# Combined PDF for the Chinese and English versions of the paper

Legendre, P. 2007. Studying beta diversity: ecological variation partitioning by multiple regression and canonical analysis. *Journal of Plant Ecology (Chinese version)*, [formerly *Acta Phytoecologica Sinica*] 31: 976-981.

Legendre, P. 2008. Studying beta diversity: ecological variation partitioning by multiple regression and canonical analysis. *Journal of Plant Ecology (English version)* 1: 3-8. [Published by Oxford University Press on behalf of the Institute of Botany, Chinese Academy of Science.]

# β-

## Pierre Legendre

Département de sciences biologiques  Université de Montréal  C.P. 6128  succursale Centre-ville  Montréal  Québec  Canada H3C 3J7

β-

β-          β-

β-               β-

β-

# STUDYING BETA DIVERSITY  ECOLOGICAL VARIATION PARTITIONING BY MULTIPLE REGRESSION AND CANONICAL ANALYSIS

Pierre Legendre

*Départment de sciences biologiques  Université de Montréal  C. P. 6128  succursale Centre-ville  Montréal  Québec  Canada H3C 3J7*

**Abstract**  *Aims*   Beta diversity is the variation in species composition among sites in a geographic region. Beta diversity is a key concept for understanding the functioning of ecosystems  for the conservation of biodiversity  and for ecosystem management. This paper describes how to analyze it from community composition and associated environmental and spatial data tables.

*Methods*   Beta diversity can be studied by computing diversity indices for each site and testing hypotheses about the factors that may explain the variation among sites. Or  one can carry out a direct analysis of the community composition data table over the study sites  as a function of sets of environmental and spatial variables. These analyses are carried out by the statistical method of partitioning the variation of the diversity indices or the community composition data table with respect to environmental and spatial variables. Variation partitioning is briefly described in this paper.

*Important findings*   Variation partitioning is a method of choice for the interpretation of beta diversity using tables of environmental and spatial variables. Beta diversity is an interesting″currency″ for ecologists to compare either different sampling areas  or different ecological communities co-occurring in an area. Partitioning must be based upon unbiased estimates of the variation of the community composition data table that is explained by the various tables of explanatory variables. The adjusted coefficient of determination provides such an unbiased estimate in both multiple regression and canonical redundancy analysis. After partitioning  one can test the significance of the fractions of interest and plot maps of the fitted values corresponding to these fractions.

**Key words**  adjusted coefficient of determination  beta diversity  biodiversity  canonical redundancy analysis  community composition  variation partitioning

-

β-　　　　　　　　　　　　　　　　　　　Whit-
taker　1960　1972　Legendre *et al*.　2005

Shannon

**1**

$y \sim X \mid W$　　　　　$n$　　　　　$y$　　　　$X$ $m$
　　　　　　　　　　　　　$X$　　　　　　　　$W$ $q$
RDA　Rao　1964　　　　　　　　$\mid W$　　　　$W$ $X$
CCA　ter Braak　1986　1987a　　　　1 $X$　　$W$　　　　　　$X$
1987b　　　　　　　　　　　　　Legen-　$X_{res\ W}$　2 $y$　　$X_{res\ W}$　　　　　　　　$R^2$
dre & Legendre　1998

Canoco　　　　ter Braak & Smilauer　2002　　　　　　　　　　　　　　$R^2$　　　　　　　　F
R　　　　　　R Development Core Team　2007　　　　　$R^2$
′vegan′　Oksanen *et al*.　2007　　　　　SS　　　　　　　　　　　　　　　　+

$R^2_{y \sim X \mid W} = $　SS　$y \sim X \mid W$　　　　/ SS　　　　　+
SS　　　　　　　　　　　　　　　　　　　　　　1
1　　　　　　　　$R^2_{y \sim X \mid W} = $　a / a +
d



図1　　　　　　　　　　　y　　　　　　　Y　　　　　　　　X　W　Venn　Legendre　1993
Fig. 1　Venn diagram representing the partition of the variation of a response variable Y or a response matrix Y
between two sets of explanatory variables X and W　Legendre　1993
100% y　Y　　　　　b　　　　　　X　　　　　　　W　　　　　　　　　　　　Legendre 1993 The rectan-
gle represents 100% of the variation in y or Y. Fraction　b　is the intersection　not the interaction　of the amounts of variation explained by linear models of X
and W

*F*

$q$　　　　　　　　　　　　　　$q = 0$　　　　　　　　$F = $　SS　$y \sim X \mid W$　　　　/m /　SS　　　/
$R^2$ $F$　　　　　　　　　　　　　　　　　$n - 1 - m - q$　　　　　　　　　　　3
$F = $　$R^2_{y \sim X \mid W}/m$ /　1 $- R^2_{y \sim X \mid W}$ / $n - 1 - m$　　　1
$- q$　　　　　　　　　　　　　　　2　　　　　$F = $　a /m /　d / $n - 1 - m - q$　　　4

F          F

Permutation test

Manly

1997   Legendre   Legendre   1998

$y \sim X | W$    $y \sim W | X$

## 2

$Y \sim X | W$                    $Y$   $n$

$\times p$      $m$                        X

RDA   X                        W   $q$

| W              W   X

F

$Y \sim X | W$    $Y \sim W$

| X

## 3

$R^2$

$$R^2 = \frac{\text{regression SS}}{\text{total SS}} = \frac{\sum \ \hat{y}_i - \bar{y}\ ^2}{\sum \ y_i - \bar{y}\ ^2} = 1 - \frac{\text{residual SS}}{\text{total SS}}$$

