



Promoting the Science of Ecology

Confronting Multicollinearity in Ecological Multiple Regression

Author(s): Michael H. Graham

Source: *Ecology*, Vol. 84, No. 11 (Nov., 2003), pp. 2809-2815

Published by: Ecological Society of America

Stable URL: <http://www.jstor.org/stable/3449952>

Accessed: 20/11/2009 09:35

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/action/showPublisher?publisherCode=esa>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



Ecological Society of America is collaborating with JSTOR to digitize, preserve and extend access to *Ecology*.

<http://www.jstor.org>

STATISTICAL REPORTS

Ecology, 84(11), 2003, pp. 2809–2815
© 2003 by the Ecological Society of America

CONFRONTING MULTICOLLINEARITY IN ECOLOGICAL MULTIPLE REGRESSION

MICHAEL H. GRAHAM¹

Moss Landing Marine Laboratories, 8272 Moss Landing Road, Moss Landing, California 95039 USA

Abstract. The natural complexity of ecological communities regularly lures ecologists to collect elaborate data sets in which confounding factors are often present. Although multiple regression is commonly used in such cases to test the individual effects of many explanatory variables on a continuous response, the inherent collinearity (multicollinearity) of confounded explanatory variables encumbers analyses and threatens their statistical and inferential interpretation. Using numerical simulations, I quantified the impact of multicollinearity on ecological multiple regression and found that even low levels of collinearity bias analyses ($r \geq 0.28$ or $r^2 \geq 0.08$), causing (1) inaccurate model parameterization, (2) decreased statistical power, and (3) exclusion of significant predictor variables during model creation. Then, using real ecological data, I demonstrated the utility of various statistical techniques for enhancing the reliability and interpretation of ecological multiple regression in the presence of multicollinearity.

Key words: *confounding factors; multicollinearity; multiple regression; principal components regression; sequential regression; structural equation modeling.*

INTRODUCTION

Ecologists often use multiple regression to develop models that describe the regulation of particular aspects of organismal, population, and community ecology (dependent or response variables) by various environmental and biological factors (independent or explanatory variables) (James and McCulloch 1990). Multiple regression analyses, however, can be hindered by the complex nature of ecological data, in which targeted ecological responses are linked to many explanatory variables that are often correlated among each other (multicollinear). Multicollinear explanatory variables are difficult to analyze because their effects on the response can be due to either true synergistic relationships among the variables or spurious correlations. Ecologists often counter by designing experimental studies that break correlations among explanatory variables and analyzing the data with analyses of variance (ANOVA) that allow for the isolation of main effects and higher-order interactions among individual explanatory variables (Scheffe 1959). In practice, however, ecological explanatory variables are often not under experimental control, in which case the explanatory variables of interest may be correlated. It is under these

conditions that multiple regression is often used to analyze ecological data (James and McCulloch 1990).

The statistical and inferential problems of multicollinearity in multiple regression have been well established in the statistical literature (e.g., Cohen and Cohen 1983, Hocking 1996, Neter et al. 1996, Tabachnick and Fidell 1996, Draper and Smith 1998, Chatterjee et al. 2000), although problems specific to ecological data have rarely been discussed (James and McCulloch 1990, Phillipi 1993, Legendre and Legendre 1998; and see Mitchell-Olds and Shaw 1987 and Petraitis et al. 1996 for related discussions of fitness regression and path analysis, respectively). Yet, despite previous warnings by statisticians, only 32 of 294 (11%) papers published in *Ecology*, *Ecological Monographs*, *Functional Ecology*, *Journal of Animal Ecology*, and *Journal of Ecology* from 1993 to 1999 that used multiple regression for data analysis even discussed the potential presence of multicollinearity. Of these 32 papers, only 17 (53%) actually tested whether multicollinearity was present; of these 17 papers, 11 (65%) found significant multicollinearity, suggesting that ecological data are typically collinear. But how desperate is the problem for ecologists? The goal of this paper was twofold: (1) to quantify through numerical simulation the statistical and inferential biases caused when multicollinearity is present in multiple regression analyses; and (2) to demonstrate the utility of various statistical techniques for enhancing the reliability and interpre-

Manuscript received 26 August 2002; revised 25 March 2003; accepted 26 March 2003; final version received 23 June 2003.
Corresponding Editor: A. M. Ellison.

¹ E-mail: mgraham@mlml.calstate.edu

tation of ecological multiple regression in the presence of multicollinearity.

