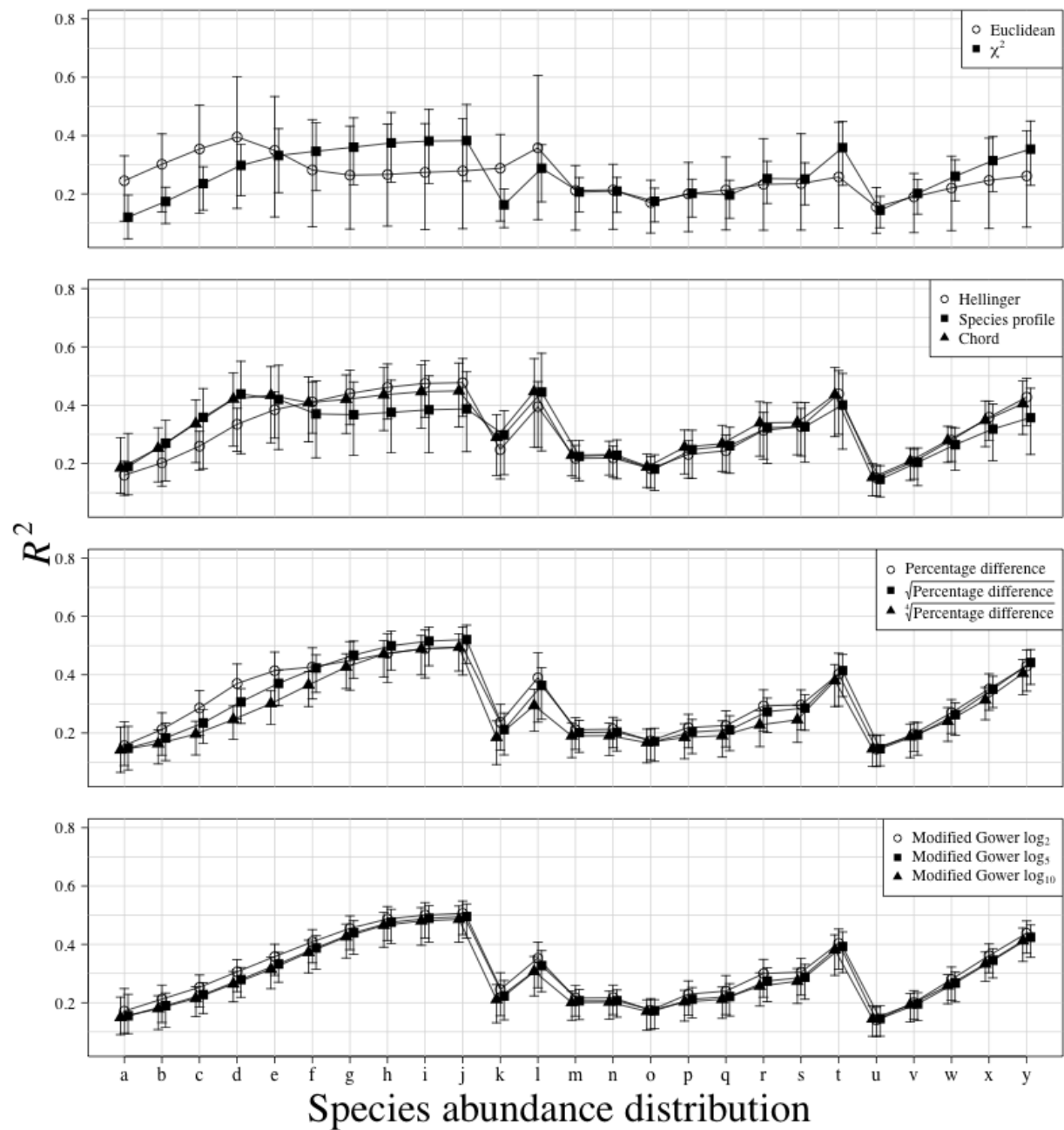


Fig. B1. Comparison of explained variance (R^2) between 11 association coefficients calculated on abundance data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.25). A thousand simulations were run for each SAD.

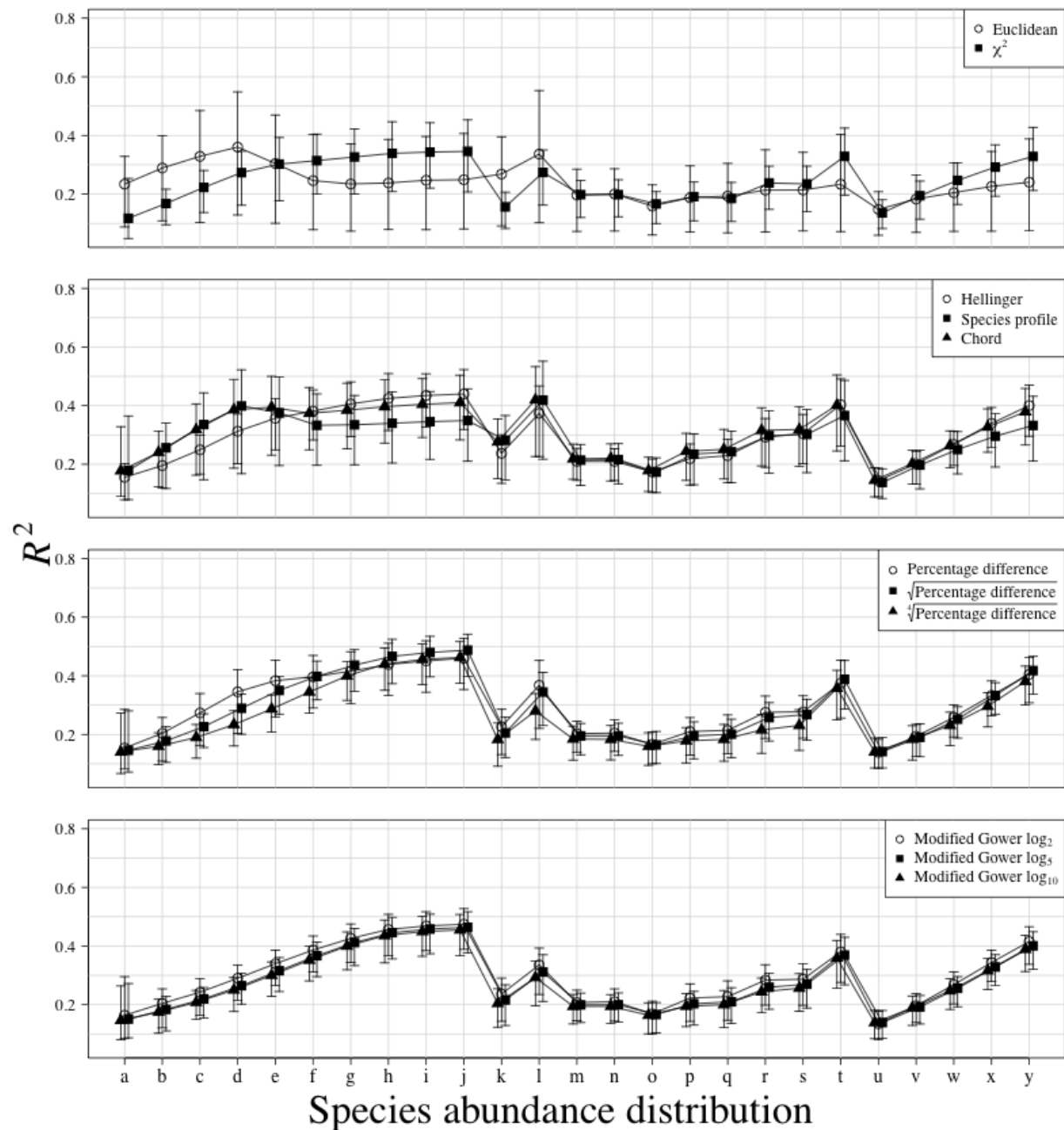


Fig. B2. Comparison of explained variance (R^2) between 11 association coefficients calculated on abundance data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.5). A thousand simulations were run for each SAD.