5

"  regression SS"

$R^2$          y

$R_a^2$ Ezekiel   1930

$$R_a^2 = 1 - \frac{\text{residual mean square}}{\text{total mean square}} = 1 - \ 1 - R^2$$

$$\frac{\text{total df}}{\text{residual df}}$$

6

$R_a^2$                    5

F                            $n - 1$

$n - m - 1$   $n$                    $m$

0  F                            $n$                    $n - m$

$R^2$                                        $n$

$m$      $m$          $n$

$R_a^2$

Ohtani   2000        $R_a^2$

X          y

$R_a^2$

RDA                    $R^2$

Bimultivariate redundancy statistic    Miller

& Farr   1971

Y

F

Permutation

$m \times p$      $p$   $n - m - 1$      Fisher-Snedecor $F$

F                    Miller   1975   Peres-Neto

2006

Hellinger                            RDA

6                    $R^2$

$R_a^2$                    X                    Y

Hellinger          5

PCA

RDA          Legendre & Gallagher   2001

$R_a^2$                                        $R_a^2$

## 4

β-                                        SS   Y

Bell   2001   Hubbell   2001   He   2005

Legendre   2005

Borcard   1992   Borcard   Legendre   1994

Y

y    1.

Y                          X   W        1              3                        y

3        3                          Y

**Table 1**  Method for calculating the adjusted fractions of variation a to d depicted in Fig. 1  Three multiple regressions or canonical analyses are required

| Canonical analyses | $R^2$ Compute $R^2$ eq. 5 | $R_a^2$ Compute $R_a^2$ eq. 6 and fractions of variation | Tested for significance |
|---|---|---|---|
| Y ~ X | $R^2$ of Y ~ X | a + b = $R_a^2$ of Y ~ X | Yes |
| Y ~ W | $R^2$ of Y ~ W | b + c = $R_a^2$ of Y ~ W | Yes |
| Y ~ X W | $R^2$ of Y ~ X W | a + b + c = $R_a^2$ of Y ~ X W | Yes |
| | | a = a + b − b | Yes |
| | | b = a + b + b + c − a + b + c | No |
| | | c = b + c − b | Yes |
| | | Residuals = d = 1 − a + b + c | No |

1  Y ~ X                          5    6

$R^2$   $R_a^2$          1                    1 $R_a^2$

a        b

2  Y ~ W                $R^2$   $R_a^2$   $R_a^2$           Borcard   1992   Borcard   Legendre  1994

1                          b    c                                    W

3  Y ~ X W                    $R^2$   $R_a^2$

$R_a^2$        1                                       Borcard   Legendre  2002   Borcard   2004

a   b      c                                           PCNM                       W

4          b   b = a + b + b + c − a +                                            PCNM

b + c

5          a   a = a + b − b                                        "

6          c   c = b + c − b                     "  DBEM   Moran           MEM   Dray *et*

7                      d   d = 1 −              *al.* 2006

a + b + c                                       2.

1        Venn

0                       Permutation              a + b   b + c      a +

X              W                       b + c   1    3                      *F*

1      Venn                                                          a    c

1      a + b       Y                                   1              b

b    X   W                         d

a                                    *F*

2      b + c       Y                             a    c                       Y ~

W   X | W   Y ~ W | X                2 ~ 4

c          W            *F*                *F*

—— Legendre

X                                   & Legendre  1998  Anderson & Legendre  1999

X                                               3.

X   Y                                Y              a + b   b + c   a + b + c   a   c

35

y                                                                                    Borcard    1992

Y          Borcard    Legendre    1994

Borcard    Legendre    2002        Borcard    2004
PCNM

a                                                    β-

Bubble plots                                β-

kriging                                                                              1

Y                                3

4                                                                                    5    24 hm$^2$

R                    ' vegan'                              Oksanen *et*

*al*.    2007                                              20 m × 20 m  40 m × 40 m

SS  Y

**5**                                                                    β-

β-

5

SS  Y    β-                                        2

20        70              β-                            β-

Principal  component                                                        β-

analysis  PCA                  Correspondence analysis  CA

Principal coordinate analysis  PCoA

20                              β-

80    90                            Canoco

ter Braak    1988    ter

Braak & Smilauer    2002

I

Power

"                                    "                              $R_a^2$

-

[1]    Legendre    1990

β-

β-
*ISI Web of*

*Knowledge* of the *Institute for Scientific Information*

603

70

---

Anderson MJ, Legendre P (1999). An empirical comparison of per-
mutation methods for tests of partial regression coefficients in a
linear model. *Journal of Statistical Computation and Simulation*,
62, 271 – 303.

Bell G (2001). Neutral macroecology. *Science*, 293, 2413 –
2418.

Borcard D, Legendre P (1994). Environmental control and spatial
structure in ecological communities: an example using oribatid
mites (Acari, Oribatei). *Environmental and Ecological Statis-
tics*, 1, 37 – 53.

Borcard D, Legendre P (2002). All-scale spatial analysis of eco-
logical data by means of principal coordinates of neighbour matri-
ces. *Ecological Modelling*, 153, 51 – 68.

Borcard D, Legendre P, Avois-Jacquet C, Tuomisto H (2004).
Dissecting the spatial structure of ecological data at multiple
scales. *Ecology*, 85, 1826 – 1832.

Borcard D, Legendre P, Drapeau P (1992). Partialling out the
spatial component of ecological variation. *Ecology*, 73, 1045 –
1055.

