# Random binary matrices in biogeographical ecology—Instituting a good neighbor policy

ARIF ZAMAN[1] and DANIEL SIMBERLOFF[2]

[1]*Lahore University of Management Sciences, LCCHS, Lahore, Pakistan*
*E-mail: arifz@lums.edu.pk*
[2]*Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville, TN 37996*

Binary matrices originating from presence/absence data on species (rows) distributed over sites (columns) have been a subject of much controversy in ecological biogeography. Under the null hypothesis that every matrix is equally likely, the distributions of some test statistics measuring co-occurrences between species are sought, conditional on the row and column totals being fixed at the values observed for some particular matrix. Many *ad hoc* methods have been proposed in the literature, but at least some of them do not provide uniform random samples of matrices. In particular, some ''swap'' algorithms have not accounted for the number of neighbors each matrix has in the universe of matrices with a set of fixed row and column sums. We provide a Monte-Carlo method using random walks on graphs that gives correct estimates for the distributions of statistics. We exemplify its use with one statistic.

## 1. Introduction

Binary matrices of a group of species' presences and absences on a set of sites, such as an archipelago of islands, have generated considerable controversy in ecological biogeography (Putman, 1994; Gotelli, 2000). Consider a matrix M of $m$ rows, each representing a species, and $n$ columns, each representing a site, with $M_{ij} = 1$ if species $i$ is present at site $j$ and $M_{ij} = 0$ otherwise. One may examine similarities (or dissimilarities) among species in their occurrences on sets of sites (e.g., Connor and Simberloff, 1979; Diamond and Gilpin, 1982; Gilpin and Diamond, 1982) or among sites in their occupancy by a group of species (e.g., Wright and Biehl, 1982; Simberloff and Connor, 1984). Such analyses are analogous, respectively, to *R*-mode and *Q*-mode analyses in numerical taxonomy (Simberloff and Connor, 1979), and both approaches have been employed to determine if the arrangement of species among sites is consistent with the hypothesis that each species' occurrences are independent of the other species' occurrences.

Here we examine methods used in *R*-mode analyses. The underlying idea is

straightforward and was presented in nascent form by Diamond (1975): if interspecific competition plays a large role in restricting species from occupying the same sites, then this fact should be reflected by the existence of pairs (or larger groups) of species that are mutually exclusive—they occupy no sites in common. The rub, according to Connor and Simberloff (1979), is that, even if species do not affect one another, one might still expect to find mutually exclusive pairs and larger groups of species. Thus the simple existence of such groups could not *ipso facto* constitute evidence for competition. Rather, one would need a null hypothesis about how many such mutually exclusive groups would be expected in the absence of species interactions.

The null hypothesis of Connor and Simberloff (1979) was that the presence/absence matrix was generated by species' colonizing the sites independently of one another and uniform randomly subject to the constraints that each species $i$ occupy the number of sites it occupies in nature $(r_i)$, and each site $j$ contain the number of species it contains in nature $(s_j)$. This null hypothesis has been attacked on several grounds, the main one of which is that a matrix generated under these constraints may incorporate historical effects of species interactions (Diamond and Gilpin, 1982; Gilpin and Diamond, 1982, 1984; Harvey *et al.*, 1983). Connor and Simberloff (1983, 1984) replied that the hypothesis is not being used to test whether interactions occurred, only whether the particular matrix is inconsistent with the hypothesis that they did not. We do not wish to argue the merits of these views here. Suffice it to say that this null hypothesis has continued to be used frequently, albeit sometimes with reservations, to study whether an observed matrix is unusual in some way (e.g., Matthews, 1982; Biehl and Matthews, 1984; Wilson, 1987, 1988; Wilson *et al.*, 1992; Guyer, 1990; Roberts and Stone, 1990; Stone and Roberts, 1990, 1992; Jackson *et al.*, 1992; Manly, 1995; Stone *et al.*, 1995; Sanderson *et al.*, 1998; Brualdi and Sanderson, 1999; Gotelli, 2000; Gotelli and Entsminger, 2000).

## 2. Generating random matrices

The strategy for using this null hypothesis (or any other) is to generate a random sample of matrices under the null hypothesis, then to compare the observed matrix to the random ones with respect to the distribution within the latter of some statistic, such as number of mutually exclusive pairs of species, or number of pairs sharing at most one site, etc. As the size of the matrix increases, the universe $U(\boldsymbol{R},\boldsymbol{S})$ of all matrices with the vector of row sums $\boldsymbol{R} = (r_1, \ldots, r_m)$ and vector of column sums $\boldsymbol{S} = (s_1, \ldots, s_k)$ can become so large so quickly (e.g., Simberloff and Connor, 1984; Snijders, 1991) that it is impractical to examine every matrix. Even counting the number of matrices in $U(\boldsymbol{R},\boldsymbol{S})$ is a difficult problem (see, e.g., Snapper, 1971), though not insurmountable (Sukhatme, 1938; Wang, 1988; Snijders, 1991), at least for fairly small matrices. Many researchers have tried to generate random samples from the set $U(\boldsymbol{R},\boldsymbol{S})$. The algorithms employed are instances of Markov Chain Monte Carlo algorithms, which have become a popular topic among statisticians over the last decade (e.g., Diaconis and Gangolli, 1995; Jerrum, 1998). We believe at least some algorithms used in biogeographical ecology have failed to produce uniform random samples (that is, all matrices in $U(\boldsymbol{R},\boldsymbol{S})$ equally likely to be included).

## 2.1 *Previous attempts*

Connor and Simberloff (1979) generated the sample of random matrices by computer simulation using pseudo-randomly generated numbers in what might be called a ''fill'' algorithm for sparse matrices. First the $r_1$ occurrences of species 1 were randomly distributed among the sites, then the $r_2$ occurrences of species 2 were randomly distributed subject to the constraint that none of the column sums $c_j$ exceeded the actual value in nature. Then the $r_3$ occurrences of species 3 were distributed, etc., through the $r_m$ occurrences of species *m*. At each successive addition of a ''1'' to the matrix, the resulting column sum was checked to ensure it did not exceed the actual value, and, if it did, this placement was not made and another random number was generated.

This procedure appeared to fill rather sparse matrices—substantially fewer 1s than 0s—adequately, but for the New Hebrides land birds, with 56 species arranged among 28 islands with a total of 887 species-occurrences, the computer algorithm often ''hung up,'' leading to a partially filled matrix that could not be completed subject to the constraints. Connor and Simberloff (1979) then turned to a ''swap'' algorithm (defined below). Additionally, it was not proved that all matrices are equally likely to be generated in the random sample.

