

Multiple Regression and Correlation Extensions of the Mantel Test of Matrix Correspondence

PETER E. SMOUSE,¹ JEFFREY C. LONG,² AND ROBERT R. SOKAL³

¹*Department of Human Genetics and Department of Biology,
University of Michigan, Ann Arbor, Michigan 48109;*

²*Department of Anthropology, CUNY-Hunter College, New York, New York 10021; and*

³*Department of Ecology and Evolution, State University of New York,
Stony Brook, New York 11794*

It is often necessary in population biology to compare two sets of distance measures. These measures can be based on genetic markers, morphological traits, geographic separation, ecological divergence, and so on. The distance measures can take various forms and frequently have unknown distributional properties. Many different procedures have been developed to compare the correspondence of one set of distances with another set. Prominent among them are the: (1) matrix correlation techniques of Sokal and Rohlf (1962), and Sneath and Sokal (1973); (2) network-matching techniques of Spielman (1973); (3) matrix dilation and rotation techniques of Gower (1971) and of Schönemann and Carrol (1970); and (4) smallest-space techniques of Lingoes (1965) and Guttman (1968). Each of these strategies has strong points, but all suffer from a difficulty in assessing the statistical significance of attained correspondence. The problem is that a set of all possible pairwise distances between k units (populations, taxa, habitats, etc.) cannot be independent.

More recently, a test of matrix correspondence—originally developed by Mantel (1967) and widely applied in geography (Hubert and Golledge, 1982) and psychometrics (Hubert, 1979a, b)—has caught the attention of population biologists (Sokal, 1979; Sokal et al., 1980; Douglas and Endler, 1982; Dow and Cheverud, 1985; Schnell et al., 1985, 1986; Sokal et al., 1986, 1987; O'Brien, 1987; Smouse and Wood, 1987). Attractions of the Mantel procedure are its wide applicability and computational simplicity. Mantel (1967)

presented a formal analysis of matrix correspondence based on the assumption of asymptotic normality for a particular test criterion. Later workers (Mielke et al., 1981) developed more general procedures for Mantel statistics that assume a Pearson Type III distributional form. The most widely used evaluation procedure, however, involves the construction of a null distribution by Monte Carlo randomization, whereby one of the matrices is held rigid and the other has its rows and corresponding columns randomly permuted (Cliff and Ord, 1981). Dietz (1983) evaluated the Mantel test as one of several permutational tests for association between distance matrices. When dimensions of two matrices are small (say $K \leq 7$), it is customary to evaluate the Mantel test criterion for all $K!$ equally likely permutations. If K is large, then a large number of random permutations are sampled with replacement.

Although the test has been useful in its present form, there are some simple modifications and extensions that would encourage even wider deployment in population work. Our purpose here is to sketch these changes and to illustrate them with an example drawn from the Yanomama Indians of lowland South America.

MODIFYING THE MANTEL TEST

Consider a pair of distance matrices X and Y . The elements of these matrices, X_{ij} and Y_{ij} , represent types of contrasts between the i th and j th units (populations, taxa, habitats, whatever), where $i, j = 1, \dots, K$. Compute the test criterion

$$\tilde{Z}_{YX} = \sum_{ij} (X_{ij}Y_{ij}), \quad (1) \quad \text{and}$$

where \sum_{ij} symbolizes summation over all ij pairs other than $i = j$, and " \sim " indicates the observed value. This test criterion is compared with the expected distribution of Z_{YX} values obtained when the corresponding elements of the two matrices are not associated in any way. Then, using an empirical null distribution derived from Monte Carlo sampling, compute the probability of obtaining a value of Z_{YX} at least as extreme as \tilde{Z}_{YX} by chance alone. The alternative hypothesis is usually that there will be a positive association between the corresponding elements of the two matrices, and we compute the empirical probability of obtaining random Z_{YX} in excess of \tilde{Z}_{YX} , $P = \Pr(Z_{YX} > \tilde{Z}_{YX})$. Thus, P is the upper tail probability for the null distribution of Z_{YX} .

