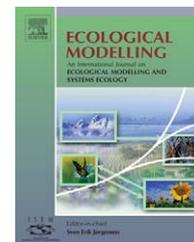


available at www.sciencedirect.comjournal homepage: www.elsevier.com/locate/ecolmodel

Spatial modelling: a comprehensive framework for principal coordinate analysis of neighbour matrices (PCNM)

Stéphane Dray^{a,b,*}, Pierre Legendre^a, Pedro R. Peres-Neto^{a,c}

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

^b Laboratoire de biométrie et biologie évolutive (UMR 5558); CNRS; Université Lyon I 43 boulevard du 11 Novembre 1918, 69622 Villeurbanne Cedex, France

^c Department of Biology, University of Regina, Regina, Saskatchewan, Canada S4S 0A2

ARTICLE INFO

Article history:

Received 4 May 2005

Received in revised form 10

December 2005

Accepted 9 February 2006

Published on line 23 March 2006

Keywords:

Ecological community

Eigenvalue

Eigenvector

Induced spatial dependence

Moran's I

Principal coordinates of neighbour matrices (PCNM)

Spatial autocorrelation

Spatial model

Spatial structure

ABSTRACT

Spatial structures of ecological communities may originate either from the dependence of community structure on environmental variables or/and from community-based processes. In order to assess the importance of these two sources, spatial relationships must be explicitly introduced into statistical models. Recently, a new approach called principal coordinates of neighbour matrices (PCNM) has been proposed to create spatial predictors that can be easily incorporated into regression or canonical analysis models, providing a flexible tool especially when contrasted to the family of autoregressive models and trend surface analysis, which are of common use in ecological and geographical analysis. In this paper, we explore the theory of the PCNM approach and demonstrate how it is linked to spatial autocorrelation structure functions. The method basically consists of diagonalizing a spatial weighting matrix, then extracting the eigenvectors that maximize the Moran's index of autocorrelation. These eigenvectors can then be used directly as explanatory variables in regression or canonical models. We propose improvements and extensions of the original method, and illustrate them with examples that will help ecologists choose the variant that will better suit their needs.

© 2006 Elsevier B.V. All rights reserved.

1. Introduction

One of the major current questions in ecology concerns the identification and explanation of the spatial variability of ecological structures (Cormack and Ord, 1979, Smith, 2002). Space can be considered either as a factor responsible for ecological structures, or as a confounding variable leading to bias when analyzing a process of particular interest. This realization leads ecologists to introduce space as either a predictor

or a covariable in statistical models. These two approaches have been used in various contexts such as the analysis of patterns of species richness (Blackburn and Gaston, 1996b, Selmi and Boulinier, 2001), species range sizes (Blackburn and Gaston, 1996a), species associations (Roxburgh and Matsuki, 1999), metacommunity analysis (Olden et al., 2001), population (Pettorelli et al., 2003) and community ecology (Borcard et al., 1992, Wagner, 2003, Borcard et al., 2004, Peres-Neto, 2004, Wagner, 2004).

* Corresponding author at: Laboratoire biométrie et biologie évolutive (UMR 5558); CNRS; University Lyon I 43 boulevard du 11 Novembre 1918, 69622 Villeurbanne Cedex, France. Tel.: +33 4 72 43 27 57; fax: +33 4 72 43 13 88.

E-mail addresses: dray@biomserv.univ-lyon1.fr (S. Dray), pierre.legendre@umontreal.ca (P. Legendre), pedro.peres-neto@uregina.ca (P.R. Peres-Neto).

0304-3800/\$ – see front matter © 2006 Elsevier B.V. All rights reserved.

doi:10.1016/j.ecolmodel.2006.02.015

Spatial structures observed in ecological communities can arise from two independent processes (Legendre, 1993, Legendre and Legendre, 1998, Section 1.1, Fortin and Dale, 2005, Chapters 1 and 5). Environmental factors that influence species distributions are usually spatially structured and then, through an indirect process, communities of species are also spatially structured; this process is called induced spatial dependence. Spatial autocorrelation can also be created directly at the community level as a result of contagious biotic processes such as growth, differential mortality, seed dispersal, or competition dynamics. In most situations, the spatial heterogeneity of communities is due to the simultaneous action of these two processes. Variation partitioning (Borcard et al., 1992, Borcard and Legendre, 1994, Méot et al., 1998) can be used to assess the importance of these two sources of spatial structure.

Incorporating spatial variation in ecological models requires tools to explicitly describe spatial relationships as predictors or covariables. Sokal (1979) used various functions of geographic distances among sites in Mantel tests to account for autocorrelation due to isolation by distance in population genetics models. Polynomial functions of the geographic coordinates can also be used as regressors to generate trend surfaces (Student [Gosset] 1914, Gittins, 1968). These spatial base functions have been used to model spatial relationships (often called “space” for short in scientific papers) in multivariate analyses such as canonical correlation analysis (Gittins, 1985, Pélissier et al., 2002, Gimaret-Carpentier et al., 2003), canonical correspondence analysis (CCA, Borcard et al., 1992, Borcard and Legendre, 1994, Méot et al., 1998) or redundancy analysis (RDA, Legendre, 1993). However, the use of trend surfaces is only satisfactory when the sampling area is roughly homogeneous, the sampling design is nearly regular, the number of spatial locations is “reasonable” (Norcliffe, 1969, Scarlett, 1972), and the spatial structure to be modelled is rather simple, such as a gradient, a single wave, or a saddle (Legendre and Legendre, 1998, Section 13.2). Moreover, the use of a trend surface introduces an arbitrary choice for the degree of the polynomial function. For instance, Wartenberg (1985a) used a second-degree polynomial while Borcard et al. (1992) used a polynomial of degree 3. In any case, polynomial trend surfaces of these degrees only allow the modelling of broad-scale spatial structures. Another problem concerns the correlations between these spatial predictors, which can be addressed by using an orthogonalization procedure in order to obtain orthogonal polynomials, but the higher-degree terms may be difficult to interpret in the case of surfaces.

Recently, a new approach called principal coordinates of neighbour matrices (PCNM) has been proposed as an alternative to trend surface analysis (Borcard and Legendre, 2002). This method has already been used with success in several ecological applications (Borcard et al., 2004, Brind’Amour et al., 2005, Legendre et al., 2005). PCNM base functions are obtained by a principal coordinate analysis (PCoA, Gower, 1966) of a truncated pairwise geographic distance matrix between sampling sites. Eigenvectors associated with the positive eigenvalues and corresponding to the Euclidean representation of the truncated distance matrix are used as spatial predictors in multivariate regression or canonical analysis (e.g., RDA, CCA). Even though this approach produces interesting and ecologi-

cally interpretable results (e.g., Borcard et al., 2004), it suffers from a lack of mathematical formalism. Indeed, these authors stated in the original description of the methods that: “This paper raises a number of mathematical questions [...] We hope that the paper will attract the interest of mathematicians who can help us understand these properties and develop methods of spatial modeling further” (Borcard and Legendre, 2002, p. 67).

In the present paper, we investigate the mathematical foundations of PCNM analysis and show that this approach is closely related to spatial autocorrelation structure functions. Using these theoretical properties, we develop improvements and extensions of the original approach. We hope this paper will help ecologists use the full potential of PCNM analysis for ecological applications and perceive the method as an extremely flexible and robust technique for the analysis of spatial problems.