Wilson (1987) modified this procedure so that occurrences are allocated at each stage either to the site with the highest number of unallocated occurrences or to the species remaining with the highest number of unallocated occurrences, whichever is the higher number. This modification was also used by Wilson (1988) and Wilson *et al.* (1992). This modification always allowed a simulated matrix to be filled, but Wilson (1987) did not demonstrate that all matrices in the universe have equal probability of inclusion in the random sample.

Snijders (1991) used a different fill algorithm to generate random binary matrices with fixed row and column sums, based on the concept of the enumeration tree (Verbeek and Kroonenberg, 1985) for contingency tables with fixed marginals. Although the matrices are not produced uniform randomly, Snijders (1991) also computed a coefficient for each generated matrix that approximates the relative frequency with which that matrix will be produced.

Pramanik (1994) provided yet another fill algorithm (also described by Rao *et al.*, 1996). However, Rao *et al.* (1996) found that it often sampled matrices far from uniformly, especially for small matrices.

To eliminate the possibility of computer gridlock for dense matrices, Connor and Simberloff (1979) developed an alternative procedure, also used by Simberloff (1986) and independently by Roberts and Stone (1990), which may be termed a ''swap'' approach (Gotelli and Entsminger, 2000). If one finds a $2 \times 2$ submatrix of *M* that is either diagonal or antidiagonal, then interchanging the 1s and 0s of such a submatrix changes neither a row nor a column sum. One thus begins with the matrix from nature, randomly chooses two rows and two columns, examines whether these produce a diagonal or antidiagonal matrix, and, if so, makes the swaps. Iterating this procedure produces a sample of matrices all with same sets of row and column sums. These authors did not prove that one could construct every matrix in $U(\boldsymbol{R}, \boldsymbol{S})$ by a finite sequence of such swaps, but, assuming this, construed their samples as uniform random ones from the relevant universe.

Of course, matrices that differ from one another by only one or a few swaps will not

differ very much from one another in statistics based on them (discussed below). That is, the null matrices examined might all be from one small section of $U(\boldsymbol{R},\boldsymbol{S})$, so statistics based on such a sample might not faithfully represent all of $U(\boldsymbol{R},\boldsymbol{S})$. Simberloff (1986) and Roberts and Stone (1990) tried to circumvent this problem by first making a large number of swaps in the observed matrix before beginning to accumulate the random sample of null matrices, then by including in the sample not every consecutive matrix but only one every $X$ swaps, where $X$ is a large number. Manly (1995) suggests that this procedure may not be necessary and that one can derive an adequate estimate of any statistic from any string of $K$ matrices generated by $K$ consecutive swaps. However, Manly's approach is an unbiased estimate only of the statistic among randomly drawn sequences of $K$ consecutive swaps; it is not an unbiased estimate of the statistic in all of $U(\boldsymbol{R},\boldsymbol{S})$. Thus, for example, it is quite possible that a particular sequence of such matrices would poorly estimate a statistic for the entire universe.

Rao *et al.* (1996) independently suggested the same swap algorithm as Roberts and Stone (1990), and they also pointed out that the matrices thus generated would not be uniform-randomly distributed across $U(\boldsymbol{R},\boldsymbol{S})$ (see below), producing an exact relative frequency for each matrix. They also addressed the problem of avoiding matrices very similar to one another (or to the original matrix) and concluded by inspection that an adequate number of swaps between matrices to be used in a ''uniform random'' sample is $3t$, where $t =$ the minimum of the number of 1s and the number of 0s in the matrix. Rao *et al.* (1996) compared computational speed for their programmed swap approach and the program for Snijder's fill approach and found the fill method to be far faster.

Sanderson *et al.* (1998; cf. Sanderson, 2000) inveighed against the entire swap approach on the grounds that it is computationally inefficient (at least if one does not adopt Manly's (1995) shortcut) and that a sequence of swaps can end up producing the same matrix more than one time. Rather, they suggested a fill algorithm based on a recursive method known as the knight's tour (Roberts, 1986). In this approach, at each step, a row and a column are each chosen randomly, and a ''1'' is placed if it does not violate row and column constraints (to this point, the knight's tour is identical to the fill method of Connor and Simberloff (1979)). However, if such a placement would violate a constraint, the algorithm goes through a series of reversals of previous placements, until it finds a previous matrix in the sequence that does permit the placement of a new, random ''1.'' Sanderson *et al.* (1998) have shown, based on Roberts (1986), that their approach can produce every matrix in $U(\boldsymbol{R},\boldsymbol{S})$. Further, by definition a knight's tour finds every member of the universe once and only once. However, they do not show that all matrices are equally likely to be included in the sample. That is, it is important to know whether their implementation of the knight's tour cause sequences of chosen matrices to be ''near'' one another in $U(\boldsymbol{R},\boldsymbol{S})$, and, in particular, whether they might tend to have similar values of whatever statistics are being tabulated.

Gotelli and Entsminger (2000) compared the swap algorithm of Roberts and Stone (1990) with their implementation of the fill algorithm of Sanderson *et al.* (1998). They found the swap algorithm to be much faster and not prone to type I or II errors. They could not compute error rates for the knight's tour algorithm. Additionally, they found the latter to produce much greater variances of the particular statistic (see below) that they measured on the various matrices.

## 2.2 *Random walks on a graph*

Ryser (1960) and Brualdi (1980) prove that a finite sequence of interchanges of the sort used by Simberloff (1986) and Roberts and Stone (1990) leads from every matrix $A$ in $U(\boldsymbol{R},\boldsymbol{S})$ to every other matrix $A'$ in $U(\boldsymbol{R},\boldsymbol{S})$. However, it is clear from both proofs that simply walking through $U(\boldsymbol{R},\boldsymbol{S})$ in this way will not ensure that each member of $U(\boldsymbol{R},\boldsymbol{S})$ is equally represented in the sample. For example, in Ryser's proof, $\bar{A}$ is the matrix $A$ with the ones and zeroes reversed, while $A^{\mathrm{T}}$ is the transpose of $A$. Define the matrix

$$D(A) = \bar{A}A^{T}. \tag{1}$$

Then $D_{ij}(A)$ is the number of times that a zero in the $i$th row has a corresponding one in the $j$th row, and so $M_{ij}(A) = D_{ij}(A)D_{ji}(A)$ is the number of $2 \times 2$ submatrices in rows $i$ and $j$ that are diagonal or antidiagonal. The total number of such submatrices in $A$ is

$$M(A) = \sum_{i<j} D_{ij}(A)D_{ji}(A). \tag{2}$$

Let $A_O$ be the matrix observed in nature, with row sum vector $\boldsymbol{R}$ and column sum vector $\boldsymbol{S}$. For any statistic $T$ defined for this matrix, we are interested in the probability that the observed value could arise by chance among all members of $U(\boldsymbol{R},\boldsymbol{S})$. The $P$-values of interest are $P\{T(A) \geq T(A_O)\}$ and $P\{T(A) \leq T(A_O)\}$, where the matrix $A$ is chosen uniform randomly from $U(\boldsymbol{R},\boldsymbol{S})$.