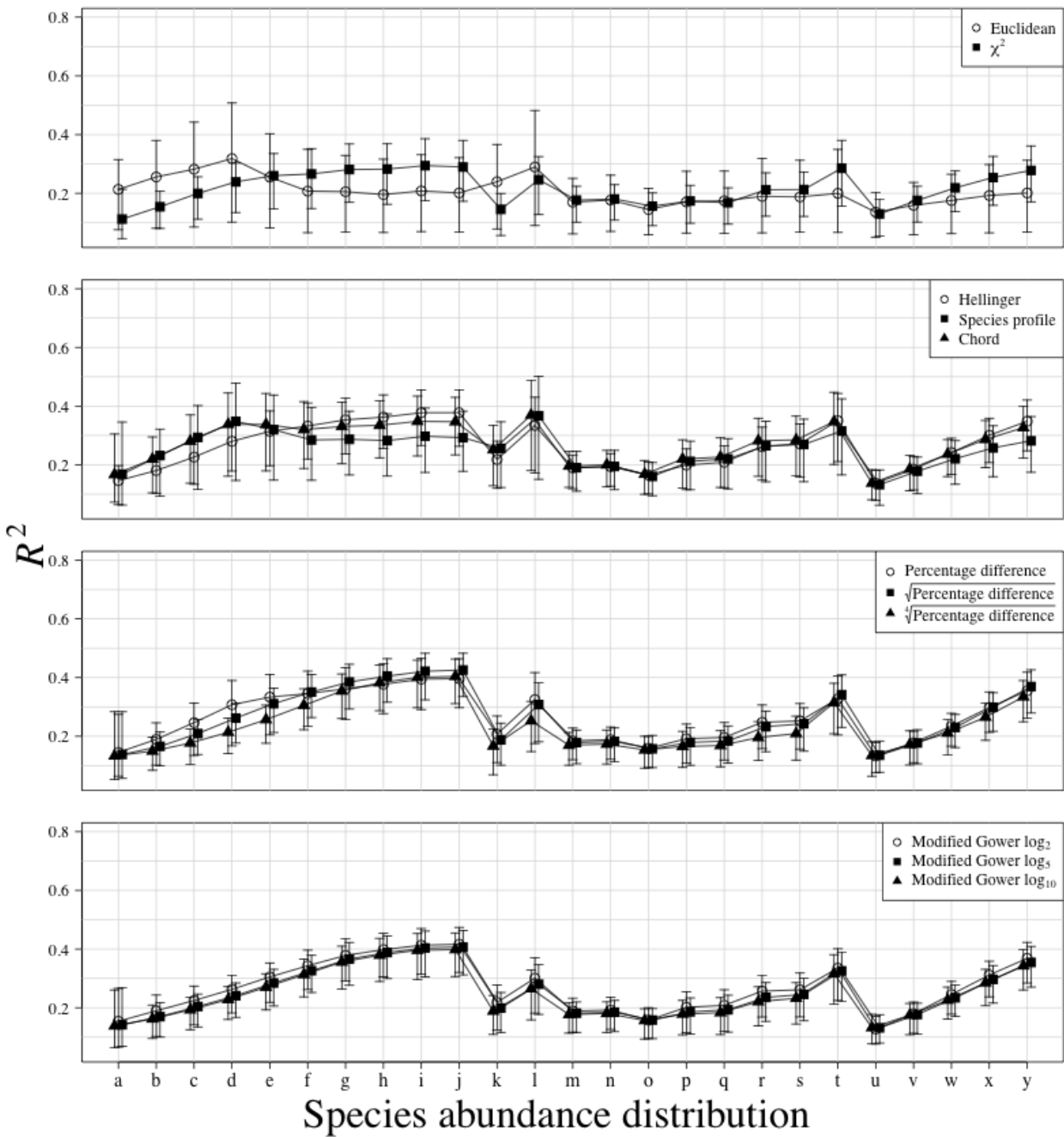


Fig. B3. Comparison of explained variance (R^2) between 11 association coefficients calculated on abundance data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 1). A thousand simulations were run for each SAD.

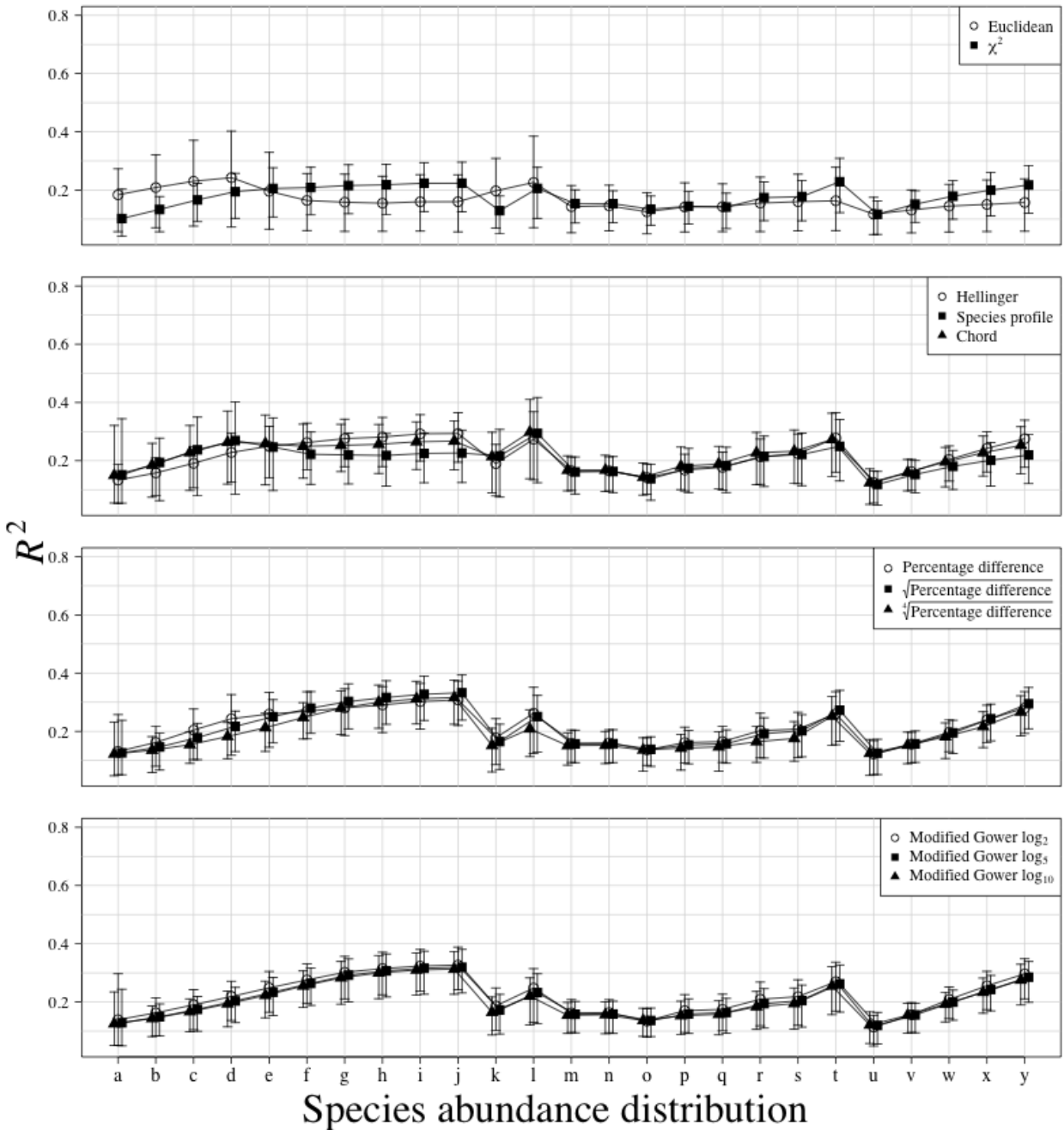


Fig. B4. Comparison of explained variance (R^2) between 11 association coefficients calculated on abundance data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 2). A thousand simulations were run for each SAD.

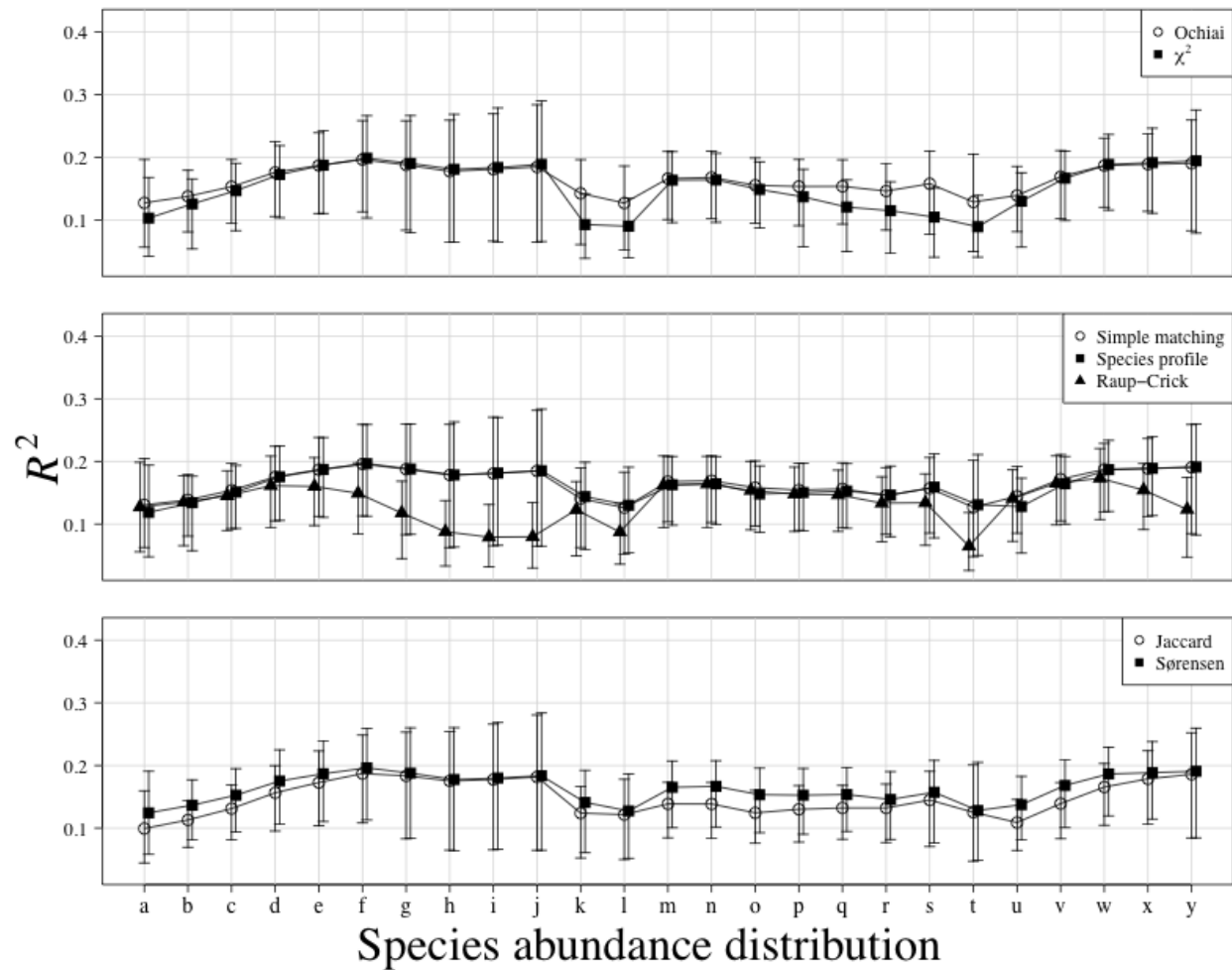


Fig. B5. Comparison of explained variance (R^2) between 6 association coefficients calculated on presence-absence data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.25). A thousand simulations were run for each SAD.