While P is a useful measure of the statistical significance of a specified departure from randomness (Mantel, 1967), Z_{YX} will be an unfamiliar measure whose scale varies from problem to problem. Note, however, that barring a correction factor, Z_{YX} is the sum of cross-products between X_{ij} and Y_{ij} . To see this, first compute the mean values of X_{ij} and Y_{ij} . Most applications employ matrices with zero diagonal elements, since the distance between any object and itself is zero. Moreover, because most distance matrices are also symmetric about this diagonal, there are only $K(K - 1)/2$ distinct entries. To be general, however, we compute:

$$\bar{X} = \sum_{ij} (X_{ij}/N) \quad (2a)$$

and

$$\bar{Y} = \sum_{ij} (Y_{ij}/N), \quad (2b)$$

where $N = K(K - 1)$ is the total number of off-diagonal elements in the matrices X and Y . The corrected sum of products is usually computed as

$$SP(X, Y) = [\tilde{Z}_{YX} - N\bar{X}\bar{Y}]. \quad (3)$$

Corresponding to this corrected sum of products, we also have a pair of corrected sums of squares, one for the elements of each matrix:

$$SS(X) = \sum_{ij} (X_{ij} - \bar{X})^2 \quad (4a)$$

$$SS(Y) = \sum_{ij} (Y_{ij} - \bar{Y})^2. \quad (4b)$$

Note that while $SP(X, Y)$ will change if one of the matrices is permuted, both $SS(X)$ and $SS(Y)$ are invariant under permutation. Combining (3) with (4a) and (4b), we obtain a regression coefficient

$$b_{YX} = SP(X, Y)/SS(X) \quad (5a)$$

and a corresponding correlation coefficient

$$r_{YX} = SP(X, Y)/[SS(X) \cdot SS(Y)]^{1/2}, \quad (5b)$$

showing that the Mantel treatment is really a regression analysis. As in standard regression analysis, we have a linear model of the simple form

$$[Y_{ij} - \bar{Y}] = b_{YX}[X_{ij} - \bar{X}] + \epsilon_{ij}. \quad (6)$$

The resulting measure of matrix correspondence (r_{YX}) is analogous to an autocorrelation coefficient, and is equivalent to a normalization of \tilde{Z}_{YX} . Similar treatments are presented in Hubert and Baker (1978), and Hubert and Subkoviak (1979). Clearly, the translation of \tilde{Z}_{YX} into b_{YX} or r_{YX} is a monotonic mapping in either direction. This normalization has been used in our own recent work (Salzano et al., 1986; Sokal et al., 1986, 1987; Smouse and Wood, 1987) and in that of O'Brien (1987). Because the elements within either of the matrices X and Y are not independent among themselves, the usual significance tests are not valid; however, since we employ a Monte Carlo null distribution, lack of independence is not a problem.

EXTENSIONS

There are occasional situations for which it is helpful to use two or more distance matrices (X_1, X_2, \dots, X_H) to predict the elements of a single "response" matrix (Y). In such situations, it will often be the case that the respective elements of the various X matrices are themselves correlated, so there is a certain redundancy of information. We need to be able to assess how well the individual X matrices predict the Y matrix, how much additional information is provided by the addition of a particular X matrix, given that others have already

been included in the analysis, and so on. It is easy to extend the above regression procedures to the inclusion of several predictive X matrices.

We first consider the treatment for a pair of X matrices; the extension to more than two X matrices will then be obvious. For the expansion of the linear model to a pair of X variables, we replace (6) with

$$[Y_{ij} - \bar{Y}] = b_{Y1}[X_{1ij} - \bar{X}_1] + b_{Y2}[X_{2ij} - \bar{X}_2] + \epsilon_{ij}. \quad (7)$$