$A_O$ has $N(A_O)$ ''neighbors,'' i.e., matrices that can be reached from $A_O$ with one swap (Fig. 1). Choose $A_1$ equiprobably from the set of neighbors of $A_O$, then choose $A_2$ equiprobably from the $N(A_1)$ neighbors of $A_1$, etc. If we consider the graph with each node consisting of an element of $U(\boldsymbol{R},\boldsymbol{S})$, and lines connecting neighbors, then the above sequence can be construed as a Markov chain with state space $U(\boldsymbol{R},\boldsymbol{S})$. Random walks on graphs have often been used for Monte Carlo simulations where the generation of a random element of the graph is difficult, but location of a random neighbor is not. Unlike a
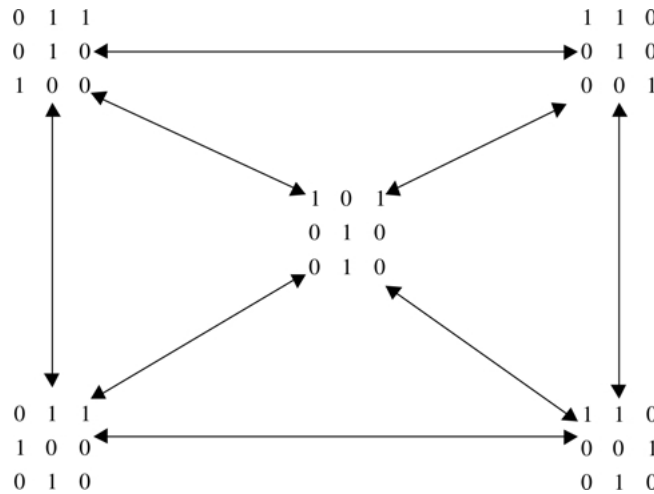


**Figure 1.** A graph of matrices in $U((2, 1, 1), (1, 2, 1))$ showing neighbor relationships.

simple random sample, a random walk will produce a dependent sequence of matrices, and with unequal sampling frequencies. Well-known results for such random walks (e.g., Aldous and Fill, 2000) state that:

a. Each node or matrix $A$ will be visited proportionally to its number of neighbors, $N(A)$. For example, in Fig. 1, the central matrix is visited 1/4 of the time while each of the remaining four matrices is sampled only 3/16 of the time.

b. The weighted average

$$n^{-1} \sum_{i=1}^{n} \frac{T(A_i)\bar{N}}{N(A_i)} \xrightarrow{\text{a.s.}} \bar{T}, \tag{3}$$

where $\bar{N}$ and $\bar{T}$ are the average values of $N_i$ and $T_i$, averaged over all $n$ graphs in $U(\boldsymbol{R},\boldsymbol{S})$.

c. The difference

$$n^{-1/2} \sum_{i=1}^{n} \left[ \frac{T(A_i)\bar{N}}{N(A_i) - n\bar{T}} \right] \xrightarrow{\text{dist.}} \quad \text{Normal} \quad (\mu = 0, \sigma^2). \tag{4}$$

While these equations are directly useful when $\bar{N}$ is known, in our case it must be estimated. It is a direct consequence of Equation (3) that

$$n^{-1} \sum_{i=1}^{n} \frac{1}{N(A_i)} \xrightarrow{\text{a.s.}} \bar{N}. \tag{5}$$

Applying Equation (5) to Equation (3), and noting that $N(A_i) \geq 1$, we can state that

$$\frac{\sum_{i=1}^{n} T(A_i)/N(A_i)}{\sum_{i=1}^{n} 1/N(A_i)} \xrightarrow{\text{a.s.}} \bar{T}, \tag{6}$$

which gives an estimate of $\bar{T}$ that does not involve any unknown quantities. This is of the form of a ratio estimate, and hence a biased estimate of its variance is available.

In simpler terminology, in order for the mean of the sample statistic $T$ to converge to the true mean $\bar{T}$ in the sample space, we need to take a weighted average; every time a matrix is produced, it is weighted inversely proportionally to its number of neighbors. In an Appendix available on website http://www.invasions.bio.utk.edu, we present some details of how to implement such a random walk on a computer.

## 3. Statistics

In most situations of biogeographical and ecological interest, there are many pairs of species (rows), and the pairs are not independent; that is, if species $i$ and $j$ tend to co-occur, and species $i$ and $k$ tend to co-occur, then species $j$ and $k$ tend to co-occur. Thus, using pairwise statistics and combining their significance levels is bound to be difficult. Sometimes it is appropriate to employ a single statistic that reflects the combined interaction for the entire matrix and find its level of signficance. Historically, many statistics of this sort have been used (Gotelli, 2000). We will present one that has not been

used in this setting, and rationalize it, but we emphasize that we could have used any of the statistics to exemplify the need to incorporate the number of neighbors in any algorithm using interchanges to generate samples of matrices.

## 3.1 *Previous statistics*

One approach has been to focus on particular pairs of species that share remarkably few sites or particular pairs of sites that share remarkably few species (e.g., Mosimann, 1968) and to assign probabilities to an observed degree of exclusion (e.g., from the hypergeometric distribution, as by Mosimann, 1968; Connor and Simberloff, 1979; Wright and Biehl, 1982). Since the choice of pairs selected is often guided by the data themselves, the reported significance levels are inflated. In addition, the significance levels for individual pairs are very sensitive to the particular sites that happen to have been included in the sample (Reddingius, 1983). The isolation of such pairs can be useful in generating hypotheses for further study, but their existence in the data set at hand cannot be construed as definitive evidence for anything.