2. The original PCNM approach

Generation of PCNM base functions is quite straightforward, requiring the following three main steps (Borcard and Legendre, 2002):

- (1) Compute a pairwise Euclidean (geographic) distance matrix between the n sampling locations ($D = [d_{ij}]$).
- (2) Choose a threshold value t and construct a truncated distance matrix using the following rule:

$$D^* = \begin{cases} d_{ij} & \text{if } d_{ij} \leq t \\ 4t & \text{if } d_{ij} > t \end{cases}$$

- (3) Perform principal coordinate analysis (PCoA) of the truncated distance matrix D^* . This analysis consists in the diagonalization of Δ where:

$$\begin{aligned} \Delta &= \left[-\frac{1}{2} (d_{ij}^{*2} - d_i^{*2} - d_j^{*2} + d^{*2}) \right] \\ &= -\frac{1}{2} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^t}{n} \right) D_2^* \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^t}{n} \right) \end{aligned}$$

where

$$\begin{aligned} d_i^{*2} &= \frac{1}{n} \sum_{j=1}^n d_{ij}^{*2}, \quad d_j^{*2} = \frac{1}{n} \sum_{i=1}^n d_{ij}^{*2}, \quad d^{*2} = \frac{1}{n} \sum_{i=1}^n d_i^{*2} \\ \text{and } D_2^* &= [(d_{ij}^*)^2] \end{aligned} \quad (1)$$

\mathbf{I} is the identity matrix and $\mathbf{1}$ is a vector containing all 1s.

After diagonalization, principal coordinates are obtained by scaling each eigenvector \mathbf{u}_k of Δ to the length $\sqrt{|\lambda_k|}$ where λ_k is the eigenvalue associated with eigenvector \mathbf{u}_k . Since the original Euclidean distance matrix has been truncated, there are negative eigenvalues, so that it is impossible to represent D^* entirely in a Euclidean space. Hence, in the original PCNM method, only the principal coordinates associated with positive eigenvalues (corresponding to the Euclidean representation of D^*) are kept and used as spatial descriptors; see the discussion in Section 4 on the usefulness of negative eigenvalues.

3. Distances, similarities, and spatial weighting matrices

PCoA is usually computed on a distance matrix but Gower (1966) showed that this analysis can also be computed from a similarity matrix (also shown in Legendre and Legendre, 1998, p. 431). For instance, consider the similarity matrix S derived from a distance matrix D :

$$S = [s_{ij}] = \left[1 - \left(\frac{d_{ij}}{\max(d_{ij})} \right)^2 \right] = \mathbf{1}\mathbf{1}^t - \frac{D_2}{\max(d_{ij})^2} \text{ with } D_2 = [d_{ij}^2] \quad (2)$$

These similarities vary between 0 (for $d_{ij} = \max(d_{ij})$) and 1 (for $d_{ij} = 0$). It is easy to show that a PCoA performed on the distance matrix D is equivalent to the diagonalization of $\Delta = (\max(d_{ij})^2/2)(I - (\mathbf{1}\mathbf{1}^t/n))S(I - \mathbf{1}\mathbf{1}^t/n)$ which provides the same eigenvectors as those obtained by diagonalization of $(I - \mathbf{1}\mathbf{1}^t/n)S(I - \mathbf{1}\mathbf{1}^t/n)$.

It is possible to interpret this similarity matrix S as a weighted graph (Chung, 1997, see Fig. 1). Each non-null value s_{ij} indicates a connection between sites i and j ; the intensity of the connection is expressed by the value s_{ij} . A more interesting interpretation is to consider S to be a spatial weighting matrix (Bavaud, 1998, Tiefelsdorf et al., 1999), which indicates the strength of the potential interaction among the spatial units. In the PCoA of D (Fig. 1a), all sites are considered to be neighbours except those corresponding to the largest distance, for which the similarity is null. The intensity of the link is expressed by $s_{ij} = 1 - (d_{ij}/\max(d_{ij}))^2$. In the original PCNM approach, the spatial weighting matrix is:

$$S^* = \mathbf{1}\mathbf{1}^t - \frac{D_2^*}{(4t)^2} \text{ with } D_2^* = [d_{ij}^{*2}] \quad (3)$$

If the distance d_{ij} is greater than the threshold value, $d_{ij}^* = 4t$ and $s_{ij}^* = 0$, sites i and j are not considered neighbours, whereas if d_{ij} is less than or equal to the threshold, the two sites are considered neighbours and the spatial link is $s_{ij}^* = 1 - (d_{ij}/4t)^2$. The weight associated with the link of a site with itself is $s_{ii}^* = 1$.

4. Moran's eigenvector maps (MEM)

In this section, we consider the n -by-1 vector $\mathbf{x} = [x_1 \dots x_n]^t$ containing measurements of a quantitative variable of interest at n sites and a n -by- n symmetric spatial weighting matrix W . The usual formulation for Moran's index of spatial autocorrelation (Moran, 1948, Cliff and Ord, 1973) is:

$$I(\mathbf{x}) = \frac{n \sum_{(2)} w_{ij} (x_i - \bar{x})(x_j - \bar{x})}{\sum_{(2)} w_{ij} \sum_{i=1}^n (x_i - \bar{x})^2} \text{ where } \sum_{(2)} = \sum_{i=1}^n \sum_{j=1}^n \text{ with } i \neq j \quad (4)$$

The values w_{ij} are weights from matrix W . In Moran's I autocorrelation analysis, we usually make $w_{ij} = 1$ for sites i and j that are within the distance class under consideration and $w_{ij} = 0$ for the sites that are outside that distance class. Note, however, that matrix W can be any non-negative spatial weighting matrix. Moran's I can be positive or negative. It can be rewritten as follows using matrix notation:

$$I(\mathbf{x}) = \frac{n}{\mathbf{1}^t W \mathbf{1}} \frac{\mathbf{x}^t (I - \mathbf{1}\mathbf{1}^t/n) W (I - \mathbf{1}\mathbf{1}^t/n) \mathbf{x}}{\mathbf{x}^t (I - \mathbf{1}\mathbf{1}^t/n) \mathbf{x}} \quad (5)$$

For a spatial matrix W , de Jong et al. (1984) have shown that the upper and lower values of Moran's I are given by $(n/\mathbf{1}^t W \mathbf{1})\lambda_{\max}$ and $(n/\mathbf{1}^t W \mathbf{1})\lambda_{\min}$ where λ_{\max} and λ_{\min} are the extreme eigenvalues of $\Omega = (I - \mathbf{1}\mathbf{1}^t/n)W(I - \mathbf{1}\mathbf{1}^t/n)$; this equation has the same form as Eq. (1). Hence, the eigenvectors of Ω are vectors with unit norm maximizing Moran's I under

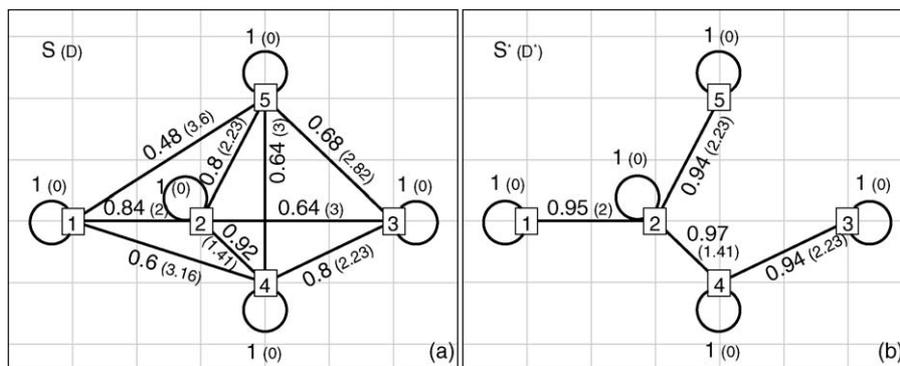


Fig. 1 – A small example to illustrate the equivalence between distances, similarities, and spatial weights. Five sites (numbers in boxes) are positioned according to their geographic coordinates. Euclidean distances (d_{ij}) among sites are computed and similarities are calculated from them. The largest distance is $d_{13} = 5$. To illustrate the property of spatial connectivity, a graph is drawn where each connection has length d_{ij} . The values on the graph connections are similarities, followed by distances in parentheses. In the case of a complete distance matrix (a), the weight given to each connection has the value of the similarity ($s_{ij} = 1 - (d_{ij}/\max(d_{ij}))^2$). No connection is represented between sites 1 and 3 because the weight is null ($s_{13} = s_{31} = 0$). Loops represent the influence of points on themselves; they are labelled by similarity (distance in parenthesis) values. In the original PCNM approach (b), the matrix is truncated using $t = 2.23$; the threshold is chosen to be the smallest distance that keeps all sites connected. The similarities are computed by $1 - (d_{ij}/4t)^2$ and all connections with lengths (distances) greater than the threshold are removed.