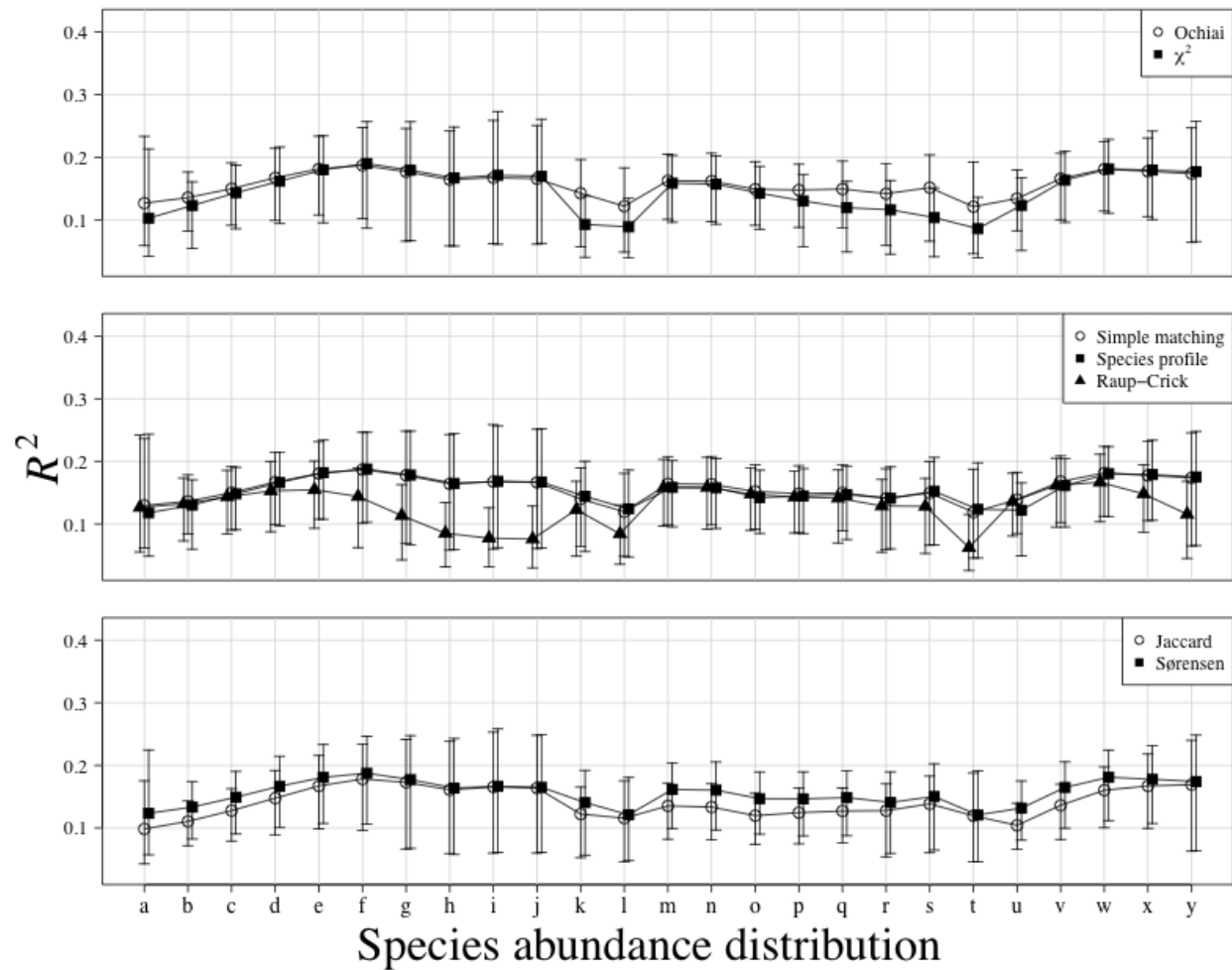


Fig. B6. Comparison of explained variance (R^2) between 6 association coefficients calculated on presence-absence data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.5). A thousand simulations were run for each SAD.

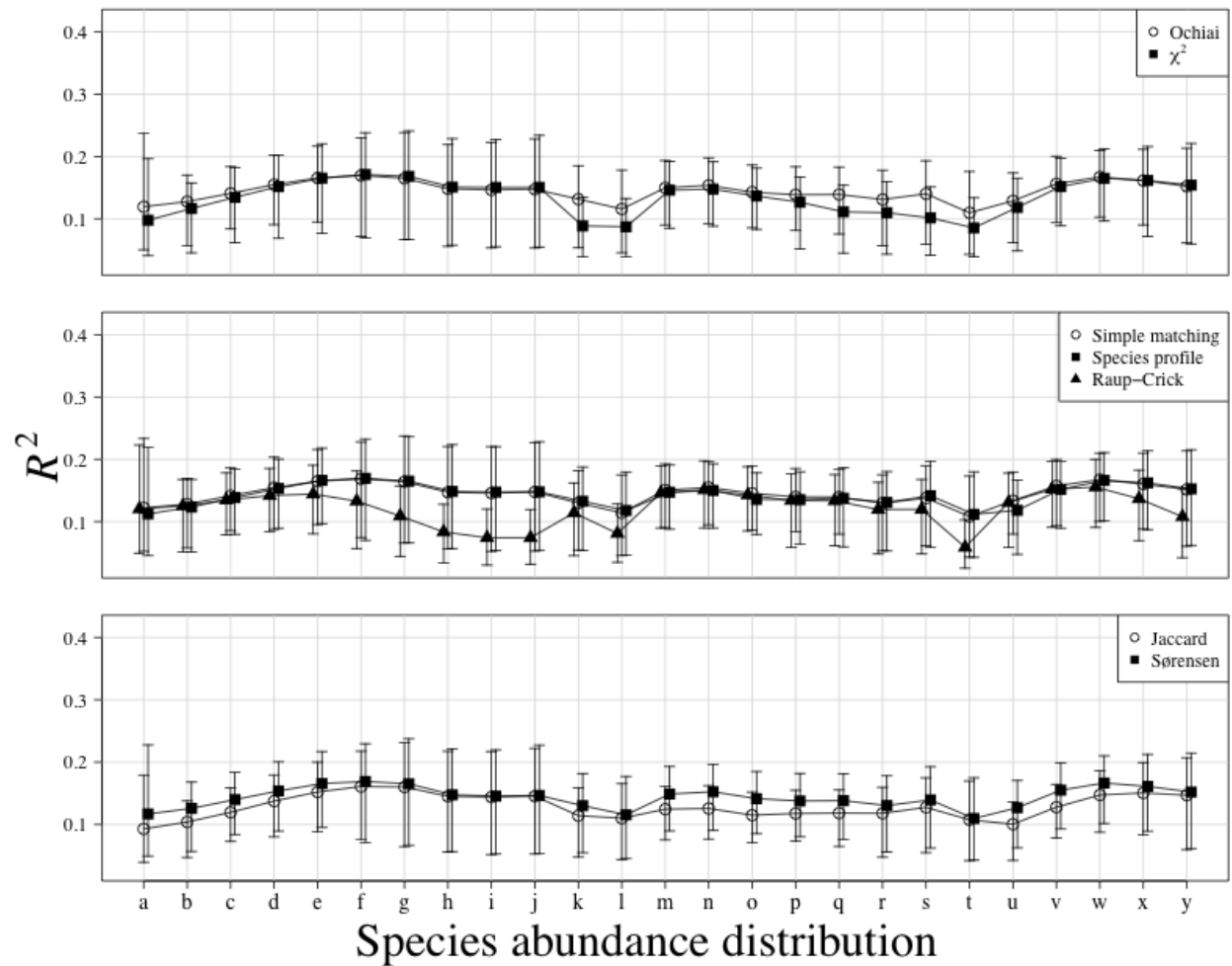


Fig. B7. Comparison of explained variance (R^2) between 6 association coefficients calculated on presence-absence data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 1). A thousand simulations were run for each SAD.