As a test of the proposition that the entire matrix is an unusual member of $U(\boldsymbol{R},\boldsymbol{S})$ with respect to mutual exclusion among species, Connor and Simberloff (1979) suggested using the number of mutually exclusive pairs (as did Gotelli and Entsminger, 2000) and also the entire distribution of numbers of pairs sharing $0, 1, \ldots, k$ sites. The former statistic is easily located in the distribution produced by a random sample from $U(\boldsymbol{R},\boldsymbol{S})$. However, this statistic is quite restricted, in that it completely ignores pairs of species that share, for example, only one or two sites out of many. On the other hand, the entire distribution, which takes into account such pairs, is problematic because it is difficult to assign a probability to the difference between the distribution for the observed matrix and the ''average'' distribution for a random sample. The degrees of freedom for a chi-squared test are obviously lower than the nominal degrees of freedom because of dependence between pairs (Connor and Simberloff, 1979), but it is uncertain how few degrees of freedom actually obtain. Roberts and Stone (1990) suggest a way to estimate the degrees of freedom (and for a sample matrix their estimate is much lower than the nominal number).

Roberts and Stone (1990) suggest using as test statistic $S^2$, the sum, over all species pairs, of the square of the number of shared sites. This metric was criticized by Sanderson *et al*. (1998) but defended by Gotelli and Entsminger (2000). Stone and Roberts (1990) suggest as a statistic $C$, the mean number of ''checkerboard units'' per species pair for the community. A checkerboard unit for species pair $(i,j)$ is defined as the number of island pairs on which species $i$ and $j$ are mutually exclusive—in other words, the number of diagonal or antidiagonal $2 \times 2$ submatrices of $M$ for which the rows are $i$ and $j$. They compare this statistic to its distribution in a sample that is random (achieved by walking through $U(\boldsymbol{R},\boldsymbol{S})$ as outlined above) but not uniform (because the number of neighbors is not taken into account). Sanderson *et al*. (1998) used the entire distribution of numbers of pairs sharing $0, 1, \ldots, k$ sites, as did Connor and Simberloff (1979), albeit with a different test. Pielou and Pielou (1968) suggest yet another statistic: the total number of observed species combinations.

## 3.2 *The cross-product ratio for species pairs*

Another possible statistic measuring interactions between pairs, and one we will use for exemplary purposes, is the cross-product ratio (Bishop *et al.*, 1975). The log cross-product ratio $\lambda_{ij}$ between species $i$ and $j$ is defined as

$$\lambda_{ij}(A) = \log\left(\frac{D_{00}(A)D_{11}(A)}{D_{01}(A)D_{10}(A)}\right),$$

where $D_{00}$ is the number of sites with neither species, $D_{11}$ is the number with both species, $D_{01}$ is the number with species $j$ only, and $D_{10}$ is the number with species $i$ only. The cross-product ratio tends toward unity when there is no particular association (negative or positive) between the two species, exceeds unity when the species are positively associated, and is less than unity when they are negatively associated. In order to avoid infinities and zeroes, we use a modifed version of the cross product ratio, namely

$$\lambda_{ij}^* = \log\left(\frac{(D_{00} + 1)(D_{11} + 1)}{(D_{01} + 1)(D_{10} + 1)}\right).$$

For any pair of species, the distribution of this cross-product ratio can be determined for a uniform random sample from $U(\boldsymbol{R},\boldsymbol{S})$, using the weighted random walk described above. Specifically, to find the significance of a positive interaction between two species, we would let $T(A) = I\{\lambda_{ij}^*(A) > \lambda_{ij}^*(A_O)\}$ and use Equation (3) to find an estimate for $T$. But $T$ is exactly the $P$-value we are seeking. Notice that any statistic calculated for a particular pair of species will depend on the four numbers $D_{00}, D_{11}, D_{01}$, and $D_{10}$, and these four numbers are constrained by three equations:

$$D_{00} + D_{11} + D_{01} + D_{10} = k, \quad D_{10} + D_{11} = r_i, \quad D_{01} + D_{11} = r_j.$$

Thus, given any of the four numbers, one can calculate the other three when $\boldsymbol{R}$ and $\boldsymbol{S}$ are known. Therefore, any pairwise statistic will usually be a monotonic function of the cross-product ratio and hence will give the same $P$-values. For example, the statistic $C_{ij}$ of Stone and Roberts (1990) is $D_{10} \times D_{01}$, the denominator of the cross-product ratio. The difference comes in how the statistics for all the pairs of species are combined to produce a single statistic based on the entire matrix.

The statistic

$$\lambda_i(A) = \sum_{j \neq i} \log \lambda_{ij}^*(A), \tag{7}$$

can be used to measure the association of row $i$ with all the other rows. If this is large, the row is similar to many other rows; if small, the row is different from most others.

If these row statistics are again combined, we get the single statistic

$$\Lambda(A) = \sum_i \lambda_i(A), \tag{8}$$

which gives an overall measure of the similarity of the rows in the entire matrix.

Finally, if we want to consider one-sided tests to detect any interaction, either positive or negative, we may use the statistics

$$
\begin{array}{ccccccccc}
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
\end{array}
$$

**Figure 2.** Matrix with a highly skewed distribution of number of neighbors.

$$
\lambda_i^+(A) = \sum_{j \neq i} \left| \lambda_{ij}^*(A) \right|,
$$

$$
\Lambda^+(A) = \sum_i \lambda_i^+(A).
$$

The $\Lambda^+$ statistic is a non-negative number that is small where there is no association among species in the matrix and large when there is either positive or negative association in one or more pairs. The probability is assessed by placing $\Lambda^+$ in its distribution among a random sample generated as described above, with each element in the sample weighted inversely proportionally to its number of neighbors.

The effect of neglecting the number of neighbors depends very much on the particular structure of the matrix. For example, a matrix with a form similar to that in Fig. 2 would have many more neighbors than its neighbors would. In general, an $n \times n$ matrix of this shape would have $(n-1)^2$ neighbors, while each of its neighbors would have only $2n - 3$ neighbors. (Note that the central matrix of Fig. 1 is also of this type, with $n = 3$.) It is interesting to note that most statistics that measure interaction are computed on the basis of $D_{ij}(A)$, and that same quantity also determines the number of neighbors of a matrix. This fact suggests that the effect of neglecting the weighting by number of neighbors may introduce definite biases. We note in passing that we tested the performance of Snijder's fill algorithm in this regard on the matrix of Fig. 1 by generating 400 random matrices along with the weighting factor he describes for each of the five members of $U(\boldsymbol{R}, \boldsymbol{S})$. The weighted numbers of the five members did not differ from a uniform distribution ($\chi^2 = 2.206$, $df = 4$, $0.75 > P > 0.5$), that is, all matrices equally likely.

## 3.3 *Starting the walk*

We discard the first few matrices and wait for the sequence to reach approximate stationarity. The furthest distance between any two matrices on the graph can be at most km/4 steps, so taking some km initial steps is probably far enough. If, on the other hand, in estimating a *P*-value, the initial matrices are used, values similar to the observed matrix will be over-represented in the early sample, so the estimate will always start somewhat conservatively. Because in any large sample size the effects of the starting point will vanish, we preferred to use the conservative approach in our examples.