the constraint of orthogonality. The eigenvalues of this matrix are equal to Moran's I coefficients of spatial autocorrelation (post-multiplied by a constant), and similarly they can be either positive or negative. Eigenvectors associated with high positive (or negative) eigenvalues have high positive (or negative) autocorrelation and describe global (or local) structures; (this point of view of local and global structures was developed in Thioulouse et al., 1995). The eigenvectors associated with eigenvalues with extremely small absolute values correspond to spatial autocorrelation with low intensity and are not suitable for defining spatial structures. In this context, the PCNM base functions are particular cases of Moran's eigenvector maps (hereafter referred to as MEM), thus maximizing spatial autocorrelation with respect to the spatial weighting matrix defined by S^* . As the construction of the matrix S^* is deduced from the computation of the matrix of geographical distances, PCNM can be seen as distance-based eigenvectors maps (DBEM). MEM is a general framework which consists in the diagonalization of a spatial weighting matrix, DBEM is a particular case of this framework where the spatial weighting matrix is defined with distances and PCNM method is a particular case of DBEM (PCNM \subset DBEM \subset MEM).

5. Notes on the original PCNM approach

We have shown that the PCNM approach is closely related to Moran's index of spatial autocorrelation. This observation provides elements that will help improve the original method proposed by Borcard and Legendre (2002). As shown in the previous section, negative eigenvalues correspond to negative autocorrelation and their associated eigenvectors can be used to describe local structures. These structures can be produced by biotic processes such as species territoriality and competition. Hypotheses involving these types of ecological processes have rarely been tested, however. Scaling the eigenvectors to lengths $\sqrt{\lambda_k}$, as done in PCoA, must be avoided in order to retain these "negative" eigenvectors which, otherwise, would become imaginary. In any case, there is no need for such scaling of the PCNM base functions because predictors can be standardized or rescaled to unity before regression or multivariate analysis without any effect on the explanatory power (measured by R^2 or AIC) or the fitted values. Eigenvectors associated with eigenvalues having small absolute values are lightly spatially structured. Their role as spatial predictors is dubious and it could be useful to remove them before subsequent analyses. Since the eigenvalues are linearly related to Moran's I , one might consider testing this index (by permutation procedures) for each eigenvector and keeping only those that represent significant spatial autocorrelation.

A second aspect concerns the diagonal terms of the spatial weighting matrix in PCNM analysis. In the classical formulation of Moran's I , the summation of quantities is limited to $i \neq j$. This means that only spatial weighting matrices that have zeros on the diagonal are considered. Bavaud (1998) related spatial weighting matrices to Markov chains and considered the possibility of having non-zero values on the diagonal. Note, however, that this option introduces the influence of a point on itself when computing autocorrelation. Although, from a

theoretical point of view, this statement may seem difficult to justify in most cases, in practice, the influence of a non-zero diagonal can easily be assessed and discounted. Consider a symmetric matrix W with zero diagonal: the diagonalization of $\Omega = (I - 11^t/n)W(I - 11^t/n)$ produces eigenvectors u_i associated with eigenvalues λ_i . If we consider that a point may have influence on itself, as in the original PCNM approach (the values s_{ii}^* on the loops are 1 in Fig. 1b), the matrix to consider becomes $(I - (11^t/n))(W + I)(I - (11^t/n))$ which can be rewritten as $(I - (11^t/n))W(I - (11^t/n)) + (I - (11^t/n)) = \Omega + (I - (11^t/n))$. The matrix $\Omega + (I - (11^t/n))$ is doubly centred therefore their eigenvectors are centred. It can easily be shown that the eigenvectors of $\Omega + (I - (11^t/n))$ are u_i associated with eigenvalues that are $\lambda_i + 1$. Hence, eigenvectors with moderate negative autocorrelation ($\lambda_i < -1$) computed from W correspond to positive autocorrelation when the spatial weighting matrix is defined as $W + I$. This explains why the original PCNM approach produces a maximum of $2n/3$ positive eigenvalues in the case of n equidistant sites along a straight line, while one would expect to obtain an equal number of positive and negative eigenvalues.

Lastly, our interpretation of the PCNM approach facilitates the understanding of the proposition to multiply the threshold value t by a factor of 4. Borcard and Legendre (2002) justify this choice by the fact that they "observed that beyond a factor of four times the threshold for the 'large' distances, the principal coordinates remain the same to within a multiplicative constant". In the PCNM approach, the spatial link between two neighbours i and j (from Eq. (3)) is expressed by

$$s_{ij}^* = 1 - \left(\frac{d_{ij}}{4t} \right)^2 \quad (6)$$

The stability empirically observed by Borcard and Legendre (2002) is due to the fact that when the factor is sufficiently high (say 4), the second term of the spatial weighting function (Eq. (6)) is very small, and the spatial link between two neighbours tends to 1. Hence, PCNM eigenvectors are very close to MEM of a binary weighting matrix defined using a distance criterion ($w_{ij} = 1$ if $d_{ij} \leq t$, and 0 otherwise).

6. Choice of a spatial weighting matrix

Our interpretation of PCNM base functions as a particular case of MEM generalizes the original approach because "the use of a generalised weighting matrix [...] allows the investigator to choose a set of weights which he deems appropriate from prior considerations. This allows great flexibility" (Cliff and Ord, 1973, p. 12). Indeed, the spatial weighting matrix can be defined in different ways according to particular ecological hypotheses of interest and their spatial interactions (Sokal, 1979). The spatial weighting matrix $W = [w_{ij}]$ can be seen as the Hadamard product (element-wise product) of a connectivity matrix $B = [b_{ij}]$ by a weighting matrix $A = [a_{ij}]$ (i.e., $[w_{ij}] = [b_{ij}a_{ij}]$). The connectivity matrix B is binary, where a connection value of 1 is given for two sites that are connected (neighbours) and 0 otherwise. It can be constructed using distance criteria (select a distance threshold and connect all points that are within that distance of each other), or more sophisticated procedures such

as the Delaunay triangulation, Gabriel graph, relative neighbourhood graph, sphere of influence, or minimum spanning tree (Jaromczyk and Toussaint, 1992, Legendre and Legendre, 1998, p. 752). Matrix **A** can be used to weight the connections defined in **B** and make **W** more realistic. For instance, we can introduce the notion of geographic similarity in **A** using $1 - (d_{ij}/\max(d_{ij}))$, which varies linearly with the geographic distance (Aubry, 2000). In population genetics, Sokal (1979) tried several non-linear transformations of the distances as weights $A(d_{ij}^k, d_{ij}^{-k}, \ln(d_{ij}))$.

The choice of a spatial weighting matrix is a critical step because it can greatly influence the results of spatial analyses (Tiefelsdorf et al., 1999). In the case of regular sampling (e.g., a regular grid), structures defined by eigenvectors are roughly similar for different definitions of **W**. For irregular distributions of sites, however, the number of positive/negative eigenvalues and the spatial structures described by their associated eigenvectors are greatly influenced by the spatial relationships defined in **W**. For a given sampling scheme, various types of connections (**B**) and weighting functions (**A**) may lead to very different spatial structures displayed by eigenvectors. This point is illustrated by Fig. 2, which is described in the next paragraph. In a study reviewing the use of different forms of weighting matrices for spatial autoregressive modelling, Griffith (1995) and Griffith and Lagona (1998) have shown that a

parsimonious specification (i.e., small number of neighbours) of the relationships among sites is to be preferred. In the original PCNM approach, **B** is based on a distance criterion and **A** is defined by the function $1 - (d_{ij}/4t)^2$.