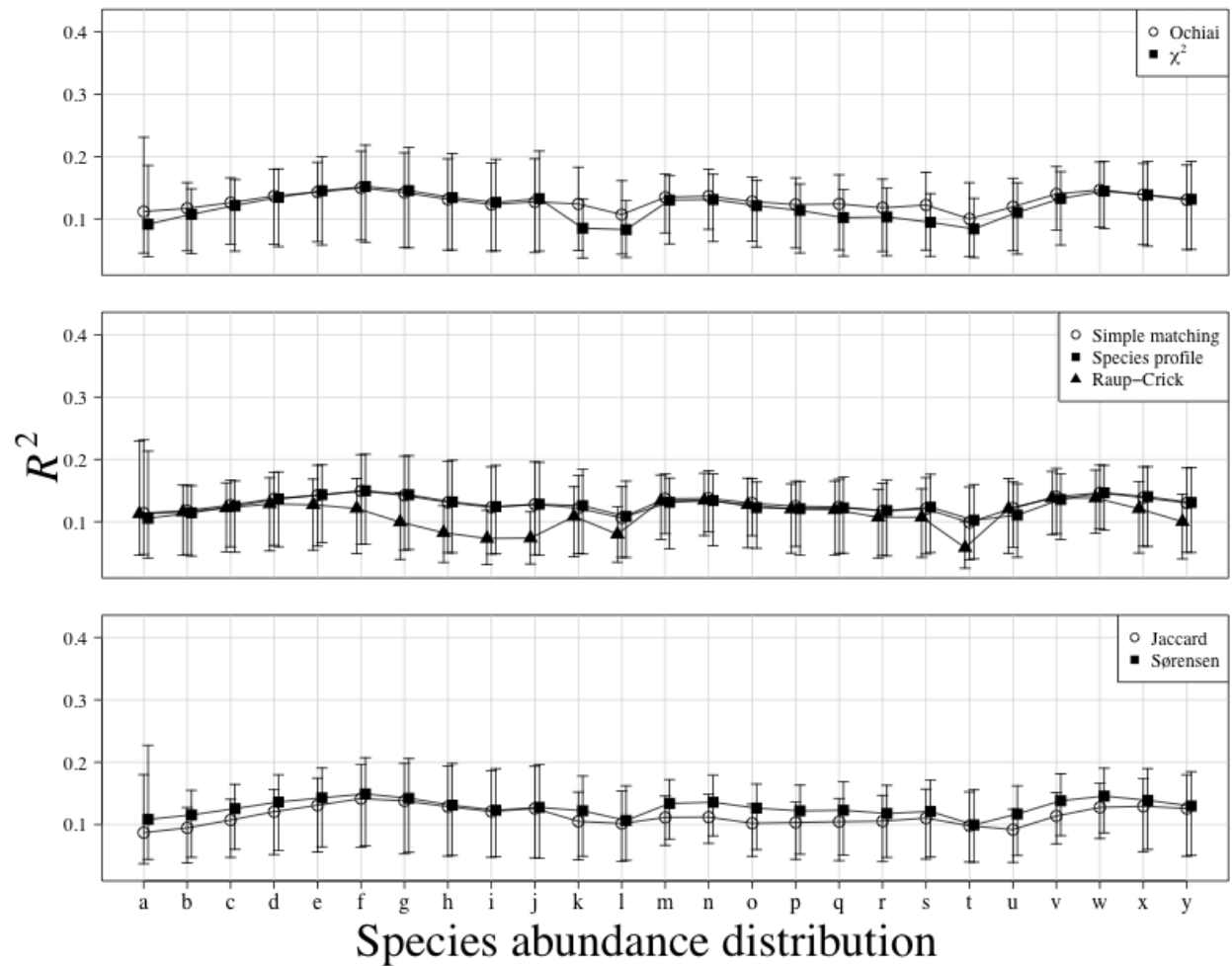


Fig. B8. Comparison of explained variance (R^2) between 6 association coefficients calculated on presence-absence data. Only the significant ($P \leq 0.05$ after 999 permutations) canonical axes were conserved to calculate R^2 . Points are R^2 means of all simulations and error bars represent 95% confidence intervals. Association coefficients are presented in different panels for visual clarity. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 2). A thousand simulations were run for each SAD.

APPENDIX C

Ecological Archives EXXX-XXX-A3

COMPARISON OF CONSENSUS RDA CONSTRUCTED USING ONLY SIGNIFICANT CANONICAL AXES
WITH CONSENSUS RDA CONSTRUCTED WITH ALL CANONICAL AXES. NINE FIGURES (FIGS. C1, C2,
C3, C4, C5, C6, C7, C8, AND C9)

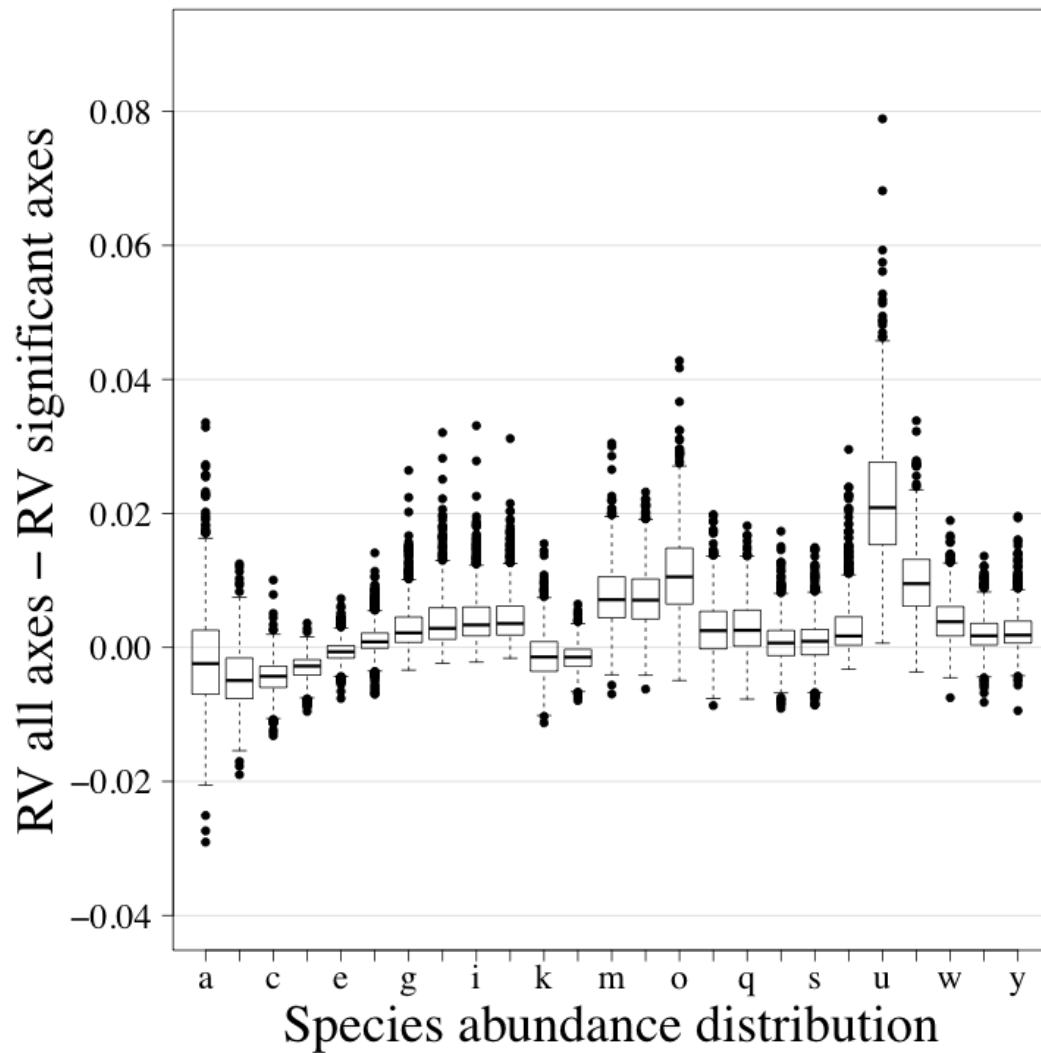


FIGURE C1. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from abundance data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.001). A thousand simulations were run for each SAD.