## 4.  Examples

We exemplify the method with three matrices discussed in the literature. First is a $13 \times 6$ matrix representing small Ozark fishes in the watersheds of the White River drainage (Matthews, 1982; Biehl and Matthews, 1984). The second is a hypothetical $20 \times 20$ matrix with all row and column sums equal to 10 (Fig. 3), one of a set of such matrices discussed by Gilpin and Diamond (1982, 1984) and Connor and Simberloff (1984). Finally, we examined the $56 \times 28$ matrix discussed above and by Connor and Simberloff (1979), Gilpin and Diamond (1984), Wilson (1987), Roberts and Stone (1990), Stone and Roberts (1990), Sanderson *et al*. (1998), and Gotelli and Entsminger (2000) representing land and freshwater birds of the New Hebrides. We constructed this matrix from the data of Diamond and Marshall (1976).

```
0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0
0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0
0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1
1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0
1 1 1 0 0 0 1 1 1 0 0 0 1 1 1 0 0 0 1 0
0 0 0 1 1 1 0 0 0 1 1 1 0 0 0 1 1 1 0 1
1 1 1 1 0 0 0 0 1 1 1 1 0 0 0 0 1 1 0 0
0 0 0 0 1 1 1 1 0 0 0 0 1 1 1 1 0 0 1 1
1 1 1 1 1 0 0 0 0 0 1 1 1 1 1 0 0 0 0 0
1 0 0 1 0 0 1 0 0 1 0 0 1 0 0 1 1 1 1 1
0 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 0 0 0 0
1 1 0 0 0 1 0 0 0 1 1 0 0 0 1 1 0 1 1 1
0 0 1 1 1 0 1 1 1 0 0 1 1 1 0 0 1 0 0 0
1 0 1 0 1 0 1 1 1 0 0 0 1 0 1 0 1 0 1 0
0 1 0 1 0 1 0 0 0 1 1 1 0 1 0 1 0 1 0 1
0 0 0 0 0 1 1 1 1 1 0 0 0 0 0 0 1 1 1 1
1 1 1 0 1 0 1 0 1 1 0 0 1 1 0 0 0 0 0 1
0 0 0 1 0 1 0 1 0 0 1 1 0 0 1 1 1 1 1 0
```

**Figure 3.**  A sample $20 \times 20$ matrix with 10 exclusively distributed pairs of species (rows).

The simulations were done using a computer program written in C, available from the junior author. The program provides estimates of the $P$-values associated with testing for unusual similarity/dissimilarity in the rows using all the variants of the $\lambda$ statistics described above. It is described in detail in the Appendix, available on website http://www.invasions.bio.utk.edu.

In the first matrix, $13 \times 6$, with 30,000 simulations, no pair of species stands out as unusually associated. For the entire matrix, the $P\{\Lambda(A) \geq \Lambda(A_O)\} = 0.0587$. Matthews (1982) and Biehl and Matthews (1984), generating matrices in the manner of Connor and Simberloff (1979) described above, found no reason to think the matrix unusual, although they did not associate a particular probability level with their result. Simply for comparison, the result of 725,000 simulations, in which matrices were not weighted according to their numbers of neighbors, was computed for the same data and found to yield $P = 0.0652$, about 12% higher than the correct value determined by the weighted average. The use of an unweighted average corresponds to the method used by Simberloff (1986) and Roberts and Stone (1990).

The hypothetical $20 \times 20$ matrix (Fig. 3) is highly unusual (Fig. 4), with an overall significance of 0.0002% for the observed $\Lambda$ (1,705,050 simulations). In fact, this is a conservative estimate as explained above, and it comes from having observed values as extreme as the observed matrix about three times in 1.7 million simulations. All three of these observations occurred right at the beginning of the simulation, so the $P$-value may be far smaller than $2 \times 10^{-6}$.

The pairwise $P$-values given in Fig. 4 easily pick out the ten mutually exclusive species pairs. Overall, the rest of the matrix does not appear to have unusual structure. Five other pairwise results are individually significant at the 0.01 to 0.05 level, but perhaps this is not surprising in a matrix of 190 entries. Connor and Simberloff (1984) also detected this matrix as unusual because of the exclusive pairs but could not associate a probability level with their result.

The New Hebrides $56 \times 28$ matrix is also highly unusual (Fig. 5), with an overall $P$-value of 0.0001% (4,134,448 simulations). The comment about the conservative nature of the $P$-value applies here as well. From the distribution of isolated pairwise cross-products, it is clear that certain species are associated negatively or positively with a particularly large number of other species. For example, species 51, *Lichmera incana*, has at least marginally significant negative associations with many of the other species. Oddly, this is one of the minority of species listed by Diamond and Marshall (1977) as not undergoing any niche shifts in response to interactions with other species. Connor and Simberloff (1979) found neither the number of exclusive pairs nor the entire distribution of number of pairs sharing $0, 1, \ldots, 28$ islands unusual, whereas Gilpin and Diamond (1984), using a swap algorithm, found $0.10 < P < 0.25$ by $\chi^2$ test of the entire distribution. Wilson (1987) found significant departures from a random expectation, although he did not assign a probability to the matrix as a whole. Stone and Roberts (1990), using the checkerboard statistic $C$, found the matrix to depart from a random one at $P < 0.001$, while Roberts and Stone (1990), using the shape of the entire distribution and an estimated number of degrees of freedom for a chi-squared test, found $P = 0.000005$. Roberts and Stone (1990) also used the $S^2$ statistic, and found $P < 0.001$. Sanderson *et al.* (1998) looked at the shape of the entire distribution and concluded that species pairs co-occurring exactly 2, 9, or 10 times are not within a 99% confidence interval, while the rest of the distribution is within the interval. On the other hand, using the $S^2$ statistic and their knight's tour algorithm to