The structure of **B** is very sensitive to the distribution of sites and the shape of the PCNM base functions may, then, be greatly influenced by the sampling design. We illustrate this observation by considering two irregular samples of 100 sites along a transect (Fig. 2). The PCNM base functions obtained for these two samples display very different spatial structures and are only slightly correlated ($r_{1a,1b} = 0.28, r_{2a,2b} = 0.30$). A parsimonious specification of **B** (two consecutive points on the transect are neighbours) with $A = [1 - (d_{ij}/\max(d_{ij}))]$, as suggested by Griffith and Lagona (1998), produces eigenvectors that are more robust relative to the spatial distribution of the sampling sites. Indeed, Pearson correlations show that a first-neighbour connection (of the minimum spanning tree type) is less sensitive to variations in the sampling design ($r_{3a,3b} = 0.93, r_{4a,4b} = 0.98$).

Recommendation. The choice of the spatial weighting matrix **W** is the most critical step in spatial analysis. This matrix is a model of the spatial interactions recognized among the sites, all other interactions being excluded. In some cases, a theory-driven specification can be adopted, and the spatial weighting matrix can be constructed based upon biological

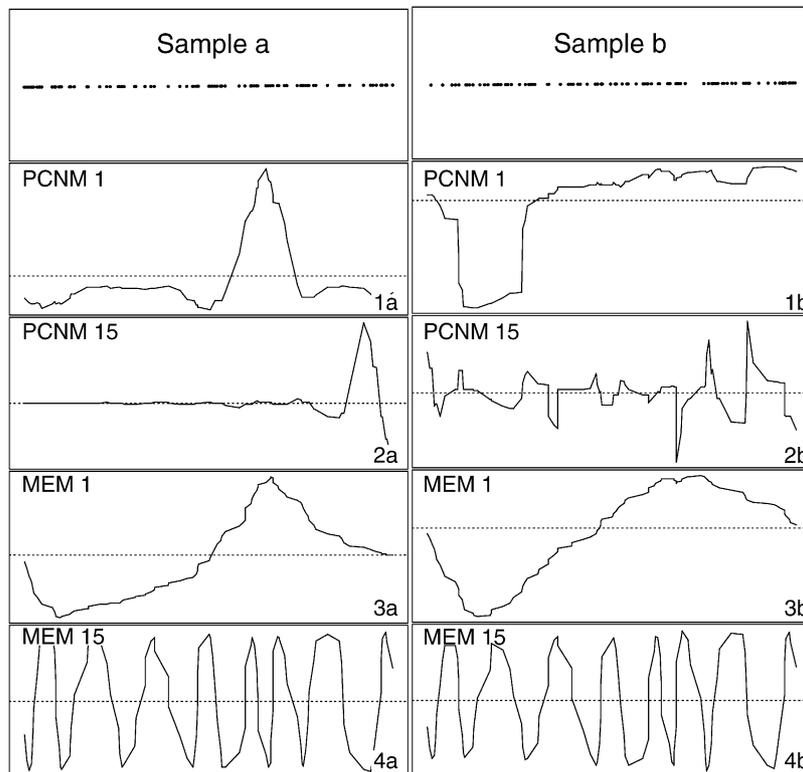


Fig. 2 – Comparisons of the eigenvectors obtained for two different spatial weighting matrices. Two irregular samples of 100 sites randomly positioned along a straight line are considered (a and b). For each sample, the first (1) and fifteenth (2) eigenvectors obtained by the original PCNM approach are presented. The first (3) and fifteenth (4) Moran’s eigenvectors (MEM) are also given for a spatial weighting matrix, corresponding to a minimum spanning tree (MST), where only consecutive points along the transect are considered neighbours, with weights equal to $1 - (d_{ij}/\max(d_{ij}))$. $r_{1a,1b} = 0.28, r_{2a,2b} = 0.3, r_{3a,3b} = 0.93, r_{4a,4b} = 0.98$. In each panel, the dotted line indicates the ordinate zero. Note how site positioning is less influential in the case of MST.

considerations (propagation process, Sokal and Oden, 1978, e.g., patch size, Hanski, 1994, dispersion capability, Knapp et al., 2003). Olden et al. (2001), for instance, discussed a variety of isolation distances; some of them may reflect better the challenges encountered by fish while dispersing. In most situations, however, the choice of a particular matrix may become rather difficult and a data-driven specification could then be applied. Under this latter approach, the objective is to select a configuration of \mathbf{W} that results in the optimal performance of the spatial model (Getis and Aldstadt, 2004). In the case of multiple linear regression of a response variable \mathbf{y} on a set of explanatory variables \mathbf{X} , the efficiency of the model can be assessed, for instance, by the Akaike information criterion (AIC). For standard linear models, we have:

$$AIC = -n \log \left(\frac{RSS}{n} \right) + 2K \quad (7)$$

where RSS is the residual sum of squares, K is the number of parameters in the fitted model and n represents the number of sites. When the sample size n is small, a bias correction is needed (Hurvich and Tsai, 1989) and the corrected AIC becomes:

$$AIC_c = \frac{AIC + 2K(K + 1)}{(n - K - 1)} \quad (8)$$

The lowest value of AIC_c identifies the best model. For the case of a multivariate response \mathbf{Y} (e.g., multiple species) as found in canonical analysis (e.g., RDA, CCA), an AIC-like criterion has recently been proposed (Godinez-Dominguez and Freire, 2003). It requires a multivariate analogue to the univariate RSS, which is computed for redundancy analysis (RDA) by

$$RSS = \text{trace}(\mathbf{Y}^t \mathbf{Y}) - \text{trace}(\hat{\mathbf{Y}}^t \hat{\mathbf{Y}}) \quad (9)$$

where $\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}$ contains the regression fitted values.

A procedure based on AIC_c can easily be implemented to decide which is the best weighting matrix among a number of suitable options. Our proposed approach consists of three main steps:

- (1) Define a set of possible spatial weighting matrices.
- (2) For each candidate:
 - Compute MEM and store them, by columns, in a matrix \mathbf{U} ($n \times p$).
 - Select the model with the lowest AIC_c . Note that this task is facilitated by the orthogonality of MEM because we can rewrite $RSS = \text{trace}(\mathbf{Y}^t \mathbf{Y}) - \sum_{i=1}^p \text{trace}(\mathbf{Y}^t \mathbf{u}_i \mathbf{u}_i^t \mathbf{Y})$. Each eigenvector \mathbf{u}_i reduces RSS by the value $\text{trace}(\mathbf{Y}^t \mathbf{u}_i \mathbf{u}_i^t \mathbf{Y})$. To obtain the best model for a given \mathbf{U} , we must:

Sort the eigenvectors in descending order according to their associated values $\text{trace}(\mathbf{Y}^t \mathbf{u}_i \mathbf{u}_i^t \mathbf{Y})$.

Enter the eigenvectors from the sorted list, one by one, in the model. We obtain p models with $K = 1, 2, \dots, p$ variables.

Compute AIC_c for each model.

Select the model with the lowest AIC_c .

- (3) Select the spatial weighting matrix corresponding to the model with the lowest AIC_c .

This data-driven approach selects the best matrix \mathbf{W} according to AIC_c . A complete example of this approach is presented next.

7. Ecological illustration

Here we illustrate the use of MEM and the data-driven process for selection of the spatial weighting matrix with a real data set. We re-examine data concerning the distribution of oribatid mites in the peat blanket of a bog lake. This data set has been used to illustrate the variation partitioning method with space modelled as a third order polynomial of geographic coordinates (Borcard et al., 1992; Borcard and Legendre, 1994) as well as the original PCNM approach (Borcard and Legendre, 2002; Borcard et al., 2004). The community has been described by the abundances of 35 mite species in 70 soil cores. Prior to the analyses reported here, the species data were Hellinger-transformed (Legendre and Gallagher, 2001) and detrended by multiple linear regression on geographic coordinates to remove the effect of a linear gradient; see Borcard et al. (2004) for details.