Given the three matrices, Y , X_1 and X_2 , we compute the sums of squares

$$SS(X_1) = \sum_{ij} (X_{1ij} - \bar{X}_1)^2, \quad (8a)'$$

$$SS(X_2) = \sum_{ij} (X_{2ij} - \bar{X}_2)^2, \quad (8b)$$

and

$$SS(Y) = \sum_{ij} (Y_{ij} - \bar{Y})^2, \quad (8c)$$

and the sums of cross-products

$$SP(X_1, Y) = [\tilde{Z}_{Y1} - N\bar{X}_1\bar{Y}], \quad (9a)$$

$$SP(X_2, Y) = [\tilde{Z}_{Y2} - N\bar{X}_2\bar{Y}], \quad (9b)$$

and

$$SP(X_1, X_2) = [\tilde{Z}_{12} - N\bar{X}_1\bar{X}_2], \quad (9c)$$

where

$$\tilde{Z}_{Y1} = \sum_{ij} (X_{1ij}Y_{ij}), \quad (10a)$$

$$\tilde{Z}_{Y2} = \sum_{ij} (X_{2ij}Y_{ij}), \quad (10b)$$

and

$$\tilde{Z}_{12} = \sum_{ij} (X_{1ij}X_{2ij}). \quad (10c)$$

We next define a matrix $X'X$, describing variation within and covariation between the X matrices:

$$X'X = \begin{bmatrix} SS(X_1) & SP(X_1, X_2) \\ SP(X_1, X_2) & SS(X_2) \end{bmatrix}, \quad (11)$$

as well as a vector $X'Y$, describing the covariation between each of the X matrices and the Y matrix,

$$X'Y = \begin{bmatrix} SP(X_1, Y) \\ SP(X_2, Y) \end{bmatrix}. \quad (12)$$

The regression coefficients are obtained from the vector equation

$$b = [b_1, b_2]' = (X'X)^{-1}X'Y, \quad (13)$$

which is the analogue of (5a). We obtain the correlation coefficients r_{Y1} , r_{Y2} , and r_{12} by appropriate substitutions of the elements of equations (8a) through (9c) into (5b).

With this treatment, we also obtain all of the usual regression output. The partial regression coefficient of Y on X_1 (for fixed values of X_2) is given by

$$b_{Y1.2} = \frac{SP(X_1, Y) - b_{Y2} \cdot SP(X_1, X_2)}{SS(X_1) - b_{12} \cdot SP(X_1, X_2)}, \quad (14a)$$

while that for Y on X_2 (for fixed values of X_1) is given by

$$b_{Y2.1} = \frac{SP(X_2, Y) - b_{Y1} \cdot SP(X_1, X_2)}{SS(X_2) - b_{21} \cdot SP(X_1, X_2)}. \quad (14b)$$

Similarly, the partial correlation of Y and X_1 (for fixed values of X_2) is given by

$$r_{Y1.2} = \frac{r_{Y1} - r_{Y2}r_{12}}{[(1 - r_{Y2}^2)(1 - r_{12}^2)]^{1/2}}, \quad (15a)$$

while that for Y and X_2 (for fixed values of X_1) is given by

$$r_{Y2.1} = \frac{r_{Y2} - r_{Y1}r_{12}}{[(1 - r_{Y1}^2)(1 - r_{12}^2)]^{1/2}}. \quad (15b)$$

Finally, the coefficient of multiple determination R^2 is given by

$$R^2 = 1 - (1 - r_{Y1}^2)(1 - r_{Y2.1}^2) \\ = 1 - (1 - r_{Y2}^2)(1 - r_{Y1.2}^2). \quad (16)$$

These are the standard results for a pair of X variables, and are described more fully in Sokal and Rohlf (1981:chapter 16). With three or more X variables, additional (but standard) matrix methods are required.

The only tricky aspect is significance testing. The procedure of choice depends on whether X_1 and X_2 are viewed as predictors of Y , or merely as correlated measures. For the usual regression treatment, the X matrices are treated as fixed and measured without error. The point of extending the analysis to a pair of X matrices, rather than treating each separately, is that the two X matrices are themselves not independent. The dependence between X_1 and X_2 is the important feature of the situation, and we need a permutational null

distribution that treats this dependence as fixed. We accomplish this by permuting Y , holding both X_1 and X_2 constant, and re-computing any desired statistics each time. The test criteria of interest can then be evaluated against the corresponding statistics obtained from the Monte Carlo procedures.