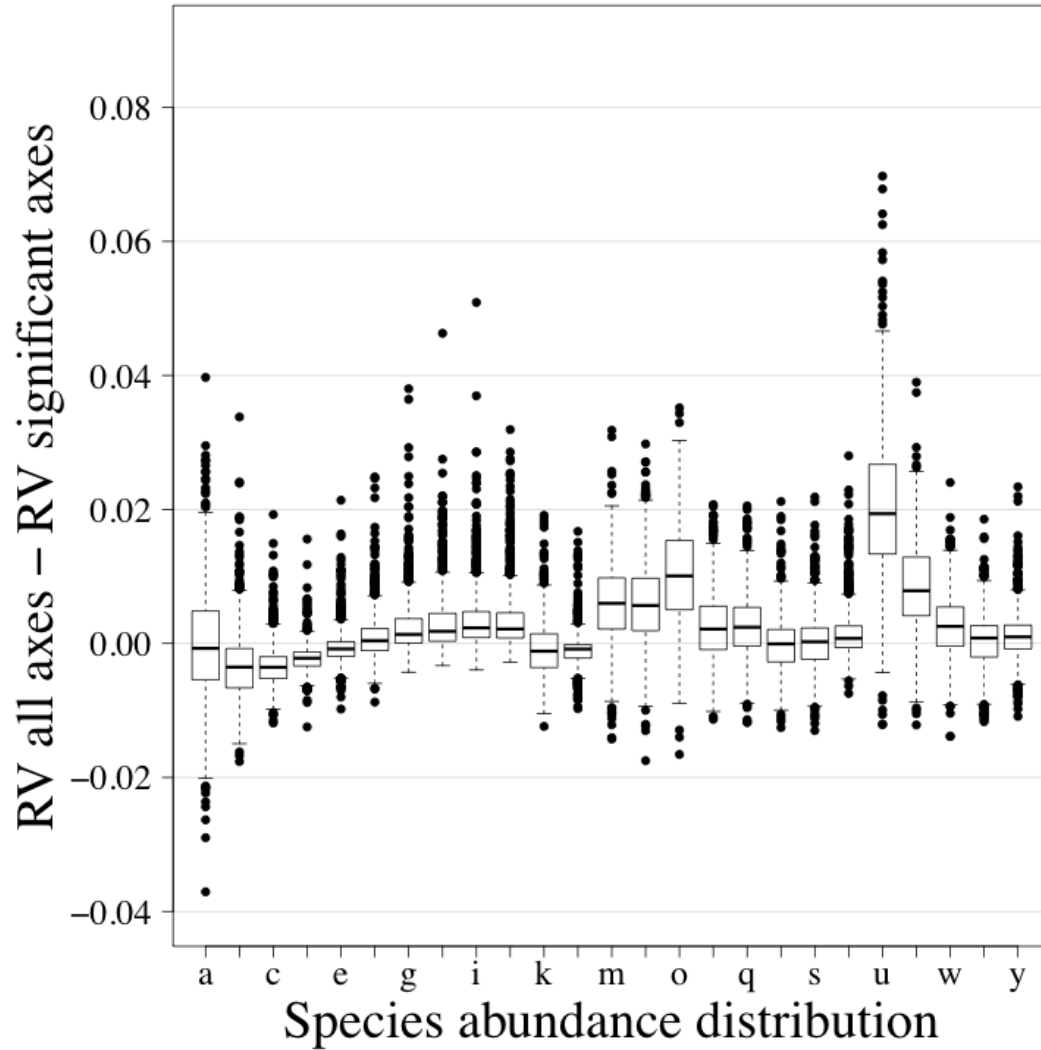


FIGURE C2. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from abundance data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.25). A thousand simulations were run for each SAD.

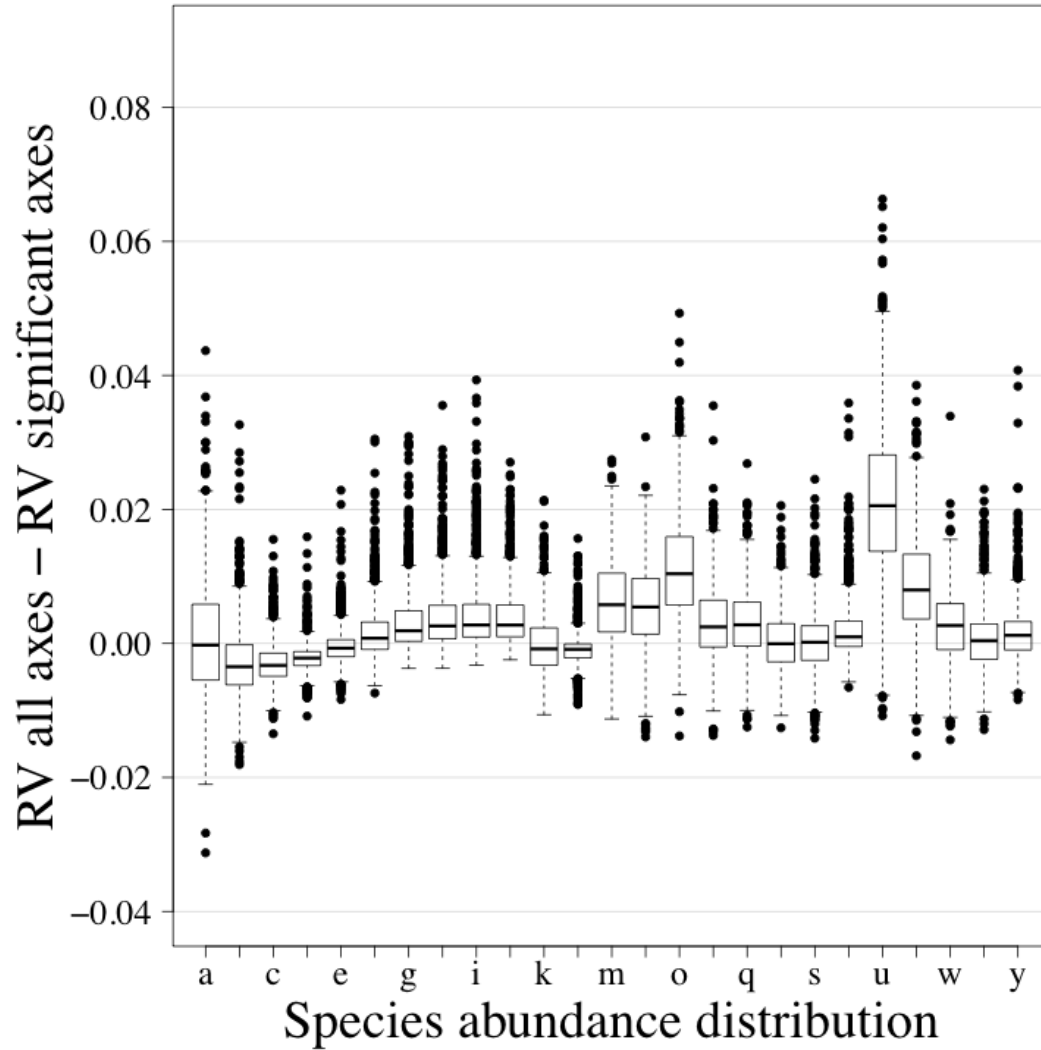


FIGURE C3. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from abundance data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.5). A thousand simulations were run for each SAD.

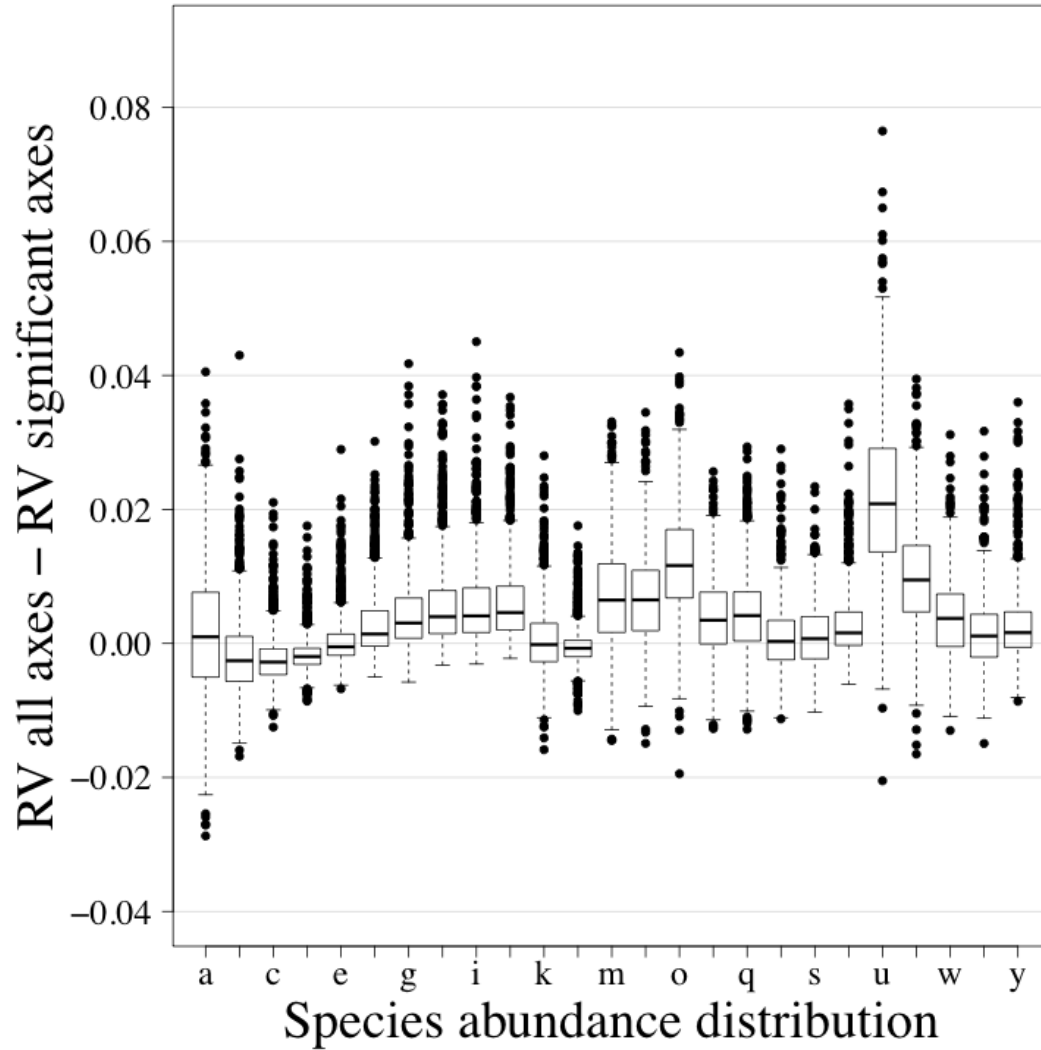


FIGURE C4. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from abundance data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 1). A thousand simulations were run for each SAD.