We tested different types of spatial weighting matrices using the data-driven procedure presented above. Five ways of defining neighbouring (Jaromczyk and Toussaint, 1992) relationships were used (matrix \mathbf{B}): Delaunay triangulation (**tri**), Gabriel graph (**gab**), relative neighbourhood graph (**rel**), minimum spanning tree (**mst**) and distance criterion (**dnn**). For the last approach, two sites i and j were considered as neighbours if $d_{ij} < \gamma$. In this case, ten values of the parameter γ evenly distributed between 1.011 m and 4 m were considered. The lowest value in this range (1.011 m) is the lowest value that keeps all sites connected; it corresponds to the longest edge of the minimum spanning tree constructed using geographic distances. The highest value (4 m) has been deduced from an empirical multivariate variogram (Wagner, 2003) of the transformed species data (Fig. 3). This multivariate extension is simply the sum of the empirical univariate variograms for all species. The chosen value corresponds to the highest distance at which the variogram is significant.

We assume that the ecological similarity between two sites is higher for site pairs that are spatially closer. This assumption must be taken into account when constructing the spatial weighting matrix \mathbf{W} . To achieve this goal, spatial weights are defined using monotonic decreasing functions varying with distance. Three functions have been considered in this example: a linear ($f_1 = 1 - d_{ij}/\max(d_{ij})$), a concave-down ($f_2 = 1 - (d_{ij}/\max(d_{ij}))^\alpha$), and a concave-up function ($f_3 = 1/d_{ij}^\beta$). We considered the sequence of integers between 2 and 10 for α , and between 1 and 10 for β . The case $\alpha = 1$ corresponds to the linear weighting function f_1 .

We computed MEM for the five types of binary connectivity matrices (**bin**) and for all combinations of these connectivity matrices with the three weighting functions. We used our data-driven specification procedure to identify the best spatial weighting matrix according to AIC_c . For each combination of a type of connectivity and a weighting option, we only report

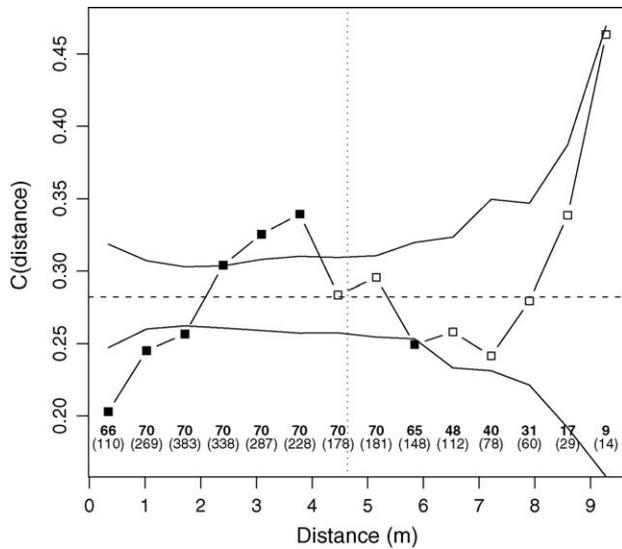


Fig. 3 – Empirical multivariate variogram of the oribatid mite community data. The variogram was computed with 14 distances classes of 0.69 m in width. The dotted line corresponds to half the maximum distance between sites, beyond which the variogram is not interpretable. The dashed line indicates the total inertia of the oribatid mite community data. For the k -th distance class, the envelope corresponds to the 0.025/14 and 1-0.025/14 quantiles (Bonferroni correction) of the distribution of 1000 variograms obtained after permutation of the species data. For each distance class, the numbers in bold above the x-axis indicate the number of sites used in the computation of the variogram; numbers in parentheses indicate the number of unique pairs of sites that fall into each distance class; significant values are indicated by filled symbols.

the values of α , β and/or γ producing the best model, as well as the corresponding value of AIC_c . We also included in our procedure the original PCNM (**pcnm**), a third-order polynomial of the geographic coordinates (**poly**), and a set of random orthogonal vectors (**rnd**).

The results are summarized in Table 1. The best spatial weighting matrix, corresponding to the lowest value of AIC_c , was constructed with a distance criterion (**dnn**) using $\gamma=2.67$ m as the cut-off value and with the concave-down weighting function f_2 with $\alpha=3$. Results of original PCNM analysis are better than those of a third-order polynomial, but many of the other spatial weighting matrices produced a better fit. The highest value of AIC_c was obtained with orthogonal random vectors, indicating that all methods provide fits better than chance alone. The best model for the **dnn-bin** option according to AIC_c was obtained with $\gamma=2.01$. This value corresponds to the last distance where significant positive autocorrelation was obtained in the variogram (Fig. 3). This element confirms the ability of AIC_c to select the “best” model. For the other weighting functions with **dnn** connectivity, no relationship can be found between values of γ and the variogram because the computation of this variogram is based on a binary weighting function.

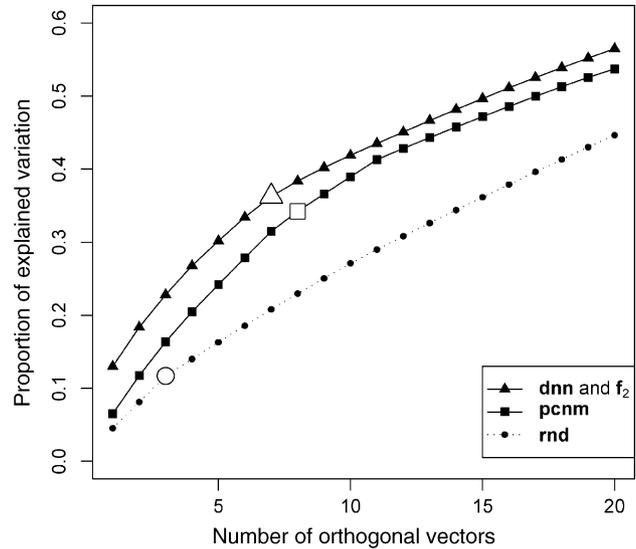


Fig. 4 – Proportion of the explained variation of the oribatid community as a function of the number of orthogonal vectors included in the model. Results are represented for a set of random orthogonal vectors (**rnd**), original PCNM (**pcnm**), and MEM of the best spatial weighting matrix selected by the data-driven procedure (**dnn** and f_2). For each set of predictors, the best model according to AIC_c is represented by a large empty symbol. Predictors are included in the model according to their associated values $\text{trace}(Y^t u_i u_i^t Y)$ ensuring that the fit obtained with k explanatory variables is the best one for a model of that rank.

The improvement due to the new approach is illustrated in Fig. 4. The data-driven procedure selected a model that increased the part of variation explained by space compared to the original PCNM approach. For a model containing 7 explanatory variables (number of variables of the best model according to AIC_c), space explained 31.4% of the variation of the oribatid mite community with the original PCNM approach and 36.2% with the **dnn- f_2** option. For a model with 8 variables (number of variables of the best model with PCNM as predictors, according to AIC_c) these values increase to 34.2% and 38.4% respectively.

8. Relationships with other eigenvector-based approaches

The new interpretation of the PCNM approach provided in this paper highlights relationships with other existing approaches. For instance, if all sites are connected (i.e., $\forall i, j [b_{ij}] = 1$) and $A = [1 - (d_{ij}/\max(d_{ij}))]$, the approach is equivalent to a PCoA based on $\sqrt{d_{ij}}$, proposed by Critchley (1978) as an alternative to multi-dimensional scaling. In the context of spatial analyses, Méot et al. (1993) diagonalized $D_w - W$ where $D_w = \text{Diag}(p_i)$ is a diagonal matrix containing the row sums of W ($p_i = \sum_{j=1}^n W_{ij} = \sum_{i=1}^n W_{ij}$). This diagonalization allowed

Table 1 – Results of the procedure for the data-driven specification of the spatial weighting matrix for the oribatid mite data set

Connectivity	Weighting function	AIC _c	Best model Number of variables	Values of parameters
tri	bin	–95.19	5	
	f1	–96.71	6	
	f2	–96.22	6	$\alpha = 10$
	f3	–96.83	7	$\beta = 1$
gab	bin	–96.09	9	
	f1	–97.70	9	
	f2	–98.55	8	$\alpha = 3$
	f3	–92.88	9	$\beta = 1$
rel	bin	–99.22	8	
	f1	–97.10	9	
	f2	–97.01	8	$\alpha = 5$
	f3	–94.01	9	$\beta = 1$
mst	bin	–98.92	6	
	f1	–97.02	8	
	f2	–98.29	5	$\alpha = 5$
	f3	–95.43	7	$\beta = 1$
dnn	bin	–100.56	5	$\gamma = 2.01$
	f1	–101.99	6	$\gamma = 3.66$
	f2	–102.70*	7	$\gamma = 2.67; \alpha = 3$
	f3	–100.51	8	$\gamma = 1.01; \beta = 1$
Other approaches				
pcnm		–97.85	8	
poly		–95.78	6	
rnd		–89.62	3	