Dray S, Legendre P, Peres-Neto PR (2006). Spatial modelling: a
comprehensive framework for principal coordinate analysis of
neighbour matrices (PCNM). *Ecological Modelling*, 196, 483 –
493.

Ezekiel M (1930). *Methods of Correlation Analysis*. John Wiley
and Sons, New York.

He F (2005). Deriving a neutral model of species abundance from
fundamental mechanisms of population dynamics. *Functional E-
cology*, 19, 187 – 193.

Hubbell SP (2001). *The Unified Neutral Theory of Biodiversity and
Biogeography*. Princeton University Press, Princeton, New Jer-
sey.

Legendre P (1990). Quantitative methods and biogeographic analy-
sis. In: Garbary DJ, South RG eds. *Evolutionary Biogeography
of the Marine Algae of the North Atlantic*. NATO ASI Series, Vol.
G 22. Springer-Verlag, Berlin, 9 – 34.

Legendre P (1993). Spatial autocorrelation: trouble or new
paradigm? *Ecology*, 74, 1659 – 1673.

Legendre P, Borcard D, Peres-Neto PR (2005). Analyzing beta di-
versity: partitioning the spatial variation of community composition
data. *Ecological Monographs*, 75, 435 – 450.

Legendre P, Gallagher ED (2001). Ecologically meaningful trans-
formations for ordination of species data. *Oecologia*, 129, 271 –
280.

Legendre P, Legendre L (1998). *Numerical Ecology* 2nd English
edn. Elsevier Science BV, Amsterdam.

Manly BJF (1997). *Randomization, Bootstrap and Monte Carlo
methods in biology* 2nd edn. Chapman and Hall, London.

Miller JK (1975). The sampling distribution and a test for the sig-
nificance of the bimultivariate redundancy statistic: a Monte Carlo
study. *Multivariate Behavioral Research*, 10, 233 – 244.

Miller JK, Farr SD (1971). Bimultivariate redundancy: a compre-
hensive measure of interbattery relationship. *Multivariate Behav-
ioral Research*, 6, 313 – 324.

Ohtani K (2000). Bootstrapping $R^2$ and adjusted $R^2$ in regression
analysis. *Economic Modelling*, 17, 473 – 483.

Oksanen J, Kindt R, Legendre P, O'Hara RB (2007). Vegan:
community ecology package version 1.8 – 5. URL http://cran.
r-project.org/.

Peres-Neto PR, Legendre P, Dray S, Borcard D (2006). Variation
partitioning of species data matrices: estimation and comparison
of fractions. *Ecology*, 87, 2614 – 2625.

Rao CR (1964). The use and interpretation of principal component
analysis in applied research. *Sankhyaá, Series A, Indian Journal
of Statistics*, 26, 329 – 358.

R Development Core Team (2007). *R: A Language and Environ-
ment for Statistical Computing*. R Foundation for Statistical Com-
puting, Vienna, Austria. URL http://www.R-project.org.

ter Braak CJF (1986). Canonical correspondence analysis: a new
eigenvector technique for multivariate direct gradient analysis. *E-
cology*, 67, 1167 – 1179.

ter Braak CJF (1987a). The analysis of vegetation-environment re-
lationships by canonical correspondence analysis. *Vegetatio*, 69,
69 – 77.

ter Braak CJF (1987b). Ordination. In: Jongman RHG, ter Braak
CJF, van Tongeren OFR eds. *Data Analysis in Community and
Landscape Ecology*. Pudoc, Wageningen (reissued in 1995 by
Cambridge University Press, Cambridge), 91 – 173.

ter Braak CJF (1988). *CANOCO — a FORTRAN program for
Canonical Community Ordination by [ partial ] [ detrended ]
[ canonical ] Correspondence Analysis, Principal Component Anal-
ysis and Redundancy Analysis (version 2.1)*. Agricultural Mathe-
matics Group, Ministry of Agriculture and Fisheries, Wagenin-
gen.

ter Braak CJF, Smilauer P (2002). *Canoco Reference Manual and
CanoDraw for Windows User's Guide: Software for Canonical
Community Ordination (version 4.5)*. Microcomputer Power,
Ithaca, New York.

Whittaker RH (1960). Vegetation of the Siskiyou mountains, Ore-
gon and California. *Ecological Monographs*, 30, 279 – 338.

Whittaker RH (1972). Evolution and measurement of species diver-
sity. *Taxon*, 21, 213 – 251.

# Studying beta diversity: ecological variation partitioning by multiple regression and canonical analysis

*Pierre Legendre*

*Département de sciences biologiques, Université de Montréal, CP 6128, succursale Centre-ville, Montréal, Québec, Canada H3C 3J7*
E-mail: Pierre.Legendre@umontreal.ca

# Abstract

## Aims

Beta diversity is the variation in species composition among sites in a geographic region. Beta diversity is a key concept for understanding the functioning of ecosystems, for the conservation of biodiversity and for ecosystem management. The present report describes how to analyse beta diversity from community composition and associated environmental and spatial data tables.

## Methods

Beta diversity can be studied by computing diversity indices for each site and testing hypotheses about the factors that may explain the variation among sites. Alternatively, one can carry out a direct analysis of the community composition data table over the study sites, as a function of sets of environmental and spatial variables. These analyses are carried out by the statistical method of partitioning the variation of the diversity indices or the community composition data table with respect to environmental and spatial variables. Variation partitioning is briefly described herein.