On the other hand, when X_1 and X_2 are themselves variable, possibly measured with error, and associated with (rather than predictive of) Y , then it is preferable to denote the three matrices as Y_1 , Y_2 and Y_3 . The partial correlation coefficient ($r_{ij \cdot k}$) should be computed as follows: regress Y_i on Y_k , and compute a matrix $D_{i \cdot k}$ whose elements are the residual deviations $d_{i \cdot k}$. Similarly, compute the residual matrix $D_{j \cdot k}$ with elements $d_{j \cdot k}$, by regressing Y_j on Y_k . The correlation between corresponding elements of $D_{i \cdot k}$ and $D_{j \cdot k}$ is the partial correlation $r_{ij \cdot k}$ of the matrices Y_i and Y_j , given Y_k . The significance of the coefficient can be evaluated by random permutation of one of the residual matrices, while holding the other constant. This approach can also be used to compute partial correlations in which more than one variable is kept constant (e.g., see Sokal and Thomson, 1987).

The formulation employed above extracts only the linear component of the association between pairs of variables. Non-linearities can be dealt with by means of suitable transformations, and some attention can profitably be directed toward the choice of input measures and their scaling. Given that the Mantel approach is frequently used when there is little a priori knowledge of the precise functional relations among variables, a certain amount of trial and error in the definition of variables is probably not out of place in most applications. One can always, of course, reduce the elements of the various distance matrices to ranks, thus estimating and testing a series of rank correlations. However, one is then implicitly extracting the linear component of the associations between two or more sets of ranks. The real value of the Mantel procedure is that it provides a useful and valid test, even

when little is known about the elements of the distance matrices.

Two alternative extensions of the Mantel test to three matrices have recently been proposed. Based on earlier work by Hubert and Golledge (1981), Dow and Cheverud (1985) computed a Mantel test for the matrices Y and $(X_1 - X_2)$, after first comparably scaling X_1 and X_2 . The sign of a statistically significant association demonstrates whether X_1 or X_2 has the greater influence on Y . In a second application, Hubert (1985) computed a Mantel test for the matrices Y and $X_1 X_2$, where $X_1 X_2$ is the Hadamard (element-by-element) product of X_1 and X_2 . The question posed by Hubert was whether Y is significantly influenced by the interaction of X_1 and X_2 . Since the alternative procedures discussed above address rather different questions, both can be viewed as complementary to the regression/correlation approach advocated here. Each reduces the problem of three matrices to that of a pair of matrices, using the traditional \bar{Z} measure. All three procedures can be extended to multiple matrices (see Hubert and Golledge, 1981), although that extension is a bit cumbersome for either the Hubert (1985) or the Dow and Cheverud (1985) approach. By contrast, extension of the regression/correlation approach advocated here is both straightforward and familiar. Moreover, it can be extended to both multiple X and multiple Y matrices, capitalizing on the full panoply of traditional least squares methodology.

AN ILLUSTRATION WITH THE YANOMAMA

We have recently completed an analysis of the pattern of genetic distances among a set of 50 Yanomama villages from southern Venezuela and northern Brazil (Sokal et al., 1986). As predictors of the genetic distance matrix Y , we used two different X matrices: the first based on geographic distances (X_1), the second based on linguistic affinity and fission history (X_2). An association between Y and X_1 is interpretable in terms of clinal patterns of gene flow among villages; an association between Y and X_2 is ascribable to the fission history

of the tribe. The two X matrices are themselves strongly associated, by virtue of the fact that the fission history of the tribe has a marked geographic component. It is important to determine how much of the geographic pattern of the genetic distances is accounted for by fission history. Conversely, some of the current hierarchical subdivision within the tribe may be more a consequence of geographically structured gene flow than of fission history.