For each combination of a type of connectivity and a weighting option, we only report the values of α , β and/or γ producing the best model according to AIC_c. See text for details. The lowest value of AIC_c (*) identifies the best spatial weighting matrix.

them to maximize (or minimize) Geary's index of spatial autocorrelation (Geary, 1954). When the neighbouring weights are uniform ($\forall i, p_i = 1/n$), their method is equivalent to minimizing (or maximizing) Moran's *I*. In this approach, only binary spatial weighting matrices are considered (i.e., $\forall i, j [a_{ij}] = 1$). Another attempt due to Griffith (1996) is based on the diagonalization of Ω to obtain eigenvectors maximizing Moran's *I*. Hence, our interpretation of PCNM is strictly equivalent to Griffith's approach. Note, however, that his work mainly used binary spatial weighting matrices.

9. MEM and spatial modelling

Autocorrelation is often related to a statistical problem because it introduces biases in standard statistical inference methods. Because the value observed at one site is influenced by the values at neighbouring sites, these values are not independent of one another. Since individual observations convey information about their neighbours, the number of degrees of freedom for a given set of observations may be reduced. That is the reason why, in the presence of positive autocorrelation, statistical tests become too liberal (the null hypothesis is rejected more often than it should). This problem has been well studied in the context of correlation analysis (Bivand, 1980, Clifford et al., 1989, Dutilleul, 1993), linear regression (e.g., Cliff and Ord, 1981) or analysis of variance (Legendre et al., 1990). The problem of spatial autocorrelation can be han-

dled in two ways: (i) the statistical method can be modified in order to take autocorrelation into account or (ii) the spatial dependency between observation can be removed; following that, a classical statistical method is used (spatial filtering).

In the context of linear regression, various approaches have been developed to take spatial autocorrelation into account. The linear model can be fully described in a matrix form: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$ with $\mathbf{e} \sim N(0, \mathbf{V})$ and $\mathbf{V} = \sigma^2\mathbf{I}$. The ordinary least squares (OLS) procedure allows the estimation of \mathbf{b} by minimizing $\sum e_i^2$. An explicit assumption of the OLS procedure is that the variance-covariance matrix \mathbf{V} is of the form $\sigma^2\mathbf{I}$. The generalized least squares (GLS) procedure allows users to deal with the more general situation in which $\mathbf{V} = \sigma^2\mathbf{D}$, where $\mathbf{D} \neq \mathbf{I}$. If \mathbf{V} is a symmetric matrix with a positive determinant, GLS can be treated as an OLS problem $\mathbf{y}^* = \mathbf{X}^*\mathbf{b} + \mathbf{e}$ using the transformed variables $\mathbf{y}^* = \mathbf{P}^{-1}\mathbf{y}$ and $\mathbf{X}^* = \mathbf{P}^{-1}\mathbf{X}$ where $\mathbf{D} = \mathbf{P}^t\mathbf{P}$. In the case of spatially autocorrelated errors, a GLS procedure can be used if the autocorrelation structure is known. The simultaneous scheme for error autocorrelation (SAR-error) assumes that $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$ where $\mathbf{e} = \rho\mathbf{W}\mathbf{e} + \mathbf{u}$ and $\mathbf{u} \sim N(0, \sigma^2\mathbf{I})$. The SAR-error model can be rewritten $\mathbf{y} = \mathbf{X}\mathbf{b} + \rho\mathbf{W}\mathbf{y} - \rho\mathbf{W}\mathbf{X}\mathbf{b} + \mathbf{u}$ and is equivalent to the GLS problem $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$ where $\mathbf{e} \sim N(0, \mathbf{V})$ and $\mathbf{V} = \sigma^2((\mathbf{I} - \rho\mathbf{W}^t)(\mathbf{I} - \rho\mathbf{W}^t))^{-1}$. Removing the term $\rho\mathbf{W}\mathbf{X}\mathbf{b}$ from a SAR-error model results in the simpler autoregressive response model (SAR-lag) $\mathbf{y} = \mathbf{X}\mathbf{b} + \rho\mathbf{W}\mathbf{y} + \mathbf{u}$ where $\mathbf{u} \sim N(0, \sigma^2\mathbf{I})$. An alternative to SAR is the conditional scheme model (CAR) in which $\mathbf{V} = \sigma^2(\mathbf{I} - \rho\mathbf{W}^t)^{-1}$. The reader should consult Wall (2004) for a comprehensive description of SAR and CAR,

and Lichstein et al. (2002) for an ecological application. Usually, the spatial weighting matrix W is row-standardized (row totals equal to 1) so that ρ is restricted within the range -1 to 1 . The SAR-error model could be fitted if the value of ρ was known: V could be estimated easily and GLS procedure could be used. However, in general, ρ is not known and must be estimated using a maximum likelihood procedure. A way to avoid this relatively complex numerical procedure is to guess a value of ρ . For instance, one can assign the maximum possible value $\rho = 1$ and the SAR-error model can then be rewritten as $(I - W)y = (I - W)Xb + u$. This approach is often called 'spatial differencing' because the transformed variables $(I - W)y$ and $(I - W)X$ are the differences between observed values and weighted averages (as reflected by W) of neighbouring values (e.g., Ord, 1975). If ρ has to be estimated, the maximum likelihood approach must be carried out. This procedure requires to compute the Jacobian determinant of the n -by- n matrix $I - \rho W$ and is computationally intensive when n is large. Many approaches have been proposed to speed up this step using eigendecomposition (Ord, 1975, improved in Smirnov and Anselin, 2001), LU or Cholesky decomposition (Pace, 1997, Pace and Barry, 1997a, 1997b), or approximation functions (e.g., Barry and Pace, 1999). Theoretically, autoregressive models can be fitted to a variety of response distributions, including binary (autologistic), and Poisson (auto-Poisson). However, the auto-Poisson model can only have negatively autocorrelated errors and is therefore of limited practical use. The autologistic model has been used in several ecological applications (e.g., Wu and Huffer, 1997) but estimation of parameters is impaired by several technical problems (Griffith, 2004).

In the context of multivariate analysis (multivariate response Y), spatial autocorrelation has rarely been introduced explicitly in statistical models. Lebart (1969) introduced spatial differencing in PCA, Wartenberg (1985b) proposed a PCA that finds linear combinations of variables maximizing Moran's index of autocorrelation, while developments of multivariate autoregressive models are recent (e.g., Jin et al., 2005).

Griffith (2000b, 2003) and Getis and Griffith (2002) suggested the use of eigenvectors as spatial predictors in multiple regression as an alternative to autoregressive models; this suggestion is similar to that of Borcard and Legendre (2002) to use PCNM base functions as spatial predictors in multiple regression or canonical analysis. The spatial filtering techniques convert variables that are spatially autocorrelated into spatially independent variables in an OLS regression framework. Advantages of this approach are linked to the use of an OLS technique: OLS regression is simpler than GLS regression, has a well-developed theory, and has available a battery of diagnostic tools that make interpretations easier. Moreover, the eigenvector approach presents the advantage of decomposing the variation into spatial and non-spatial components, and the orthogonality of eigenvectors facilitates the analysis and the understanding of the spatial part. Getis and Griffith (2002) demonstrated the efficiency of two approaches of spatial filtering (including the eigenvector approach) and the consistency of their results with those obtained by a spatial autoregressive model. The eigenvector approach has been introduced in the generalized linear modelling framework using Poisson

(Griffith, 2002) or binomial error (Griffith, 2004). In this context, spatial filtering is used as an alternative for auto-Poisson and autologistic models. It allows users to deal with negative as well as positive autocorrelation (contrary to the auto-Poisson model) and avoids computational problems related to the estimation of autoregressive models. A major drawback of the eigenvector approach is that it requires the diagonalization of the spatial weighting matrix, which is a computer intensive task for large matrices. Hopefully, Griffith (2000a) provided an analytical formulation for eigenvalues and eigenvectors (using cosines and sines) for regular square tessellations. These analytical results are very useful in the case of raster maps (e.g., satellite images) where each pixel is an observation.