## Important findings

Variation partitioning is a method of choice for the interpretation of beta diversity using tables of environmental and spatial variables. Beta diversity is an interesting 'currency' for ecologists to compare either different sampling areas or different ecological communities co-occurring in an area. Partitioning must be based upon unbiased estimates of the variation of the community composition data table that is explained by the various tables of explanatory variables. The adjusted coefficient of determination provides such an unbiased estimate in both multiple regression and canonical redundancy analysis. After partitioning, one can test the significance of the fractions of interest and plot maps of the fitted values corresponding to these fractions.

**Keywords:** Adjusted coefficient of determination • beta diversity • biodiversity • canonical redundancy analysis • community composition • variation partitioning

Received: 4 August 2006 Accepted: 4 August 2006

# Introduction

Ecologists collect community composition data (species presence–absence or abundance data) at several sites in a region of interest in order to analyse and interpret beta diversity, which is the variation in species composition among the sites (Whittaker, 1960, 1972; Legendre et al., 2005). Analysis of a synthetic descriptor such as species richness or Shannon diversity can be done by multiple regression, whereas the analysis of whole community composition data tables is carried out by canonical analysis. Results from these two types of analyses are not equivalent: analysis of the whole community composition data produces results that are much more informative since they provide information about the reactions of individual species to the environmental and spatial variables. The asymmetrical forms of canonical analysis used for this type of research are canonical redundancy analysis (RDA; Rao,

1964) and canonical correspondence analysis (CCA; ter Braak, 1986, 1987a, b). These analyses are described in several textbooks, including Legendre and Legendre (1998). They are implemented in computer packages such as Canoco (ter Braak and Smilauer, 2002) and the 'vegan' library (Oksanen et al., 2007) of the R statistical language (R Development Core Team, 2007).

Variation in species composition among sites is studied by canonical analysis of the species composition data as a function of different types of environmental variables: water or soil chemistry, geology, geomorphology, environmental impact descriptors, and so on. The study of spatial structures involves spatial variables derived from the geographic coordinates of the sampling sites, described below. Variation partitioning is a technique of choice for this type of analysis. In all cases, statistics are used to describe how successful the explanatory variables are at explaining the response variables (community composition

data). The choice of an appropriate, unbiased statistical estimator is of great importance for the correct interpretation of the results. This report will briefly describe partial linear regression and canonical analysis, the simple and adjusted forms of the coefficient of determination used in regression and canonical analysis, and finally variation partitioning.

## Partial linear regression

The notation $\mathbf{y} \sim \mathbf{X}|\mathbf{W}$ represents the partial linear regression of a response variable $\mathbf{y}$ (vector of length $n$) on a matrix $\mathbf{X}$ containing $m$ explanatory variables, while controlling for the linear effect of a matrix $\mathbf{W}$ containing $q$ covariables. Partial regression is computed in two steps: (i) regress $\mathbf{X}$ on $\mathbf{W}$ and compute the residuals $\mathbf{X}_{\text{res}(\mathbf{W})}$; and (ii) regress $\mathbf{y}$ on $\mathbf{X}_{\text{res}(\mathbf{W})}$ to obtain the partial $R^2$, the fitted values, the residuals, and so on.

The $R^2$ statistic of a partial regression that will be used to construct the $F$-statistic for the test of significance (next paragraph) is called the partial $R^2$. It is the ratio of the sum-of-squares (SS) of the fitted values of the partial regression on the sum (SS of the fitted values + SS of the residuals):

$$R^2_{\mathbf{y} \sim \mathbf{X}|\mathbf{W}} = SS(\text{fitted values of } \mathbf{y} \sim \mathbf{X}|\mathbf{W})/(SS(\text{fitted values}) + SS(\text{residuals})) \quad (1)$$

Using the graphical representation of Fig. 1, $R^2_{\mathbf{y} \sim \mathbf{X}|\mathbf{W}} = $ [a]/[a + d].

The $F$-statistic used to test the significance of the partial regression relationship takes into account the number of covariables $q$; in ordinary multiple regression, $q = 0$. The $F$-statistic is computed as follows using the partial $R^2$:

$$F = (R^2_{\mathbf{y} \sim \mathbf{X}|\mathbf{W}}/m)/((1 - R^2_{\mathbf{y} \sim \mathbf{X}|\mathbf{W}})/(n - 1 - m - q)) \quad (2)$$

It can also be computed directly from the sums-of-squares:

$$F = (SS(\text{fitted values of } \mathbf{y} \sim \mathbf{X}|\mathbf{W})/m)/((SS(\text{residuals}))/(n - 1 - m - q)) \quad (3)$$

or, using Fig. 1:

$$F = ([a]/m)/([d]/(n - 1 - m - q)) \quad (4)$$

Significance of the $F$-statistic can be tested with reference to an $F$-distribution if the condition of normality of the residuals is met (this is rarely the case for ecological data), or by a permutation test if it is not (this is the most common case). Permutation tests are described in several textbooks, including Manly (1997) and Legendre and Legendre (1998). In the application to variation partitioning described below, both $\mathbf{y} \sim \mathbf{X}|\mathbf{W}$ and $\mathbf{y} \sim \mathbf{W}|\mathbf{X}$ will be computed and tested for significance.