We discovered that both X matrices are useful predictors of the Y matrix, and that the pattern of genetic distance among villages has both clinal and fission components. The correlation coefficients show that geographic proximity ($r_{Y1} = 0.415$, $P < 0.01$) is more highly correlated with genetic distance than is fission history ($r_{Y2} = 0.365$, $P < 0.01$). The interesting feature is that geographic distance and fission history distance are themselves highly correlated ($r_{12} = 0.773$, $P < 0.01$). The partial correlation of geography and genetic distance is significant ($r_{Y1.2} = 0.224$, $P < 0.05$), while that of fission history and genetic distance is not ($r_{Y2.1} = 0.077$, $P > 0.05$). Once having fit geographic distance ($R^2 = r^2_{Y1} = 0.172$), there was almost nothing to be gained by adding fission distance, since the coefficient of multiple determination was only increased to $R^2 = 0.177$. On the other hand, having fit fission distance first ($R^2 = r^2_{Y2} = 0.133$), some additional information was obtained by adding the geographic locations of the villages in question ($R^2 = 0.177$). These results stem from the fact that fission history is geographically structured. In an earlier analysis with log-linear regression procedures, Ward and Neel (1976) described clinal variation within the Yanomama and ascribed it to the geographic pattern of tribal fission and expansion. They did not test the hypothesis, but the results presented here and at greater length by Sokal et al. (1986) confirm that earlier conjecture.

ACKNOWLEDGMENTS

This paper is Contribution No. 598 in Ecology and Evolution from the State University of New York at

Stony Brook. The authors would like to thank R. Chakraborty, J. Cheverud, M. Dow, N. Oden, A. Templeton, J. Wood, and a pair of anonymous reviewers for many helpful comments. The computations were carried out by M. C. Wooten, who developed a computer program MULTMAN to this end. That program has been included in the spatial-analysis package available from Daniel Wartenberg. PES was supported by NIH-R01-GM-32589 and DE-86-ER60089, JCL by NIH-GM-T32-07544, and RRS by NIH-R01-GM-28262.

REFERENCES

- CLIFF, A. D., AND J. K. ORD. 1981. *Spatial processes: Models and applications*. Pion, London.
- DIETZ, E. J. 1983. Permutation tests for association between two distance matrices. *Syst. Zool.*, 32:21-26.
- DOUGLAS, M. E., AND J. A. ENDLER. 1982. Quantitative matrix comparisons in ecological and evolutionary investigations. *J. Theor. Biol.*, 99:777-795.
- DOW, M. M., AND J. M. CHEVERUD. 1985. Comparison of distance matrices in studies of population structure and genetic microdifferentiation: Quadratic assignment. *Am. J. Phys. Anthropol.*, 68:367-373.
- GOWER, J. C. 1971. Statistical methods of comparing different multivariate analyses of the same data. Pages 138-149 in *Mathematics in the archaeological and historical sciences* (F. R. Hodson, D. G. Kendall, and P. Tautu, eds.). Edinburgh Univ. Press, Edinburgh.
- GUTTMAN, L. 1968. A general nonmetric technique for finding the smallest coordinate space for a configuration of points. *Psychometrika*, 33:469-506.
- HUBERT, L. J. 1979a. Comparisons of sequences. *Psychol. Bull.*, 56:1098-1106.
- HUBERT, L. J. 1979b. Matching models in the analysis of cross-classifications. *Psychometrika*, 44:21-41.
- HUBERT, L. J. 1985. Combinatorial data analysis: Association and partial association. *Psychometrika*, 50:449-467.
- HUBERT, L. J., AND F. B. BAKER. 1978. Evaluating the conformity of sociometric measurements. *Psychometrika*, 43:31-41.
- HUBERT, L. J., AND R. G. GOLLEDGE. 1981. A heuristic method for the comparison of related structures. *J. Math. Psychol.*, 23:214-226.
- HUBERT, L. J., AND R. G. GOLLEDGE. 1982. Measuring association between spatially defined variables: Tjøstheim's index and some generalizations. *Geogr. Anal.*, 14:273-278.
- HUBERT, L. J., AND M. J. SUBKOVIAK. 1979. Confirmatory inference and geometric models. *Psychol. Bull.*, 86:361-370.
- LINGOES, J. C. 1965. An IBM 7090 program for Guttman-Lingoes smallest space analysis—I. *Behav. Sci.*, 10:183-184.
- MANTEL, N. A. 1967. The detection of disease clustering and a generalized regression approach. *Cancer Res.*, 27:209-220.