10. Future directions

This paper provides new insights on the original formulation of the PCNM method, and introduces it in the framework of Moran's eigenvector maps. This formalism extends the original PCNM approach by allowing various definitions of spatial weighting matrices and other aspects related to this definition, as well as making it possible to consider negative spatial autocorrelation. Some questions remain to be solved, however. The first one concerns the choice of the eigenvectors to be introduced as spatial predictors in a statistical model. In all applications of PCNM analysis, this choice is done by classical forward selection, which tends to minimize the residual sum of squares (RSS). Forward selection is known to underestimate the residual variance (Freedman et al., 1992), becoming too liberal when there are a large number of candidate regressors (Westfall et al., 1998). In other words, even if the response variable is not spatially structured, the probability of selecting at least one MEM is greater than the chosen significance level. To solve this problem, Copas and Long (1991) proposed a correction of the residual degrees of freedom in the case of orthogonal regression, whereas multiple testing procedures are used by Westfall et al. (1998). Tiefelsdorf and Griffith (submitted) considered the possibility of basing model selection on the minimization of the spatial autocorrelation of residuals instead of RSS. In their approach, the retained model is the one that presented a Moran's I for the residual variation below an established threshold. A problem related to this approach is that some eigenvectors included in the model by the Moran's I minimization procedure can be non-significant using the RSS criterion. A more appropriate method may be to select the relevant spatial predictors with a criterion that will take into account both the fit of the model (RSS) and the spatial structure (Moran's I). Further work is required to evaluate these different procedures and to extend them to the case of multivariate response variables as in canonical analyses, which are routinely used in ecological research (Birks et al., 1996).

Lastly, our approach to Moran's eigenvectors maps is only suitable for symmetric spatial weighting matrices; we assume that the influence of site i on site j is equal to that of site j on site i . It could be interesting to extend this approach to non-symmetric spatial weighting matrices. This could be very useful in some ecological problems. For instance, we could take into account different upstream and downstream con-

nectivities in river networks. There are many possible avenues for expanding the applications of Moran's eigenvectors analysis. This study is an initiative to clarify the mathematical properties of PCNM analysis and show ecologists that it is a flexible and robust analytical tool for considering ecological data in a spatial context.

11. Supplement

An R package "spacemaker" containing functions to perform the analyses presented in the paper is available online. It includes a detailed documentation indicating how to create and manage spatial weighting matrices, compute their Moran's eigenvectors, and use the model selection procedure.

Acknowledgements

We would like to thank Daniel Borcard and the two reviewers for their comments on our manuscript. This research was supported by NSERC grant OGP0007738 to P. Legendre.

REFERENCES

- Aubry, P., 2000. Le traitement des variables régionalisées en écologie. Apports de la géomatique et de la géostatistique. Thèse de doctorat. Université Lyon I, Lyon.
- Barry, R., Pace, K., 1999. Monte Carlo estimates of the log determinant of large sparse matrices. *Linear Algebra Applications* 289, 41–54.
- Bavaud, F., 1998. Models for spatial weights: A systematic look. *Geogr. Anal.* 30, 153–171.
- Birks, H.J.B., Peglar, S.M., Austin, H.A., 1996. An annotated bibliography of canonical correspondence analysis and related constrained ordination methods 1986–1993. *Abstr. Bot.* 20, 17–36.
- Bivand, R., 1980. A Monte Carlo study of correlation coefficient estimation with spatially autocorrelated observations. *Quaest. Geogr.* 6, 5–10.
- Blackburn, T.M., Gaston, K.J., 1996a. Spatial patterns in the geographic range sizes of bird species in the New World. *Phil. Trans. Roy. Soc. Lond. Ser. B – Biol.* 351, 897–912.
- Blackburn, T.M., Gaston, K.J., 1996b. Spatial patterns in the species richness of birds in the New World. *Phil. Trans. Roy. Soc. Lond. Ser. B – Biol.* 19, 369–376.
- Borcard, D., Legendre, P., 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., Legendre, P., 2002. All-scale spatial analysis of ecological data by means of principal coordinates of neighbour matrices. *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.
- Brind'Amour, A., Boisclair, D., Legendre, P., Borcard, D., 2005. Multiscale spatial distribution of a littoral fish community in relation to environmental variables. *Limnol. Oceanogr.* 50, 465–479.
- Chung, F.K.R., 1997. Spectral graph theory. American Mathematical Society.
- Cliff, A.D., Ord, J.K., 1973. Spatial autocorrelation. Pion, London.
- Cliff, A.D., Ord, J.K., 1981. Spatial processes. Pion, London.
- Clifford, P., Richardson, S., Hémon, D., 1989. Assessing the significance of the correlation between two spatial processes. *Biometrics* 45, 123–134.
- Copas, J.B., Long, T., 1991. Estimating the residual error variance in orthogonal regression with variable selection. *The Statistician* 40, 51–59.
- Cormack, R.M., Ord, J.K., 1979. Spatial and temporal analysis in ecology. International Co-operative Publishing House, Fairland.
- Critchley, F., 1978. Multidimensional scaling: a short critique and a new method. In: Corsten, L.C.A., Hermans, J. (Eds.), *COMPSTAT 1978: Proceedings in computational statistics*. Physica-Verlag, Leiden.
- de Jong, P., Sprenger, C., van Veen, F., 1984. On extreme values of Moran's I and Geary's c. *Geogr. Anal.* 16, 17–24.
- Dutilleul, P., 1993. Modifying the t-test for assessing the correlation between two spatial processes. *Biometrics* 49, 305–314.
- Fortin, M.-J., Dale, M.B., 2005. Spatial analysis: a guide for ecologists. Cambridge University Press, Cambridge.
- Freedman, L.S., Pee, D., Midthune, D.N., 1992. The problem of underestimating the residual error variance in forward stepwise regression. *The Statistician* 41, 405–412.
- Geary, R.C., 1954. The contiguity ratio and statistical mapping. *The Incorporated Statistician* 5, 115–145.
- Getis, A., Aldstadt, J., 2004. Constructing the spatial weights matrix using a local statistic. *Geogr. Anal.* 36, 90–104.
- Getis, A., Griffith, D.A., 2002. Comparative spatial filtering in regression analysis. *Geographical Analysis* 34, 130–140.
- Gimaret-Carpentier, C., Dray, S., Pascal, J.-P., 2003. Broad-scale biodiversity pattern of the endemic tree flora of the Western Ghats (India) using canonical correlation analysis of herbarium records. *Ecography* 26, 429–444.
- Gittins, R., 1968. Trend-surface analysis of ecological data. *J. Ecol.* 56, 845–869.
- Gittins, R., 1985. Canonical Analysis, A Review with Applications in Ecology. Springer Verlag, Berlin.
- Godinez-Dominguez, E., Freire, J., 2003. Information-theoretic approach for selection of spatial and temporal models of community organization. *Mar. Ecol.-Prog. Ser.* 253, 17–24.
- Gower, J.C., 1966. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53, 325–338.
- Griffith, D.A., 1995. Some guidelines for specifying the geographic weights matrix contained in spatial statistical models. In: Arlinghaus, S.L. (Ed.), *Practical Handbook of Spatial Statistics*. CRC Press, Boca Raton, pp. 65–82.
- Griffith, D.A., 1996. Spatial autocorrelation and eigenfunctions of the geographic weights matrix accompanying geo-referenced data. *Can. Geogr.* 40, 351–367.
- Griffith, D.A., 2000a. Eigenfunction properties and approximations of selected incidence matrices employed in spatial analyses. *Linear Algebra Appl.* 321, 95–112.
- Griffith, D.A., 2000b. A linear regression solution to the spatial autocorrelation problem. *J. Geogr. Syst.* 2, 141–156.
- Griffith, D.A., 2002. A spatial filtering specification for the auto-Poisson model. *Stat. Prob. Lett.* 58, 245–251.
- Griffith, D.A., 2003. Spatial autocorrelation and spatial filtering: gaining understanding through theory and scientific visualization. Springer-Verlag, Berlin.
- Griffith, D.A., 2004. A spatial filtering specification for the autologistic model. *Environ. Planning A* 36, 1791–1811.
- Griffith, D.A., Lagona, F., 1998. On the quality of likelihood-based estimators in spatial autoregressive models