## Partial canonical analysis

Similarly, the notation $\mathbf{Y} \sim \mathbf{X}|\mathbf{W}$ represents the partial canonical redundancy analysis (partial RDA) of a response data ma-
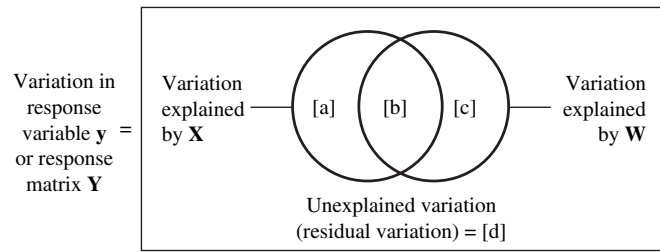


**Figure 1** Venn diagram representing the partition of the variation of a response variable $\mathbf{y}$ or a response matrix $\mathbf{Y}$ between two sets of explanatory variables $\mathbf{X}$ and $\mathbf{W}$. The rectangle represents 100% of the variation in $\mathbf{y}$ or $\mathbf{Y}$. Fraction [b] is the intersection (not the interaction) of the amounts of variation explained by linear models of $\mathbf{X}$ and $\mathbf{W}$. Adapted from Legendre (1993).

trix $\mathbf{Y}$ of size ($n \times p$) on a matrix $\mathbf{X}$ containing $m$ explanatory variables, while controlling for the linear effect of a matrix $\mathbf{W}$ containing $q$ covariables. Partial canonical analysis is computed in the same way as partial linear regression and uses the same $F$-statistic for significance testing (see below for details). In the application to variation partitioning described below, both $\mathbf{Y} \sim \mathbf{X}|\mathbf{W}$ and $\mathbf{Y} \sim \mathbf{W}|\mathbf{X}$ will be computed and tested for significance.

## Unadjusted and adjusted coefficients of determination

The coefficient of multiple determination (unadjusted $R^2$) estimates the forecasting potential of a multiple regression equation:

$$R^2 = \frac{\text{regression SS}}{\text{total SS}} = \frac{\Sigma(\hat{y}_i - \bar{y})^2}{\Sigma(y_i - \bar{y})^2} = 1 - \frac{\text{residual SS}}{\text{total SS}} \quad (5)$$

where 'regression SS' is the sum-of-squares of the fitted values of the regression equation. It measures the proportion of the variation of $\mathbf{y}$ about its mean that is explained by the regression equation.

In multiple regression, an alternative measure of determination is the adjusted coefficient of multiple determination $R^2_a$ (Ezekiel, 1930):

$$R^2_a = 1 - \frac{\text{residual mean square}}{\text{total mean square}} = 1 - (1 - R^2)\left(\frac{\text{total d.f.}}{\text{residual d.f.}}\right) \quad (6)$$

The right-hand parentheses of equation 6 shows that $R^2_a$ takes into account the numbers of degrees of freedom associated with the numerator and denominator of equation 5. In ordinary multiple regression, the total degrees of freedom of the $F$-statistic are $(n - 1)$ and the degrees of freedom of the residuals are $(n - m - 1)$ where $n$ is the number of observations and $m$ is the number of explanatory variables in the model. In multiple regression through the origin, where the intercept is forced to zero, the total degrees of freedom of the $F$-statistic are

$n$ and the residual degrees of freedom are $(n - m)$. In both cases, the correction takes into account the number of objects $n$ and the number of explanatory variables $m$; the correction is light when $m$ is small when compared with $n$. $R_a^2$ is a suitable measure of goodness-of-fit for comparing regression equations fitted to different data sets, with different numbers of objects and explanatory variables. Using simulated data with normal error, Ohtani (2000) has shown that $R_a^2$ is an unbiased estimator of the contribution of a set of explanatory variables **X** to the explanation of **y**. The $R_a^2$ statistic cannot be directly computed for partial linear regression because the number of degrees of freedom to use in the correction is then unknown.

In RDA, the canonical $R^2$ is called the bimultivariate redundancy statistic (Miller and Farr, 1971) or the canonical coefficient of determination. It is computed in the same way as in multiple regression: it is the ratio of the sum of each response variable's regression (or fitted values) SS to the sum of all response variables' total SS. In canonical analysis, the significance of the $F$-statistic is always tested by permutation, except in the very restrictive case where the variables in **Y** are standardized and the residuals are multinormal. These conditions are almost never met with ecological data; in the rare cases where they are, the $F$-statistic is tested using the Fisher–Snedecor $F$-distribution with $(m \times p)$ and $p(n - m - 1)$ degrees of freedom (Miller, 1975). Using numerical simulations, Peres-Neto et al. (2006) have shown that, for normally distributed data or Hellinger-transformed species abundances in RDA, the adjusted bimultivariate redundancy statistic $R_a^2$, obtained by applying equation 6 to the canonical $R^2$, produced unbiased estimates of the real contributions of the variables in **X** to the explanation of a response matrix **Y**. The Hellinger transformation is one of five transformations that make community composition data containing many zeros suitable for analysis by linear methods such as principal component analysis (PCA) or RDA (Legendre and Gallagher, 2001).

Adjusted coefficients of determination in multiple regression and canonical analysis can, on occasion, take negative values. For large data sets, $R_a^2$ is zero when the explanatory variables explain no more variation than random normal variables would. Negative values of $R_a^2$ are interpreted as zeros; they correspond to cases where the explanatory variables explain less variation than random normal variables would.

# Variation partitioning

The technique of variation partitioning is used when two or more complementary sets of hypotheses can be invoked to explain the variation of an ecological response variable. For example, the abundance of a species could vary as a function of biotic and abiotic factors. In the study of beta diversity, the total variation of the community composition data table, denoted SS(**Y**), can be partitioned among one or more sets

of environmental variables and a table describing the spatial relationships among the sampling sites. Fitting the community composition data to spatial variables, as described below, allows researchers to establish that there are significant spatial patterns, perhaps at various scales, present in the species data. The presence of significant spatial patterns in the response data can be invoked as support either for a neutral model (Bell, 2001, Hubbell, 2001, He, 2005) or for environmental control since environmental data are often spatially structured. The presence of significant relationships between the species and environmental variables would strongly support the hypothesis of environmental control, which is not in opposition to a hypothesis of neutral process, as discussed by Legendre et al. (2005).