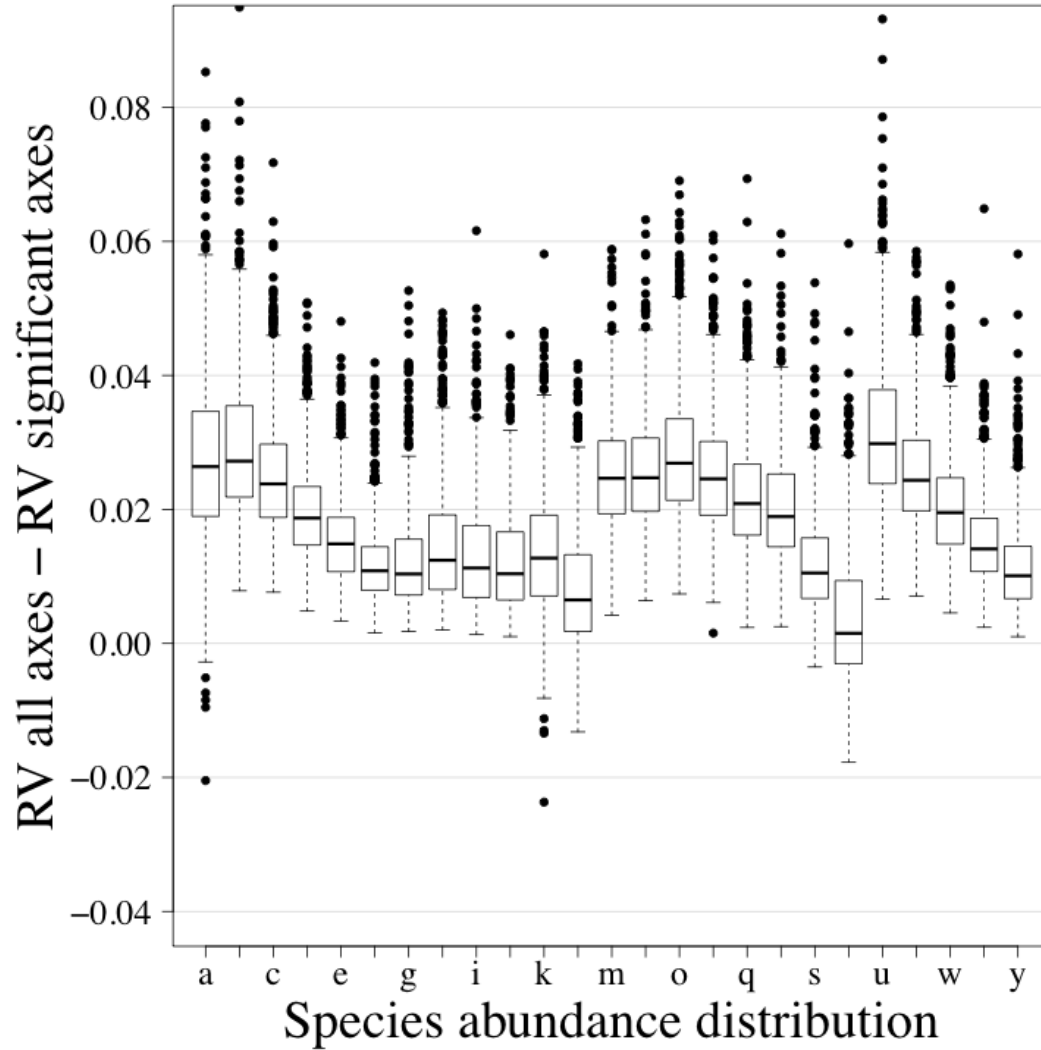


FIGURE C5. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from presence-absence data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.001). A thousand simulations were run for each SAD.

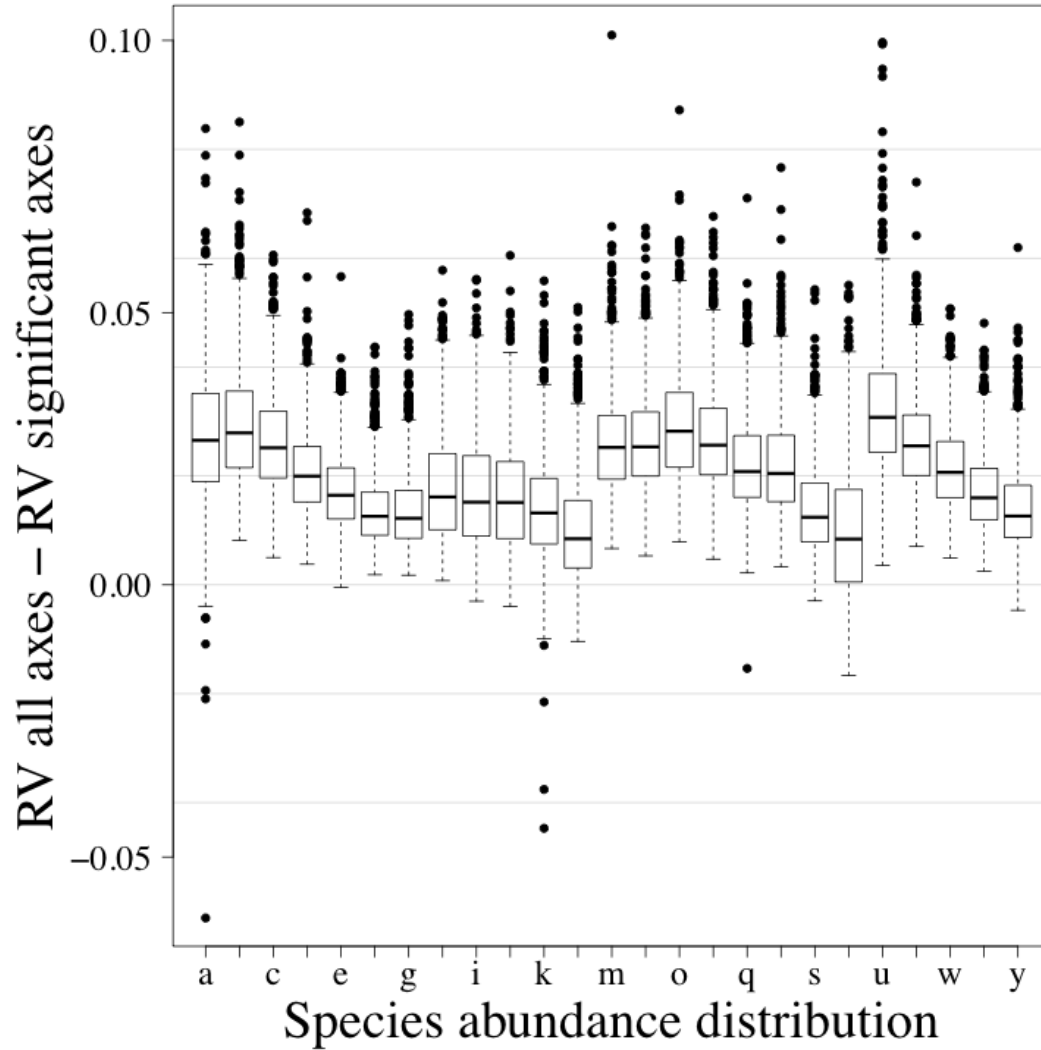


FIGURE C6. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from presence-absence data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.25). A thousand simulations were run for each SAD.

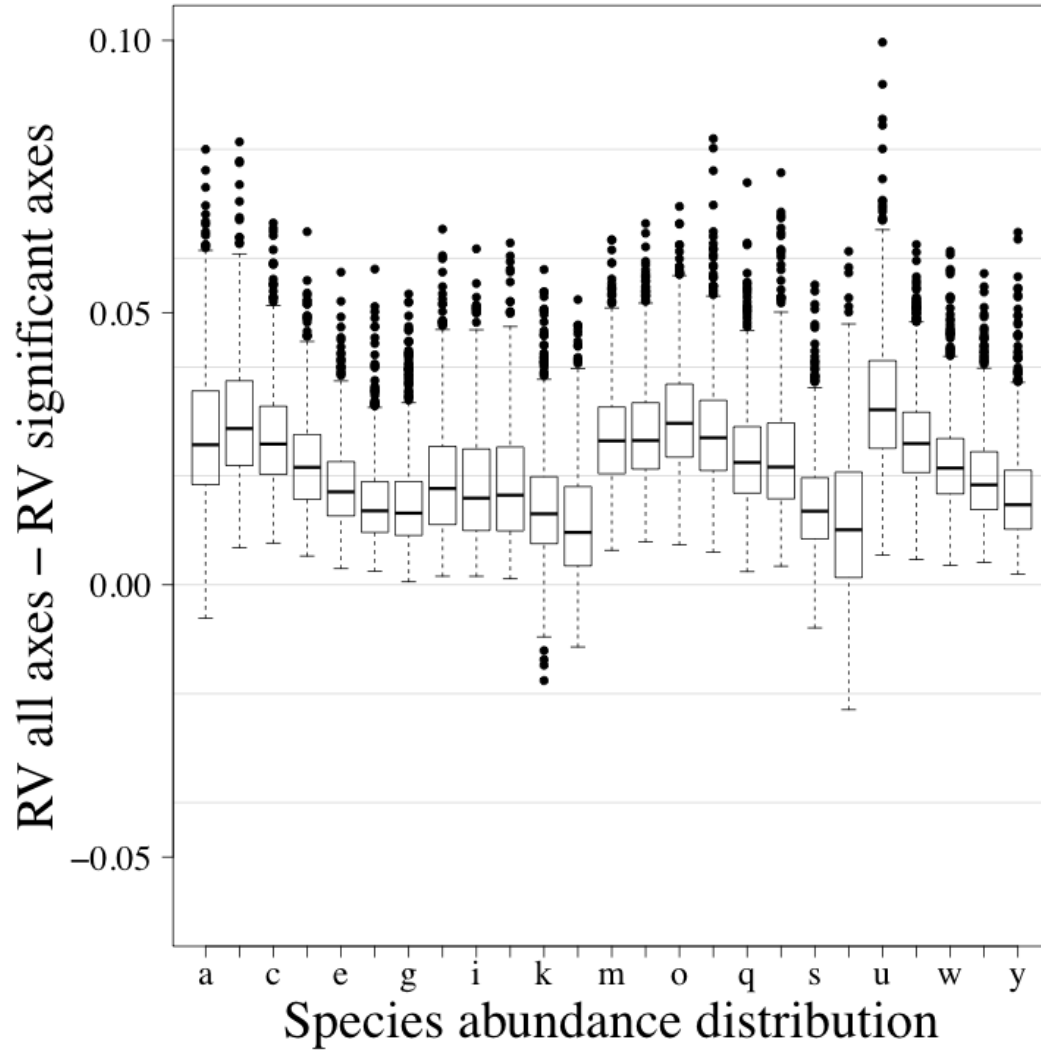


FIGURE C7. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from presence-absence data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.5). A thousand simulations were run for each SAD.