- MIELKE, P. W., K. J. BERRY, AND G. W. BRIER. 1981. Application of multi-responses: Permutation procedures for examining seasonal changes in monthly mean sea-level pressure patterns. *Month. Weather Rev.*, 109:120-126.
- O'BRIEN, E. 1987. The correlation between population structure and genetic structure in the Hutterite population. In *Mammalian dispersal patterns: The effects of social structure on population genetics* (B. D. Chepko-Sade and Z. Halpin, eds.). Univ. Chicago Press, Chicago (in press).
- SALZANO, F. M., H. GERSHOWITZ, H. MOHRENWEISER, J. V. NEEL, P. E. SMOUSE, M. A. MESTRINER, T. A. WEIMER, M. H. L. P. FRANCO, A. L. SIMÕES, J. CONSTANS, A. E. OLIVEIRA, AND M. J. DE MELO E FREITAS. 1986. Gene flow across tribal barriers and its effect among the Amazonian Içana River Indians. *Am. J. Phys. Anthropol.*, 69:3-14.
- SCHNELL, G. D., D. J. WATT, AND M. DOUGLAS. 1985. Statistical comparison of proximity matrices: Applications in animal behavior. *Anim. Behav.*, 33: 239-253.
- SCHNELL, G. D., M. E. DOUGLAS, AND D. J. HOUGH. 1986. Geographic patterns of variation in offshore spotted dolphins (*Stenella attenuata*) of the eastern tropical Pacific Ocean. *Mar. Mamm. Sci.*, 2:186-213.
- SCHÖNEMANN, P. H., AND R. M. CARROL. 1970. Fitting one matrix to another under choice of a central dilation and a rigid motion. *Psychometrika*, 35: 245-255.
- SMOUSE, P. E., AND J. W. WOOD. 1987. The genetic demography of the Gainj of Papua New Guinea. III. Functional models of migration and their genetic implications. In *Mammalian dispersal patterns: The effects of social structure on population genetics* (B. D. Chepko-Sade and Z. Halpin, eds.). Univ. Chicago Press, Chicago (in press).
- SNEATH, P. H. A., AND R. R. SOKAL. 1973. Numerical taxonomy. Freeman, San Francisco.
- SOKAL, R. R. 1979. Testing statistical significance of geographic variation patterns. *Syst. Zool.*, 28:227-231.
- SOKAL, R. R., J. BIRD, AND B. RISKA. 1980. Geographic variation in *Pemphigus populicaulis* (Insecta: Aphididae) in eastern North America. *Biol. J. Linn. Soc.*, 14:163-200.
- SOKAL, R. R., N. L. ODEN, AND J. F. S. BARKER. 1987. Spatial structure in *Drosophila buzzatii* populations: Simple and 2-dimensional spatial autocorrelation. *Am. Nat.* (in press).
- SOKAL, R. R., AND F. J. ROHLF. 1962. The comparison of dendrograms by objective methods. *Taxon*, 11: 33-40.
- SOKAL, R. R., AND F. J. ROHLF. 1981. *Biometry*. Second edition. Freeman, San Francisco.
- SOKAL, R. R., P. E. SMOUSE, AND J. V. NEEL. 1986. The genetic structure of a tribal population, the Yanomama Indians. XV. An effort to find pattern by autocorrelation analysis. *Genetics*, 114:259-287.
- SOKAL, R. R., AND J. D. THOMSON. 1987. Applications of spatial autocorrelation in ecology. In *Numerical ecology* (P. Legendre, ed.). Springer-Verlag, New York (in press).
- SPIELMAN, R. S. 1973. Differences among Yanomama Indian villages: Do the patterns of allele frequencies, anthropometrics and map locations correspond? *Am. J. Phys. Anthropol.*, 39:461-480.
- WARD, R. H., AND J. V. NEEL. 1976. The genetic structure of a tribal population, the Yanomama Indians. XIV. Clines and their interpretation. *Genetics*, 82:103-121.

Received 23 June 1986; accepted 23 July 1986.