Variation partitioning among environmental and spatial components was first described by Borcard et al. (1992) and Borcard and Legendre (1994). Variation partitioning will be presented in the context of the analysis of a response community composition data table **Y**. It can also be applied to a single response variable **y** since the algebra of partial linear regression is the same as that of partial canonical analysis.

Variation partitioning of a response data table **Y** with respect to two matrices of explanatory variables **X** and **W** involves the following three steps, which correspond to different research objectives.

## Obtaining the Fractions of Variation

The calculations, based upon three multiple regressions (for a single variable **y**) or three canonical analyses (for a multivariate response table **Y**), are summarized in Table 1.

(i) Compute the canonical analysis of **Y** with respect to the first table of explanatory variables **X**. Compute the $R^2$ and $R_a^2$ using equations 5 and 6. Assuming that the rectangle has a surface area normalized to 1, the $R_a^2$ corresponds to the surface area of the left-hand circle in Fig. 1. It contains the adjusted fractions [a] and [b].

(ii) Compute the canonical analysis of **Y** with respect to the second table of explanatory variables **W**. Compute the $R^2$ and $R_a^2$ using equations 5 and 6. The $R_a^2$ corresponds to the surface area of the right-hand circle in Fig. 1. It contains the adjusted fractions [b] and [c].

(iii) Compute the canonical analysis of **Y** with respect to the union of tables **X** and **W**. Compute the $R^2$ and $R_a^2$ using equations 5 and 6. The $R_a^2$ corresponds to the union of the two circles in Fig. 1. It contains the adjusted fractions [a], [b] and [c].

(iv) From these first results, compute fraction [b] by subtraction: [b] = [a + b] + [b + c] − [a + b + c].

(v) Compute fraction [a] by subtraction: [a] = [a + b] − [b].

(vi) Compute fraction [c] by subtraction: [c] = [b + c] − [b].

(vii) Compute fraction [d], which represents the residual variation, by subtraction: [d] = 1 − [a + b + c].

These values can be added to a Venn diagram such as the one shown in Fig. 1. Because they are based on adjusted

**Table 1** Method for calculating the adjusted fractions of variation [a] to [d] depicted in Fig. 1

| Canonical analyses | Compute $R^2$ (eq. 5) | Compute $R_a^2$ (eq. 6) and fractions of variation | Can be tested for significance |
|---|---|---|---|
| **Y~X** | $R^2$ of **Y~X** | [a + b] = $R_a^2$ of **Y~X** | Yes |
| **Y~W** | $R^2$ of **Y~W** | [b + c] = $R_a^2$ of **Y~W** | Yes |
| **Y~(X,W)** | $R^2$ of **Y~(X,W)** | [a + b + c] = $R_a^2$ of **Y~(X,W)** | Yes |
| | | [a] = [a + b] − [b] | Yes |
| | | [b] = [a + b] + [b + c] − [a + b + c] | No |
| | | [c] = [b + c] − [b] | Yes |
| | | Residuals = [d] = 1 − [a + b + c] | No |

Three multiple regressions or canonical analyses are required.

coefficients of determination, the fractions can, on occasion, take negative values. These are interpreted as zeros, as explained in the previous section.

When **X** is a matrix of environmental variables and **W** contains descriptors of the spatial relationships among the sampling sites, the Venn diagram (Fig. 1) provides the following information:

(i) The circle containing [a + b] shows how much of the variation of **Y** is explained by the environmental variables. Of that, [b] is the variation explained jointly by **X** and **W**, or the fraction of the environmentally explained variation that is spatially structured. [a] is the environmentally explained variation that is not explained by the spatial variables found in **W**.

(ii) The circle containing [b + c] shows how much of the variation of **Y** is explained by the spatial variables found in **W**. Of that, [c] is the variation explained uniquely by a linear model of the spatial variables found in **W** and not by a linear effect of the environmental variables **X**. This component may be due to spatially structured environmental variables that are not present in table **X** or to non-linear effects of the environmental variables **X** on **Y**. That variation may also be due to processes, such as competition or dispersal, in the ecological community depicted by table **Y**. In that case, it cannot be related to environmental variables.

To model broad-scale spatial patterns only, Borcard et al. (1992) and Borcard and Legendre (1994) used a third-degree polynomial function of the geographic coordinates of the sampling sites as matrix **W** in variation partitioning. More recently, Borcard and Legendre (2002) and Borcard et al. (2004) described PCNM (principal coordinate analysis of neighbour matrices) analysis, which generates a matrix **W** containing spatial descriptors that represent a spectral decomposition of the spatial relationships among the sampling sites. PCNM analysis allows researchers to model these relationships at all spatial scales. PCNM geographic functions are a type of 'distance-based eigenvector maps' (DBEMs), which belong to a general class called 'Moran's eigenvector maps' (MEMs) (Dray et al., 2006).

## Testing the Significance of the Fractions

The fractions must be tested for significance in order to support fully the reasoning described in the first paragraph of this section. The $F$-statistics of the three regressions or canonical analyses giving rise to the adjusted fractions [a + b], [b + c] and [a + b + c] (Table 1) can be tested directly by parametric or permutation tests. Individual fractions [a] and [c] cannot be tested in that way (see below), while fraction [b] cannot be tested at all, as shown in Table 1. [d] is the residual variation. Fraction [d], together with its degrees of freedom, forms the denominator of the $F$-statistics used in testing the other fractions.