|    | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|----|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 1  | • | 9 | . | . | . | . | . | . | . | 3  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 2  | . | • | . | . | . | . | . | . | 3 | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 3  | 2 | . | • | 9 | . | . | . | . | . | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 4  | . | 2 | . | • | . | . | . | . | . | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 5  | . | 2 | . | . | • | 9 | . | . | . | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 6  | 2 | . | . | . | . | • | . | . | . | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 7  | . | . | 2 | . | . | . | • | 9 | 2 | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 8  | . | . | . | 2 | . | . | . | • | . | 2  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 9  | 2 | . | . | . | . | . | . | 3 | • | 9  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 10 | . | 4 | . | . | . | . | 3 | . | . | •  | .  | .  | .  | .  | .  | .  | .  | .  | .  | .  |
| 11 | . | . | . | . | . | . | . | . | . | .  | •  | 9  | .  | .  | .  | .  | .  | .  | .  | .  |
| 12 | . | . | . | . | . | . | . | . | . | .  | .  | •  | .  | .  | .  | .  | .  | .  | .  | .  |
| 13 | . | . | 2 | . | . | . | . | . | . | .  | 2  | .  | •  | 9  | .  | .  | .  | .  | .  | .  |
| 14 | . | . | . | 2 | . | . | . | . | . | .  | .  | 2  | .  | •  | .  | .  | .  | .  | .  | .  |
| 15 | . | . | . | . | . | . | . | . | . | .  | .  | .  | 2  | .  | •  | 9  | .  | .  | .  | .  |
| 16 | . | . | . | . | . | . | . | . | . | .  | .  | .  | .  | 2  | .  | •  | .  | .  | .  | .  |
| 17 | . | . | . | . | . | . | . | . | . | .  | .  | .  | .  | .  | 2  | .  | •  | 9  | .  | .  |
| 18 | . | . | . | . | . | . | . | . | . | .  | .  | .  | .  | .  | .  | 2  | .  | •  | .  | .  |
| 19 | . | . | . | . | . | . | . | . | . | .  | .  | .  | .  | .  | .  | .  | 2  | .  | •  | 9  |
| 20 | . | . | . | . | . | . | . | . | . | .  | .  | .  | .  | .  | .  | .  | .  | 2  | .  | •  |

**Figure 4.** Matrix of *P*-values for cross-product ratios between rows of the $20 \times 20$ matrix discussed in text. Entries are $-2\log P$. Corresponding *P*-values as follows:

| Table entry $-2\log P$ | Range of P-values |
|---|---|
| 2 | 0.0562–0.1778 |
| 3 | 0.0178–0.0562 |
| 4 | 0.0056–0.0178 |
| 5 | 0.0018–0.0056 |
| 6 | 0.0006–0.0018 |
| 7 | 0.0001–0.0006 |

Entries above the diagonal correspond to $P\{\lambda_{ij}(A) \leq \lambda_{ij}(A_O)\}$, and those below the diagonal refer to the upper tail probabilities.

generate a sample of matrices, Sanderson *et al.* (1998) did not find the observed matrix to differ from expected (this is why they assail both the $S^2$ statistic and the swap algorithm).

## 5. Conclusion

The specific criticism we have leveled here against previous approaches applies only to certain swap algorithms. However, the problem of drawing a uniform random sample from $U(\textbf{R},\textbf{S})$ is a general one. For example, the knight's tour can arrive at every matrix in $U(\textbf{R},\textbf{S})$, but Sanderson *et al.* (1998) have not shown that a specific sample will not tend to

```
                    1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 4 5 5 5 5 5 5 5
            1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6
  1 * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 3 . . . . .
  2 . * . . . . . . . 3 . . . 3 . 2 . . . 5 . . . 3 . . . . . . . . . . . . . 2 . . 2 . . . . . 3 . . . 2 . . . 3 2
  3 . . * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 3 . . 3 . . . . 2
  4 . . . * . . . . 2 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 2 . . . . .
  5 . . . . * . . . . . . . . . . . . . 2 . . . . . . . . . . . 2 . 2 . . . . . . . . . . . . . . .
  6 5   4   * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 3 . . . . .
  7 . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . 2 . . . . . . 2 . . . . . . . 3 . .
  8 . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
  9 . . . 2 . . . . * 4 . . . . 4 . . . . . . . . . . . . . . . . . . . 3 2 . . . 3 . . . . . 5 2 . . 3 . .
 10 . . . . . . . . * . 6 3 . . . . . . . . . . . . . . . 2 . . . . . . . . . . . . . . . . . . .
 11 . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 2 . . . . . .
 12 . . . . . . . . . . * . . . . 2 . . . . . . . . . . 3 . . . . 3 5 . . 3 3 . . . 3 . 3 . . . . 2 . .
 13 . . . 2 . 2 . . . . . . * . . . . . . 2 . . . . . . . . 3 . . . . . 2 . . 3 . . . . . . . 3 . . . .
 14 . . . . . . . . . 2 . . * . . 3 . . . . . . . 3 3 . . . . . . . 2 . . . . . . . 3 . . . 5 . 2 . . . .
 15 . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
 16 . . . . . . . . . 3 . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
 17 . . . . . . . . . . . . . . . . * . . . . . . . . . 4 . . . . . . . . . . . . 7 . . . . .
 18 . . 2 . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . 5 . . . . 2
 19 . . . . . . 3 . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . . . 2
 20 . . . . . . . . . . . . . . . . . . . * . . . . . . 2 . . . . . . . . . . . . . . . . . .
 21 . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . . .
 22 . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . . .
 23 . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . . . . .
 24 . . . . . . . . . . . . . 2 . . . 2 . . . . . * . . . . . . . . . . . . . . . . . . . . .
 25 . . . . . . . . . . . . . . . . . . . . . . . . * . . . . 4 . . . . . . . . 3 . . . . . .
 26 . . . . . . . . . . . . . . . . . . . . . * . 2 . . . 2 . . . . . . . . . . . . . . . . .
 27 . . . . . . . . 3 . . . . 4 . . . . . . . . . * . . . . . . . . . . . . . . . 3 . . . . 2
 28 . . . . . . . . 4 . . . . . . . 2 . . . . . . . * . . . . . . . . . 3 . . . . . . . . . .
 29 . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . . .
 30 . . . . . . . . . . . . . 2 . . . . . . . . . . . . * . . . . . . . . . . . . 3 . . . . .
 31 . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . . . .
 32 . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . 3 .
 33 . . . . . . . . . . . . . . . . . . . . . . . . . . * . . 2 . . . . . . . . 3 . . 2 .
 34 . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . 2 .
 35 . . . . . . . . . . . . . . . . . . . . . . . 3 . . . * . . . . . . . . . . . . . . . . .
 36 . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . 2
 37 . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . . .
 38 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . . .
 39 . . . . . . . 3 . . . 2 . . . . . . . . . . . . . . . . . * . . . . . . . . 2 . . . . . .
 40 . . . . . . 6 . . . 3 . . . . . . . . . . 3 . . . . . * . . . . . . . . . . . . . . . . .
 41 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . . . . . . . .
 42 . . . . . . . . . . 2 . . . . . . . . . . . . . . . . . . * . . . . . . 3 . . . 2
 43 . . . . . . . 3 . . . . . . . . . . . . . . 4 3 . . * . . . . . . . . . 2 . . 2 .
 44 . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . . . 2 . . 4 .
 45 . . . . . . . . . 4 . . . . . . . . . . . . . 2 . . . * . . . . . . . . . . . . .
 46 . . . . . . 5 . . . 3 . . . . . . . . . . . . 3 3 . . . . * . . . . . . . . . .
 47 . . . . . . 3 . . . . . . . . . . . 3 . . . . . . . . * . . . . . . . . . .
 48 2 . . . . . . . . . . . . . . . . . . . . . . . . . . . * . 2 . . . . .
 49 . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . . .
 50 . . . . . . . . . 4 . . . . . . . . . . . . 2 . . 4 . . . . . * 2 . . . 4 .
 51 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . . .
 52 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . . .
 53 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . . .
 54 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * . .
 55 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . * .
 56 . . . . . . . 2 . . . . 2 . . . . . . . . . . . . . . . . 2 . . . 3 . . . . . . *
```