- when the data dependence structure is misspecified. *J. Stat. Planning Inference* 69, 153–174.
- Hanski, I., 1994. A practical model of metapopulation dynamics. *J. Animal Ecol.* 63, 151–162.
- Hurvich, C.M., Tsai, C.-L., 1989. Regression and time series model selection in small samples. *Biometrika* 76, 297–307.
- Jaromczyk, J.W., Toussaint, G.T., 1992. Relative neighborhood graphs and their relatives. *Proc. IEEE* 80, 1502–1517.
- Jin, X., Carlin, B., Banerjee, B., 2005. Generalized hierarchical multivariate CAR models for areal data. *Biometrics* 61, 950–961.
- Knapp, R.A., Matthews, K.R., Preisler, H.K., Jellison, R., 2003. Developing probabilistic models to predict amphibian site occupancy in a patchy landscape. *Ecol. Appl.* 13, 1069–1082.
- Lebart, L., 1969. *Analyse statistique de la contiguïté*. Publication de l'Institut de Statistiques de l'Université de Paris 28, 81–112.
- Legendre, P., 1993. Spatial autocorrelation: trouble or new paradigm? *Ecology* 74, 1659–1673.
- Legendre, P., Borcard, D., Peres-Neto, P.R., 2005. Analyzing beta diversity: partitioning the spatial variation of community composition data. *Ecol. Monogr.* 75, 435–450.
- Legendre, P., Gallagher, E.D., 2001. Ecologically meaningful transformations for ordination of species data. *Oecologia* 129, 271–280.
- Legendre, P., Legendre, L., 1998. *Numerical Ecology*, 2nd ed. Elsevier Science, Amsterdam.
- Legendre, P., Oden, N.L., Sokal, R.R., Vaudor, A., Kim, J., 1990. Approximate analysis of variance of spatially autocorrelated regional data. *J. Classification* 7, 53–75.
- Lichstein, J., Simons, T., Shiner, S., Franzreb, K., 2002. Spatial autocorrelation and autoregressive models in ecology. *Ecol. Monogr.* 72, 445–463.
- Méot, A., Chessel, D., Sabatier, R., 1993. In: Lebreton, J.D., Asselain, B. (Eds.), *Opérateurs de voisinage et analyse des données spatio-temporelles*. Biométrie et environnement. Masson, Paris.
- Méot, A., Legendre, P., Borcard, D., 1998. Partialling out the spatial component of ecological variation: questions and propositions in the linear modelling framework. *Environ. Ecol. Stat.* 5, 1–27.
- Moran, P.A.P., 1948. The interpretation of statistical maps. *J. Roy. Stat. Soc. Ser. B-Methodol.* 10, 243–251.
- Norcliffe, G.B., 1969. On the use and limitations of trend surface models. *Can. Geogr.* 13, 338–348.
- Olden, J.D., Jackson, D.A., Peres-Neto, P.R., 2001. Spatial isolation and fish communities in drainage lakes. *Oecologia* 127, 572–585.
- Ord, J., 1975. Estimation methods for models of spatial interaction. *J. Am. Stat. Assoc.* 70, 120–126.
- Pace, R., 1997. Performing large spatial regressions and autoregressions. *Econ. Lett.* 54, 283–291.
- Pace, K., Barry, R., 1997a. Sparse spatial autoregressions. *Stat. Probab. Lett.* 33, 291–297.
- Pace, K., Barry, R., 1997b. Quick computation of regressions with a spatially autoregressive dependent variable. *Geogr. Anal.* 29, 291–297.
- Pélicissier, R., Dray, S., Sabatier, D., 2002. Within-plot relationships between tree species occurrences and hydrological soil constraints: an example in French Guiana investigated through canonical correlation analysis. *Plant Ecol.* 162, 143–156.
- Peres-Neto, P.R., 2004. Patterns in the co-occurrence of fish species in streams: the role of site suitability, morphology and phylogeny versus species interactions. *Oecologia* 140, 352–360.
- Pettorelli, N., Dray, S., Gaillard, J.M., Chessel, D., Duncan, P., Illius, A., Guillon, N., Klein, F., Van Laere, G., 2003. Spatial variation in springtime food resources influences the winter body mass of roe deer fawn. *Oecologia* 137, 363–369.
- Roxburgh, S.H., Matsuki, M., 1999. The statistical validation of null models used in spatial association analyses. *Oikos* 85, 68–78.
- Scarlett, M.J., 1972. Problems of analysis of spatial distribution. In: Adams, W.P., Helleiner, F.M. (Eds.), *Congrès international de géographie*. University of Toronto Press, Montréal, pp. 928–931.
- Selmi, S., Boulinier, T., 2001. Ecological biogeography of southern ocean islands: the importance of considering spatial issues. *Am. Nat.* 158, 426–437.
- Smirnov, O., Anselin, L., 2001. Fast maximum likelihood estimation of very large spatial autoregression models. A characteristic polynomial approach. *Comput. Stat. Data Anal.* 35, 301–319.
- Smith, E.P., 2002. Ecological statistics. In: El-Shaarawi, A.H., Pieters, W.W. (Eds.), *Encyclopedia of Environmetrics*. John Wiley and Sons, Chichester, pp. 589–602.
- Sokal, R.R., 1979. Testing statistical significance of geographic variation patterns. *Syst. Zool.* 28, 227–232.
- Sokal, R.R., Oden, N.L., 1978. Spatial autocorrelation in biology 1. Methodology. *Biol. J. Linnean Soc.* 10, 199–228.
- Student [Gosset, W. S.], 1914. The elimination of spurious correlation due to position in time or space. *Biometrika* 10, 179–180.
- Thioulouse, J., Chessel, D., Champely, S., 1995. Multivariate analysis of spatial patterns: a unified approach to local and global structures. *Environ. Ecol. Stat.* 2, 1–14.
- Tiefelsdorf, M., Griffith, D. A., submitted. Semi-parametric filtering of spatial autocorrelation: the eigenvector approach. *Environment and Planning A*. <http://geog-www.sbs.ohio-state.edu/faculty/Tiefelsdorf/SpatialFiltering.pdf>.
- Tiefelsdorf, M., Griffith, D.A., Boots, B., 1999. A variance-stabilizing coding scheme for spatial link matrices. *Environ. Planning A* 31, 165–180.
- Wagner, H.H., 2003. Spatial covariance in plant communities: integrating ordination, geostatistics, and variance testing. *Ecology* 84, 1045–1057.
- Wagner, H.H., 2004. Direct multi-scale ordination with canonical correspondence analysis. *Ecology* 85, 342–351.
- Wall, M.M., 2004. A close look at the spatial structure implied by the CAR and SAR models. *J. Stat. Planning Inference* 121, 311–324.
- Wartenberg, D., 1985a. Canonical trend surface analysis: a method for describing geographic pattern. *Syst. Zool.* 34, 259–279.
- Wartenberg, D., 1985b. Multivariate spatial correlation: a method for exploratory geographical analysis. *Geogr. Anal.* 17, 263–283.
- Westfall, P.H., Young, S.S., Lin, D.K.J., 1998. Forward selection error control in the analysis of supersaturated designs. *Stat. Sin.* 8, 101–117.
- Wu, H., Huffer, F., 1997. Modelling the distribution of plant species using the autologistic regression model. *Environmental and Ecological Statistics* 4, 49–64.