The partial canonical analyses **Y~X|W** and **Y~W|X** have to be computed to test the significance of fractions [a] and [c], respectively. The $F$-statistics are computed following equation 2, 3 or 4. These $F$-statistics are tested using special permutation methods, called 'permutation of the residuals', described in Legendre and Legendre (1998) and Anderson and Legendre (1999).

## Mapping the Fitted Values of the Fractions

The fitted values corresponding to fractions [a + b], [b + c], [a + b + c], [a] and [c] can be computed in order to draw maps that will help in interpreting them. In the case of a single response variable **y**, the fitted values of the multiple and partial multiple regressions giving rise to these fractions provide the values that can be mapped. In the case of a multivariate response table **Y**, e.g. a community composition table, the fitted values are contained in multivariate tables of site scores produced by the canonical and partial canonical analyses. The first few axes of each of these tables, which correspond to the largest canonical eigenvalues, can be used for mapping. Point maps, such as bubble plots, should be produced for fraction [a] because that fraction is not spatially structured; the map will display the 'local innovation' at each sampling site. Interpolation mapping techniques, such as kriging, can be used for the other fractions, which contain spatially correlated values.

Variation partitioning of **Y** can be computed with respect to three or four tables of explanatory variables. The algebra,

which involves more steps, will not be explained in detail here. It is described in one of the documentation files of the package 'vegan' (Oksanen et al., 2007) of the R statistical language.

## Discussion

Analysis of the variation of a community composition data table is a widely used approach in community ecology. As stated in the Introduction, the total variation in a community composition table, denoted SS($\mathbf{Y}$), is a measure of beta diversity, which is the diversity among sites in the study area. Ordination methods such as PCA, correspondence analysis (CA) and principal coordinate analysis (PCoA) have been used since the 1970s to partition the variation of community composition data tables into orthogonal axes, which can be used to produce ordination plots or can be related to potentially explanatory variables. In the years 1980 and 1990, canonical ordination methods were made widely available to ecologists, firstly through the program Canoco (ter Braak, 1988; ter Braak and Smilauer, 2002). Canonical ordination offers the possibility of directly incorporating the environmental variables of interest in the analysis as constraints for the ordination, hence the expression 'constrained ordination methods'. Ecologists quickly took advantage of this improved methodology and applied it to all problems of species–environment relationships. (Two bibliographies on the applications of canonical analysis to ecology, covering together the period 1986 to 1996, contain a total of 804 entries. They are available from H. J. B. Birks, Botanical Institute, University of Bergen, Allégaten 41, N-5007 Bergen, Norway, and also on the URL http:/www.bio.umontreal.ca/casgrain/cca_bib/.) In 1990, Legendre proposed to use canonical analysis to model the spatial structure of community composition data, representing the spatial relationships among the sampling sites by a polynomial function of their geographic coordinates. That development led to the method of variation partitioning among environmental and spatial components, described in the previous section.

Variation partitioning has become a method of choice for the interpretation of beta diversity using tables of environmental and spatial variables. At the last count, the ISI Web of Knowledge of the Institute for Scientific Information listed 603 papers that had used the method or were referring to it. The published examples concern most groups of organisms. An example is the analysis of the spatial variation of a community of oribatid mites in the peat carpet of a peat bog. Thirty-five mite species collected in 70 soil cores were analysed by variation partitioning with respect to a set of environmental and spatial variables. In the papers of Borcard et al. (1992) and Borcard and Legendre (1994), a polynomial function of the geographic coordinates was used as the spatial representation of the spatial relationships among the soil cores. In Borcard and Legendre (2002) and Borcard et al. (2004), PCNM spatial base functions were used instead, providing a much better explanation of the spatial variation in species composition among the cores (beta diversity).

Beta diversity is an interesting 'currency' for ecologists to compare either different sampling areas, or different ecological communities co-occurring in an area. (i) For the comparison of different study areas to be meaningful, the areas must be of the same size and sampled in the same way. An example would be the comparative study of the five 24 ha forest plots that are presently monitored under the auspices of the Chinese Forest Biodiversity Monitoring Network, forming a latitudinal gradient through China. The comparison would be meaningful if all compared plots are similarly divided into cells of 20 m × 20 m, or 40 m × 40 m, etc. In the framework of variation partitioning, SS($\mathbf{Y}$) is a convenient measure of beta diversity within each area. The total beta variation can be partitioned among one or several sets of environmental variables, as well as a table of spatial variables. The resulting partitions of the five separate areas can be compared using the results of these analyses. (ii) In each of these forest plots, one could compare the beta diversity of trees with that of other vegetation strata, for example, after dividing the plot into cells of equal sizes. The method of variation partitioning would allow researchers to partition the beta variation of each community among environmental and spatial variables and determine if the factors controlling the spatial organization are the same for the different groups of organisms.

Statistical analysis of community composition data must not be taken lightly. For proper tests of hypotheses concerning the factors responsible for the creation and maintenance of beta diversity in ecosystems, it is important to use tests of significance that do not rely on unrealistic assumptions, such as multivariate normality, when the data do not support these assumptions. Tests of significance must have correct type I error rates and good power to detect effects, whether natural or anthropogenic, when these effects are present. When significant effects are identified, one should use unbiased statistics ($R_a^2$) to report their magnitude. The conclusions reached during ecological analysis will be used by practitioners to take important decisions about the management of ecosystems, so they must be grounded in good science.