**Figure 5.** Matrix of pairwise significance levels for land birds of the New Hebrides. Interpretation of the entries is as in Fig. 4.

be restricted to one part of $U(\textbf{R},\textbf{S})$. A similar problem exists for the fill algorithm of Wilson (1987), and for this algorithm it has also not been demonstrated that all matrices will be drawn equiprobably.

Though it is now possible to use either a swap or a fill algorithm to find a truly uniform random subsample of binary matrices with fixed row and column sums (or one very close to uniform randomness), it still seems unlikely that examination of such matrices can allow strong inferences about what causes a matrix to depart from independence. Pielou and Pielou (1968), who adumbrated the matrix simulation approach, and Simberloff and Connor (1981) emphasize that both negative and positive associations among species can be explained on the basis of site differences as well as species interactions, and detailed research will be required to determine the reasons for departure from randomness. In retrospect, it should not be surprising that a single datum—the matrix of one group of species' occurrences on one set of sites—can provide at most a suggestion about fruitful avenues of research. After all, so many factors are likely to weigh in the determination of which species are found in which sites that it would be remarkable to be able to detect a signal, such as an interaction between a particular pair of species, unambiguously amidst the ''noise'' generated by all the other forces.

Constructing random binary matrices with fixed row and column sums has also been suggested as an aid in testing hypotheses in cladistic biogeography (Simberloff *et al.*, 1981) and vicariance biogeography (Connor, 1988), and the above algorithm should be used for this purpose. In ecological biogeography (Wright and Biehl, 1982; Biehl and Matthews, 1984) and cladistic systematics (Archie, 1989), sets of random binary matrices with either row or column sums constrained, but not both, have been produced for comparison to real matrices (cf. Gotelli, 2000). So long as the algorithm for generating permutations of 1s and 0s within a row (or column) draws uniform-randomly from the universe, this method should not be subject to the problem we discuss here. For ecological biogeography, Wright and Biehl (1982), Biehl and Matthews (1984), and Wilson (1987) discuss the relative merits of fixing both row and column sums and fixing one set only.

## Acknowledgments

## References

Aldous, D.J. and Fill, J.A. (2000) *Reversible Markov Chains and Random Walks on Graphs*. Book in preparation, relevant chapters on website http://www.stat.berkeley.edu/users/aldous/.

Archie, J.W. (1989) A randomization test for phylogenetic information in systematic data. *Systematic Zoology*, **38**, 239–52.

Biehl, C.C. and Matthews, W.J. (1984) Small fish community structure in Ozark streams: Improvements in the statistical analysis of presence-absence data. *American Midland Naturalist*, **111**, 371–82.

Bishop, Y.M.M., Fienberg, S.E., and Holland, P.W. (1975) *Discrete multivariate analysis*, MIT Press, Cambridge, Massachusetts.

Brualdi, R.A. (1980) Matrices of zeros and ones with fixed row and column sum vectors. *Linear Algebra and Its Applications*, **33**, 159–231.

Brualdi, R.A. and Sanderson, J.G. (1999) Nested species subsets, gaps, and discrepancy. *Oecologia*, **119**, 256–64.

Connor, E.F. (1988) Fossils, phenetics, and phylogenetics: Inferring the historical dynamics of biogeographic distributions, in *Zoogeography of Caribbean Insects*, J.K. Liebherr (ed.), Cornell University Press, Ithaca, New York, pp. 254–69.

Connor, E.F. and Simberloff, D. (1979) The assembly of species communities: Chance or competition? *Ecology*, **60**, 1132–40.

Connor, E.F. and Simberloff, D. (1983) Interspecific competition and species co-occurrence patterns on islands: Null models and the evaluation of evidence. *Oikos*, **41**, 455–65.

Connor, E.F. and Simberloff, D. (1984) Neutral models of species co-occurrence patterns, in *Ecological Communities: Conceptual Issues and the Evidence*, D.R. Strong, Jr., D. Simberloff, L.G. Abele, and A.B. Thistle (eds), Princeton University Press, Princeton, New Jersey, pp. 316–31, 341–3.

Diaconis, P. and Gangolli, A. (1995) Rectangular arrays with fixed margins, in *Discrete Probability and Algorithms*, D.J. Aldous, P. Diaconis, J. Spencer, and J.M. Steele (eds), New York: Springer-Verlag, New York, pp. 15–41.

Diamond, J.M. (1975) Assembly of species communities, in *Ecology and Evolution of Communities*, M.L. Cody and J.M. Diamond (eds), Harvard University Press, Cambridge, Massachusetts, pp. 342–444.

Diamond, J.M. and Gilpin, M.E. (1982) Examination of the ''null'' model of Connor and Simberloff for species co-occurrence on islands. *Oecologia*, **52**, 64–74.

Diamond, J.M. and Marshall, A.G. (1976) Origin of the New Hebridean avifauna. *Emu*, **76**, 187–200.

Diamond, J.M. and Marshall, A.G. (1977) Niche shifts in New Hebridean birds. *Emu*, **77**, 61–72.

Gilpin, M.E. and Diamond, J.M. (1982) Factors contribution to non-randomness in species co-occurrences on islands. *Oecologia*, **52**, 75–84.

Gilpin, M.E. and Diamond, J.M. (1984) Are species co-occurrences on islands non-random, and are null hypotheses useful in community ecology?, in *Ecological Communities: Conceptual Issues and the Evidence*, D.R. Strong, Jr., D. Simberloff, L.G. Abele, and A.B. Thistle (eds), Princeton University Press, Princeton, New Jersey, pp. 297–315, 332–41.