THEORETICAL PROBLEMS AND EMPIRICAL CONSEQUENCES

In multiple linear regression, data are fit to a linear model that predicts values of a response (Y) as the weighted sum of explanatory variables (X_i) and random error (ε): $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_i X_i + \varepsilon$, where β s are regression coefficients. The typical goal is to build a model using the fewest variables to explain the greatest variability in the response, and to accurately parameterize regression coefficients for those variables. If all explanatory variables are independent of each other, each regression coefficient represents the total contribution of a given predictor to the response. If, however, two or more variables are collinear to any extent, partial regression coefficients need to be calculated to isolate the unique contribution of a particular explanatory variable (hereafter the predictor) from that shared with other variables (hereafter confounders). This unique contribution is the extra sums of squares. The distinction between unique and shared contributions is the crux of multiple regression's statistical and inferential problems due to multicollinearity.

When data are standardized to a mean of zero and unit variance, the partial regression coefficient for a predictor in the presence of a single confounder is defined as: $\beta = (r_{Y1} - r_{Y2}r_{12})/(1 - r_{12}^2)$, where r_{Y1} is the correlation between the response and predictor, r_{Y2} is that between the response and confounder, and r_{12} and r_{12}^2 are the correlation and coefficient of determination between the predictor and confounder (Neter et al. 1996); β reduces to $\beta^* = r_{Y1}$ in the absence of multicollinearity (i.e., r_{12} and $r_{12}^2 = 0$). As such, partial regression coefficients decrease nonlinearly with increasing multicollinearity (as shown by Petraitis et al. 1996) and deviations from β^* will occur in the presence of even the weakest multicollinearity (i.e., $\beta < \beta^*$ at all $r_{12}^2 > 0$). The marginal statistics used to test the significance of β (i.e., $H_0: \beta = 0$), which is typically used as a criterion to determine whether a given predictor is to be included in a model, is defined as $t = \beta/\text{SE}(\beta)$ (or $t = (r_{Y1} - r_{Y2}r_{12})/\sqrt{\text{MS}_{\text{residual}}}$). Here, $\text{SE}(\beta)$ is the standard error of the coefficient which increases linearly with increasing r_{12}^2 (Neter et al. 1996). Power to detect an effect as significant will therefore also decrease nonlinearly with increasing multicollinearity.

If, during stepwise variable selection, a predictor is ultimately excluded from a model due to its low apparent significance, regression coefficients and marginal statistics of the other variables will change (Mitchell-Olds and Shaw 1987, Philippi 1993, Neter et al. 1996, Petraitis et al. 1996). The use of stepwise variable selection procedures that rely on calculation of marginal statistics may even exclude explanatory variables that are actually highly correlated with the response (i.e., decrease statistical power). Furthermore,

although statistical significance and fit of a final model are not directly affected by multicollinearity (expected sums of squares and marginal statistics are not computed), interpretation of the model may be uncertain due to biased parameterization of partial regression coefficients for individual explanatory variables. Not only will the sum of r^2 for individual predictors generally differ from the R^2 of the final model, actual application of the final model to predict future values for the response can be grossly inaccurate, since none of the partial regression coefficients reflect shared contributions (Tabachnick and Fidell 1996).

These statistical difficulties in analyzing ecological data in the presence of multicollinearity were illustrated numerically by calculating marginal t statistics and P -values for a predictor in the presence of a single confounder (Fig. 1). The purpose of the simulation was to estimate the level of multicollinearity that would result in the erroneous exclusion of significant predictors from a final model. In general, (1) apparent significance (P or apparent α) decreased rapidly with increasing multicollinearity; (2) weak predictors were more vulnerable to erroneous exclusion than strong ones; (3) predictors with high true significance became more vulnerable to erroneous exclusion as the correlation between the response and confounder (r_{Y2}) increased; and, (4) even if correlations between the response and confounders were relatively weak, low levels of multicollinearity (i.e., $r_{12} \geq 0.28$ or $r_{12}^2 \geq 0.08$) resulted in significant predictors appearing insignificant.

To illustrate the negative impact of these statistical biases on the reliability and interpretation of ecological multiple regression, data were reanalyzed from a study of the effect of various environmental factors (wave orbital displacement, wave breaking depth, wind velocity, and mean tidal height) on the shallow (upper) distributional limit of the subtidal kelp *Macrocystis pyrifera* (Graham 1997). The overall severity of multicollinearity in these data was moderate, as wave orbital displacement, wave breaking depth, and wind velocity were strongly correlated among each other ($r \geq 0.6$; $\text{VIF} \geq 2$), but tidal height was only weakly correlated with the other variables ($r < 0.4$; $\text{VIF} = 1.17$). Although Neter et al. (1996) and Chatterjee et al. (2000) suggested that multicollinearity is only severe at $\text{VIFs} > 10$, it is clear from Fig. 1 that VIFs as low as 2 can have significant impacts (see also Petraitis et al. 1996). When analyzed using separate linear regressions, all of the explanatory variables were significant or marginally significant predictors of the response (i.e., $P \leq 0.1$; Table 1). Backwards stepwise multiple regression, however, suggested that only wave orbital displacement and wind velocity were important (Table 1; Appendix A); forward selection yielded the same final model. Partial regression coefficients (β in standard regressions; Table 1) were often more than 1 SE lower than the nonpartial regression coefficients (β in