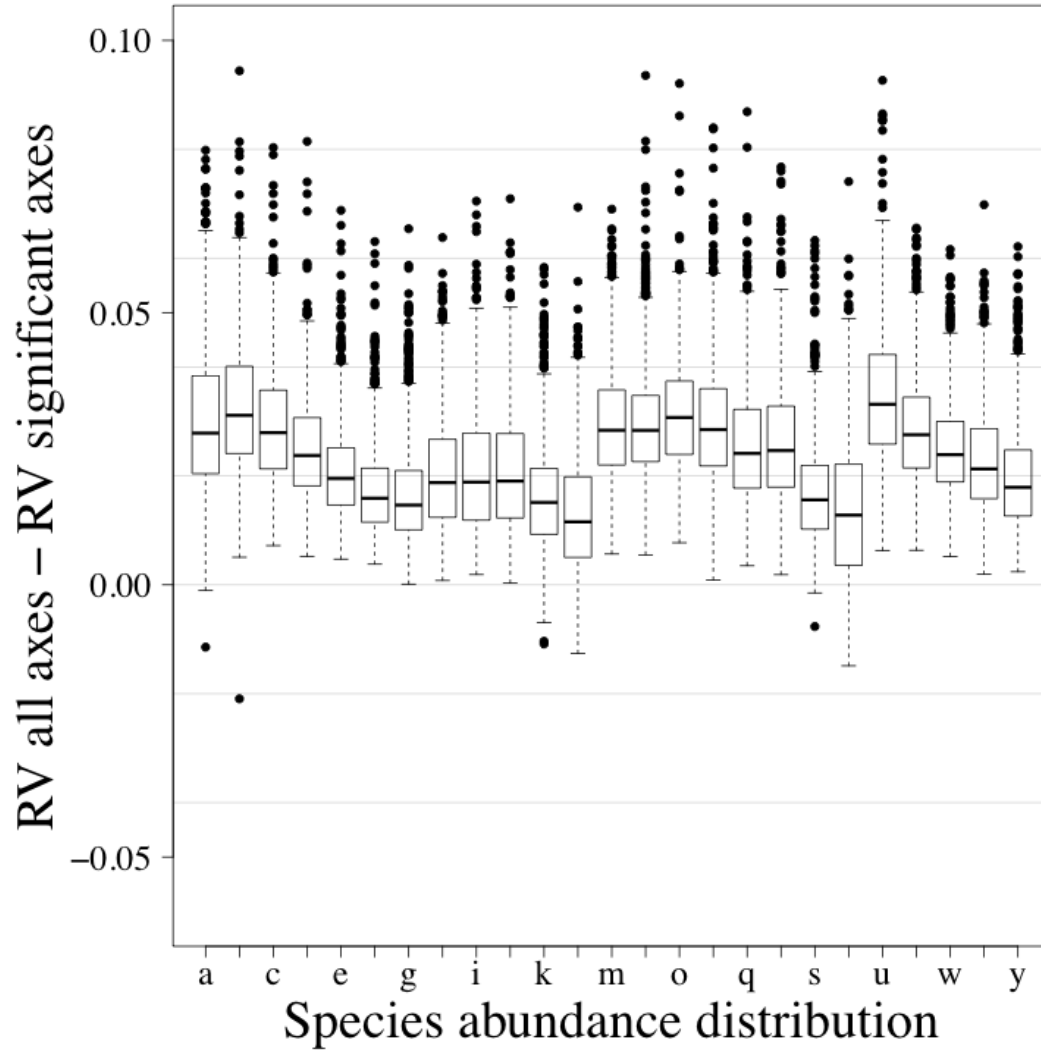


FIGURE C8. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from presence-absence data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 1). A thousand simulations were run for each SAD.

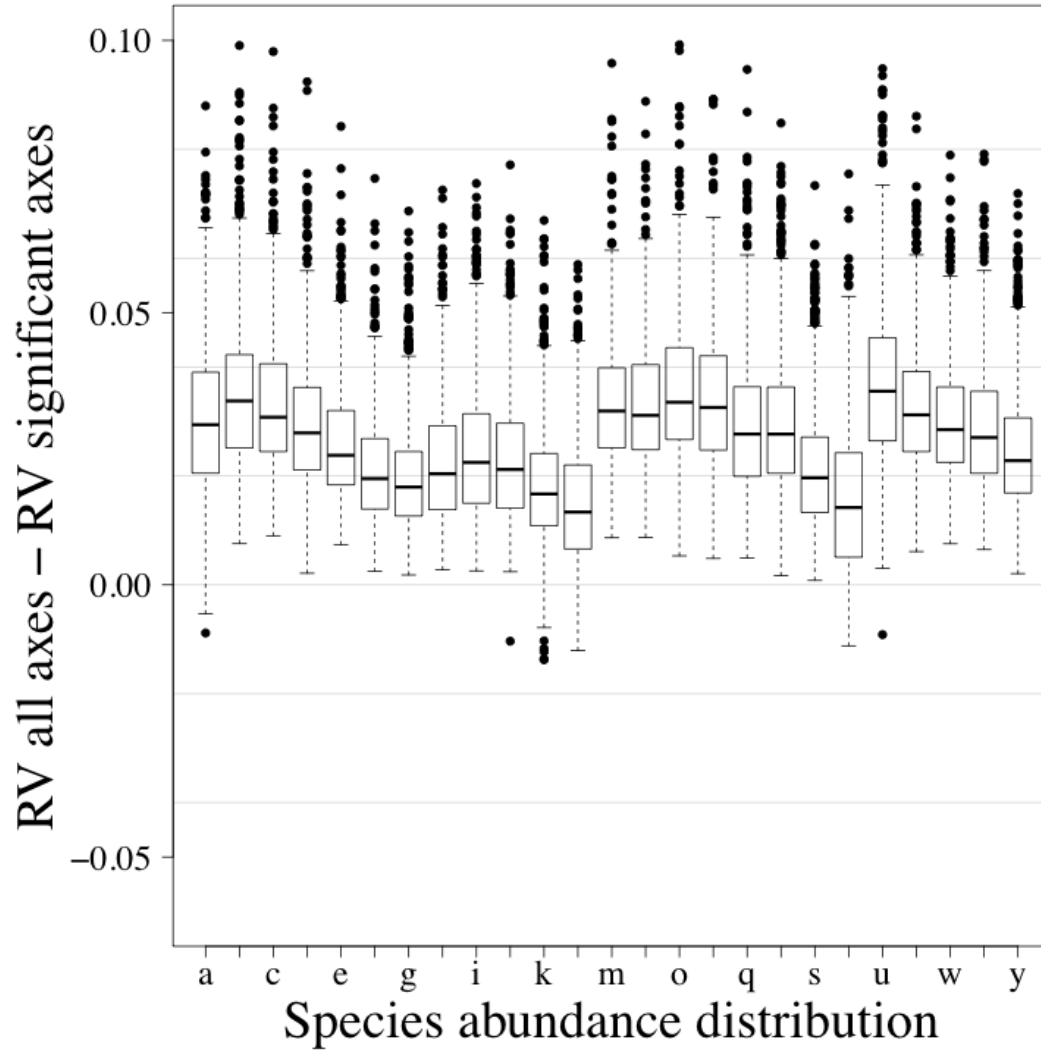


FIGURE C9. Comparison of consensus RDAs constructed using all canonical axes with consensus RDAs using only significant canonical axes. The \mathbf{Z}^* matrices calculated from presence-absence data were used in the comparison. Letters along the abscissa refer to the species abundance distribution (SAD) as presented in Figure 3.1. The ordinate presents the difference between RV coefficients calculated using all canonical axes and RV coefficients calculated using only the significant axes. The results are presented using boxplots. The upper and lower sections of the box define the first (25%) and third (75%) quartiles of the data, and the line in the middle of the box the median (50%). The lower whiskers describe the 1.5 interquartile range of the first quartile, the upper whisker stands for the 1.5 interquartile range of the third quartile, and the points indicate outliers. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 2). A thousand simulations were run for each SAD.