This report described the method of variation partitioning, which took many years to develop. Variation partitioning allows researchers to test precise hypotheses about the origin of beta diversity in ecosystems and determine how much of the spatial variation is controlled by environmental variables and how much remains unexplained. The latter fraction may be under the influence of unmeasured environmental variables, or else it may be determined by community processes such as competition or dispersal that need to be explored. In any case, the use of appropriate statistics is of foremost importance during ecological variation partitioning.

## Acknowledgements

# References

Anderson MJ, Legendre P (1999) An empirical comparison of permutation methods for tests of partial regression coefficients in a linear model. *J Stat Comp Simul* **62**:271–303.

Bell G (2001) Neutral macroecology. *Science* **293**:2413–2418.

Borcard D, Legendre P (1994) Environmental control and spatial structure in ecological communities: an example using oribatid mites (Acari, Oribatei). *Environ Ecol Stat* **1**:37–53.

Borcard D, Legendre P (2002) All-scale spatial analysis of ecological data by means of principal coordinates of neighbour matrices. *Ecol Modell* **153**:51–68.

Borcard D, Legendre P, Avois-Jacquet C, Tuomisto H (2004) Dissecting the spatial structure of ecological data at multiple scales. *Ecology* **85**:1826–1832.

Borcard D, Legendre P, Drapeau P (1992) Partialling out the spatial component of ecological variation. *Ecology* **73**:1045–1055.

Dray S, Legendre P, Peres-Neto PR (2006) Spatial modelling: a comprehensive framework for principal coordinate analysis of neighbour matrices (PCNM). *Ecol Modell* **196**:483–493.

Ezekiel M (1930) *Methods of Correlation Analysis*. New York: John Wiley and Sons.

He F (2005) Deriving a neutral model of species abundance from fundamental mechanisms of population dynamics. *Funct Ecol* **19**:187–193.

Hubbell SP (2001) *The Unified Neutral Theory of Biodiversity and Biogeography*. Princeton, NJ: Princeton University Press.

Legendre P (1990) Quantitative methods and biogeographic analysis. In: Garbary DJ, South RG (eds). *Evolutionary biogeography of the marine algae of the North Atlantic*. NATO ASI Series, Vol. G 22. Berlin: Springer-Verlag, 9–34.

Legendre P (1993) Spatial autocorrelation: trouble or new paradigm? *Ecology* **74**:1659–1673.

Legendre P, Borcard D, Peres-Neto PR (2005) Analyzing beta diversity: partitioning the spatial variation of community composition data. *Ecol Monogr* **75**:435–450.

Legendre P, Gallagher ED (2001) Ecologically meaningful transformations for ordination of species data. *Oecologia* **129**:271–280.

Legendre P, Legendre L (1998) *Numerical Ecology*, 2nd English edn. Amsterdam: Elsevier Science BV.

Manly BJF (1997) *Randomization, Bootstrap and Monte Carlo Methods in Biology*, 2nd edn. London: Chapman and Hall.

Miller JK (1975) The sampling distribution and a test for the significance of the bimultivariate redundancy statistic: a Monte Carlo study. *Multivar Behav Res* **10**:233–244.

Miller JK, Farr SD (1971) Bimultivariate redundancy: a comprehensive measure of interbattery relationship. *Multivar Behav Res* **6**:313–324.

Ohtani K (2000) Bootstrapping $R^2$ and adjusted $R^2$ in regression analysis. *Econom Modell* **17**:473–483.

Oksanen J, Kindt R, Legendre P, O'Hara RB (2007) *vegan: Community Ecology Package Version 1.8–5*. URL http://cran.r-project.org/.

Peres-Neto PR, Legendre P, Dray S, Borcard D (2006) Variation partitioning of species data matrices: estimation and comparison of fractions. *Ecology* **87**:2614–2625.

Rao CR (1964) The use and interpretation of principal component analysis in applied research. *Sankhyaá, Ser A, Indian J Stat* **26**:329–358.

R Development Core Team (2007) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL http://www.R-project.org.

ter Braak CJF (1986) Canonical correspondence analysis: a new eigenvector technique for multivariate direct gradient analysis. *Ecology* **67**:1167–1179.

ter Braak CJF (1987a) The analysis of vegetation–environment relationships by canonical correspondence analysis. *Vegetatio* **69**:69–77.

ter Braak CJF (1987b) Ordination. In: Jongman RHG, ter Braak CJF, van Tongeren OFR (eds). *Data Analysis in Community and Landscape Ecology*. Pudoc, Wageningen [reissued in 1995 by Cambridge: Cambridge University Press], 91–173.

ter Braak CJF (1988) *CANOCO—A FORTRAN Program for Canonical Community Ordination by [Partial] [Detrended] [Canonical] Correspondence Analysis, Principal Component Analysis and Redundancy Analysis (Version 2.1)*. Wageningen: Agricultural Mathematics Group, Ministry of Agriculture and Fisheries.

ter Braak CJF, Smilauer P (2002) *Canoco Reference Manual and CanoDraw for Windows User's Guide: Software for Canonical Community Ordination (Version 4.5)*. Ithaca, NY: Microcomputer Power.

Whittaker RH (1960) Vegetation of the Siskiyou mountains, Oregon and California. *Ecol Monogr* **30**:279–338.

Whittaker RH (1972) Evolution and measurement of species diversity. *Taxon* **21**:213–251.