Gotelli, N.J. (2000) Null model analysis of species co-occurrence patterns. *Ecology*, **81**, 2606–21.

Gotelli, N.J. and Entsminger, G.L. (2000) Declining the knight's tour: A re-analysis of Sanderson *et al*. (1998). Ms. submitted to *Oecologia*.

Guyer, C. (1990) The herpetofauna of La Selva, Costa Rica, in *Four Neotropical Rain Forests*, A. Gentry (ed.), Yale University Press, New Haven, Connecticut, pp. 371–85.

Harvey, P.H., Colwell, R.K., Silvertown, J.W., and May, R.M. (1983) Null models in ecology. *Annual Review of Ecology and Systematics*, **14**, 189–211.

Jackson, D.A., Somers, K.M., and Harvey, H.H. (1992) Null models and fish communities: Evidence of nonrandom patterns. *American Naturalist*, **139**, 930–51.

Jerrum, M. (1998) Mathematical foundations of the Markov chain Monte Carlo method, in *Probabilistic Methods for Algorithmic Discrete Mathematics*, no. 16 in *Algorithms and Combinatorics*, M. Habib, C. McDiarmid, J. Ramirez-Alfonsin, and B. Reed (eds), Springer-Verlag, Berlin, pp. 116–65.

Manly, B.F.J. (1995) A note on the analysis of species co-occurrences. *Ecology*, **76**, 1109–15.

Matthews, W.J. (1982) Small fish community structure in Ozark streams: Structured assembly patterns or random abundance of species? *American Midland Naturalist*, **107**, 42–54.

Mosimann, J.E. (1968) *Elementary Probability for the Biological Sciences*, Appleton-Century-Crofts, New York.

Pielou, D.P. and Pielou, E.C. (1968) Association among species of infrequent occurrence: The insect and spider fauna of *Polyporus betulinus* (Bulliard) Fries. *Journal of Theoretical Ecology*, **21**, 202–16.

Pramanik, P. (1994) Generating random (0,1)-matrices with given marginals. Presented at The Third International Conference on Lattice Path Combinatorics, New Delhi.

Putman, R.J. (1994) *Community Ecology*, Chapman & Hall, London.

Rao, A.R., Jana, R., and Bandyopadhyay, S. (1996) A Markov chain Monte Carlo method for generating random (0,1)-matrices with given marginals. *Sankhyā*, **58**(Series A), 225–42.

Reddingius, J. (1983) On species sharing islands: Comment on an article by S.J. Wright and C.C. Biehl. *American Naturalist*, **122**, 830–2.

Roberts, A. and Stone, L. (1990) Island-sharing by archipelago species. *Oecologia*, **83**, 560–7.

Roberts, E. (1986) *Thinking Recursively*, Wiley, New York.

Ryser, H.J. (1960) Matrices of zeros and ones. *Bulletin of the American Mathematical Society*, **66**, 442–64.

Sanderson, J.G. (2000) Testing ecological patterns. *American Scientist*, **88**, 332–9.

Sanderson, J.G., Moulton, M.P., and Selfridge, R.G. (1998) Null matrices and the analysis of species co-occurrences. *Oecologia*, **116**, 275–83.

Simberloff, D. (1986) Analysis of presence/absence data for species on islands: Passerine birds of the Cyclades. *Biologia Gallo-Hellenica*, **12**, 43–68.

Simberloff, D. and Connor, E.F. (1979) Q-mode and R-mode analyzes of biogeographic distributions: Null hypotheses based on random colonization, in *Contemporary Quantitative Ecology and Related Ecometrics*, G.P. Patil and M.L. Rosenzweig (eds), International Co-operative Publishing House, Fairland, Maryland, pp. 123–38.

Simberloff, D. and Connor, E.F. (1981) Missing species combinations. *American Naturalist*, **118**, 215–39.

Simberloff, D. and Connor, E.F. (1984) Inferring competition from biogeographic data: A reply to Wright and Biehl. *American Naturalist*, **124**, 429–36.

Simberloff, D., Heck, K.L., McCoy, E.D., and Connor, E.F. (1981) There have been no statistical tests of cladistic biogeographical hypotheses, in *Vicariance Biogeography: A Critique*, G. Nelson and D.E. Rosen (eds), Columbia University Press, New York, pp. 40–63.

Snapper, E. (1971) Group characters and non-negative integral matrices. *Journal of Algebra*, **19**, 520–35.

Snijders, T.A.B. (1991) Enumeration and simulation methods for 0-1 matrices with given marginals. *Psychometrika*, **56**, 397–417.

Stone, L. and Roberts, A. (1990) The checkerboard score and species distributions. *Oecologia*, **85**, 74–9.

Stone, L. and Roberts, A. (1992) Competitive exclusion, or species aggregation? *Oecologia*, **91**, 419–24.

Sukhatme, P.V. (1938) On bipartitional functions. *Philosophical Transactions of the Royal Society, Series A*, **237**, 375–409.

Verbeek, A. and Kroonenberg, P.M. (1985) A survey of algorithms for exact distributions of test statistics in r by c contingency tables with fixed margins. *Computational Statistics and Data Analysis*, **3**, 159–85.

Wang, B.Y. (1988) Precise number of (0,1)-matrices in A(R,S). *Scientia Sinica Series A*, **31**, 1–6.

Wilson, J.B. (1987) Methods for detecting non-randomness in species co-occurrences: A contribution. *Oecologia*, **73**, 579–82.

Wilson, J.B. (1988) Community structure in the flora of islands in Lake Manapouri, New Zealand. *Journal of Ecology*, **76**, 1030–42.

Wilson, J.B., James, R.E., Newman, J.E., and Myers, T.E. (1992) Rock pool algae: species composition determined by chance? *Oecologia*, **91**, 150–2.

Wright, S.J. and Biehl, C.C. (1982) Island biogeographic distributions: Testing for random, regular, and aggregated patterns of species occurrence. *American Naturalist*, **119**, 345–57.

# Biographical sketches

Arif Zaman is Professor and Dean at the School of Arts and Sciences at the Lahore University of Management Sciences, in Pakistan. He received his Ph.D. in Statistics at Stanford University and has since taught at Purdue University and Florida State University. His interests are computational and discrete puzzles and problems in mathematics and statistics. These include, but are not limited to, Markov chains, exchangeability, graph theory, pseudo-random number generation, and cryptography.

Daniel Simberloff is Nancy Gore Hunger Professor of Environmental Studies at the University of Tennessee. He received his Ph.D. in Biology from Harvard University and taught for many years at Florida State University. His interests are population and community ecology, biogeography, and invasion biology.