APPENDIX D

Ecological Archives EXXX-XXX-A4

COMPARISON OF CANONICAL ORDINATION MODELS FOR ABUNDANCE AND PRESENCE-ABSENCE DATA USING SIMULATIONS. FOUR FIGURES (FIGS. D1, D2, D3, AND D4)

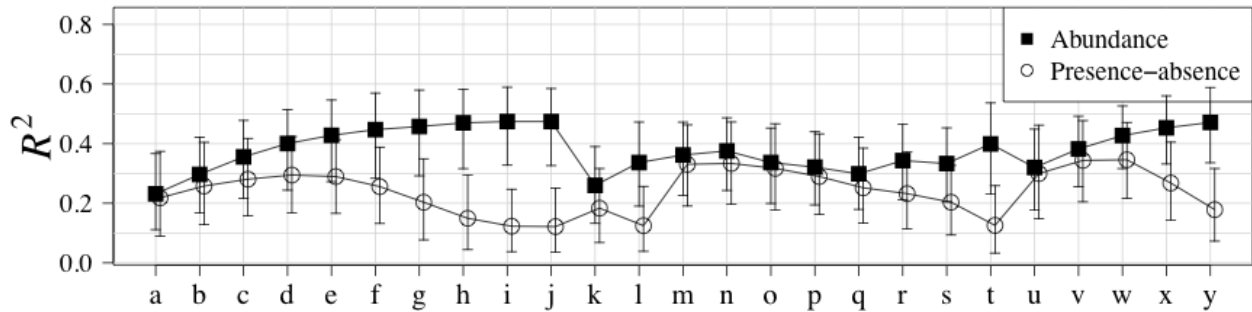


Fig. D1. RV coefficients (points) between canonical ordination model and the true species structure (equation 6 without the error term). For each data type (abundance and presence-absence), the significant canonical axes for all association coefficients (with the exception of the symmetric coefficient) were grouped. Error bars represent 95% confidence intervals. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.25). A thousand simulations were run for each SAD.

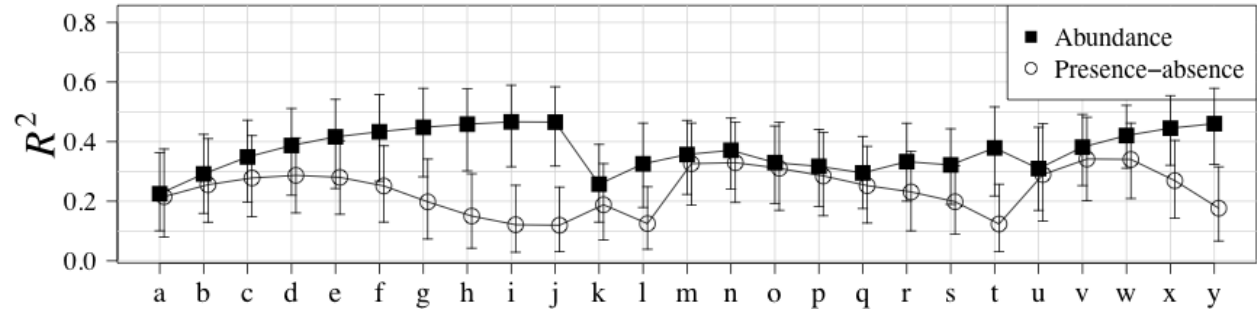


Fig. D2. RV coefficients (points) between canonical ordination model and the true species structure (equation 6 without the error term). For each data type (abundance and presence-absence), the significant canonical axes for all association coefficients (with the exception of the symmetric coefficient) were grouped. Error bars represent 95% confidence intervals. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 0.5). A thousand simulations were run for each SAD.

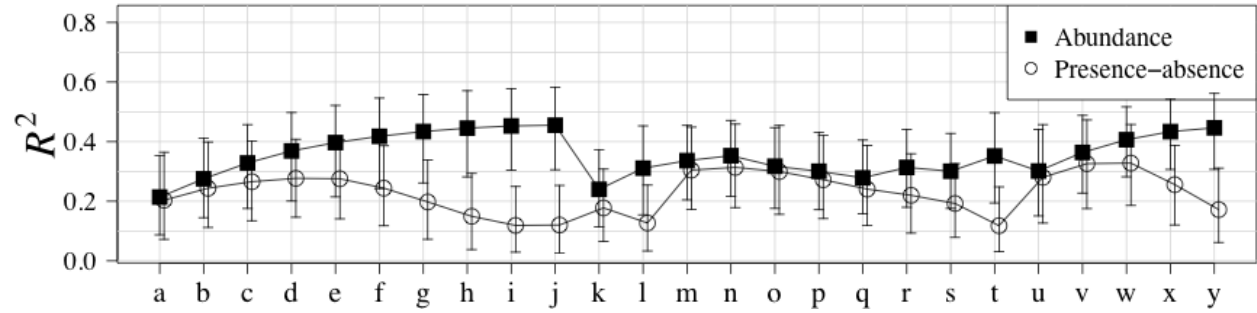


Fig. D3. RV coefficients (points) between canonical ordination model and the true species structure (equation 6 without the error term). For each data type (abundance and presence-absence), the significant canonical axes for all association coefficients (with the exception of the symmetric coefficient) were grouped. Error bars represent 95% confidence intervals. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 1). A thousand simulations were run for each SAD.

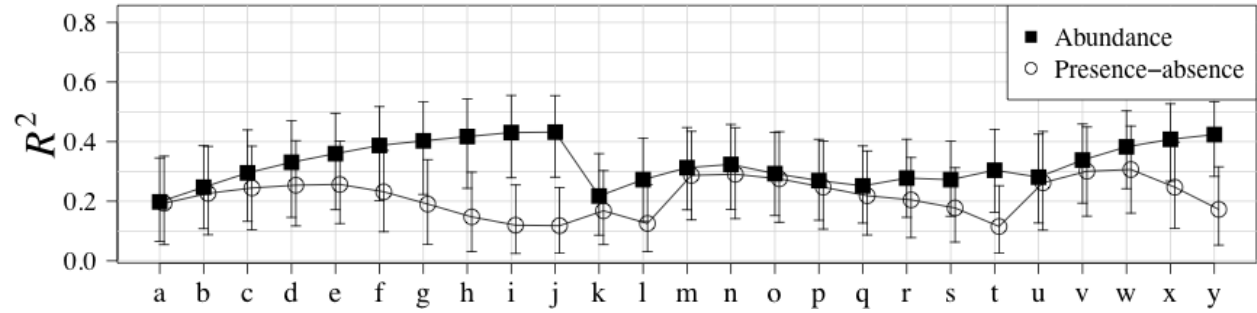


Fig. D4. RV coefficients (points) between canonical ordination model and the true species structure (equation 6 without the error term). For each data type (abundance and presence-absence), the significant canonical axes for all association coefficients (with the exception of the symmetric coefficient) were grouped. Error bars represent 95% confidence intervals. Letters on the x-axis refer to the species-abundance distribution (SAD) presented in Fig. 1. A line was drawn between each SAD of each association coefficient to ease comparisons between coefficients. Results are based on species simulated with an error term sampled from a Normal distribution (mean = 0, standard deviation = 2). A thousand simulations were run for each SAD.

APPENDIX E

Ecological Archives EXXX-XXX-A5

SPECIES CODE AND NAMES FOR CARABIDAE AND TREES SPECIES TWO TABLES (TABLES E1 AND E2)

Table D1: Species code and Latin name for Carabidae.

Code	Latin name
Agongrat	<i>Agonum gratiosum</i>
Agonplac	<i>Agonum placidum</i>
Agonretr	<i>Agonum retractum</i>
Agonsord	<i>Agonum sordens</i>
Agonsupe	<i>Agonum superioris</i>
Amarlitt	<i>Amara littoralis</i>
Amarluni	<i>Amara lunicollis</i>
Badiobtu	<i>Badister obtusus</i>
Bembgrap	<i>Bembidion grapii</i>
Bembrupi	<i>Bembidion rupicola</i>
Calaadve	<i>Calathus advena</i>
Calaingr	<i>Calathus ingratus</i>
Calofrig	<i>Calosoma frigidum</i>
Caracham	<i>Carabus chamissonis</i>
Dichcogn	<i>Dicheirotichus cognatus</i>
Elapamer	<i>Elaphrus americanus</i>
Elaplapp	<i>Elaphrus lapponicus</i>
Harpfulv	<i>Harpalus fulvilabris</i>
Loripili	<i>Loricera pilicornis</i>
Miscarct	<i>Miscodera arctica</i>
Nebrgyll	<i>Nebria gyllenhali</i>
Notibore	<i>Notiophilus borealis</i>
Notidire	<i>Notiophilus directus</i>
Patrfove	<i>Patrobus foveocollis</i>
Patrsept	<i>Patrobus septentrionis</i>
Platdece	<i>Platynus decentis</i>
Platmann	<i>Platynus mannerheimii</i>
Pteradst	<i>Pterostichus adstrictus</i>
Pterbrev	<i>Pterostichus brevicornis</i>
Pterpens	<i>Pterostichus pensylvanicus</i>
Pterpunc	<i>Pterostichus punctatissimus</i>
Pterripa	<i>Pterostichus riparius</i>
Seriquad	<i>Sericoda quadripunctata</i>
Sterhaem	<i>Stereocerus haematopus</i>
Synuimpu	<i>Synuchus impunctatus</i>
Trecapic	<i>Trechus apicalis</i>
Trecchal	<i>Trechus chalybeus</i>

Table D2: Species code, common and Latin name of trees species.

Code	Common name	Latin name
Pt	Aspen	<i>Populus tremuloides</i>
Bp	White birch	<i>Betula papyrifera</i>
Ab	Balsam fir	<i>Abie balsamea</i>
Ll	Tamarack	<i>Larix laricina</i>
Pb	Balsam poplar	<i>Populus balsamifera</i>
Pc	Lodgepole pine	<i>Pinus contorta</i>
Pm	Black spruce	<i>Picea mariana</i>
Pg	White spruce	<i>Picea glauca</i>