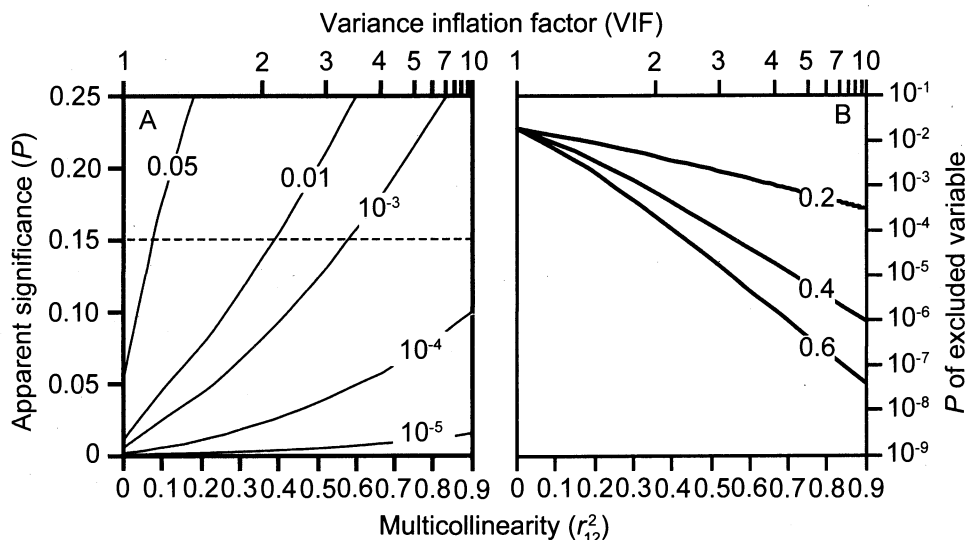


FIG. 1. (A) Effect of multicollinearity on predictor apparent significance (P or apparent α) in the presence of a single confounder. Multicollinearity was represented by r_{12}^2 and variance inflation factors ($VIF = 1/(1 - *R^2)$; $*R^2$ is the R^2 when explanatory variable i is regressed on all other variables in model). Values of r_{y1} were chosen to provide specific levels of “true” significance in the absence of multicollinearity (α ; given on each line). MS_{residual} equaled 0.1 for 35 degrees of freedom and was taken from real standardized data. Predictors with apparent P values that increased to ≥ 0.15 were considered to be negatively affected by multicollinearity; all true P values were ≤ 0.05 in this exercise. (B) Effect of multicollinearity on the exclusion of a significant predictor. The y-axis is the true significance of predictors (expressed as P) that would have been excluded during variable selection. Data were obtained by setting $t = 1.47$ (corresponding to $P = 0.15$ for 35 df) and solving for “true” significance (r_{y1}) as a function of various levels of multicollinearity (r_{12}^2 or VIF) and confounder strength (r_{12} ; given on each line).

simple regressions; Table 1), reflecting the omission of variability in the response shared among predictors. Thus, although wave-breaking depth was initially identified as important (Table 1), this was due almost entirely to variability shared with wave orbital displacement. Many would argue that the removal of wave-breaking depth was therefore necessary because it was a redundant variable, however, there was no evidence that wave-breaking depth wasn't the variable functionally responsible for the shared contribution. Clearly, for two highly collinear explanatory variables that have a strong shared contribution to the response, the decision as to which is the most important predictor, and should therefore be retained, is very ambiguous.

SOME OLD AND NOT-SO-OLD SOLUTIONS

If the entire goal of conducting a multiple regression analysis is to develop a model that best predicts variability in the response, and there is no interest in studying particular relationships between the response and explanatory variables, then the problems due to multicollinearity can be effectively ignored (i.e., “the proof is in the pudding” scenario). In most ecological studies, however, researchers are interested in examining the effects of particular explanatory variables, in which case various techniques are available for addressing the statistical pitfalls of multicollinearity. One approach is to avoid or stabilize the use of marginal statistics for variable selection. The easiest way to do

this is to simply drop collinear variables from analysis (Philippi 1993, Legendre and Legendre 1998). Variable exclusion, however, ignores the unique contribution of the omitted variable and can result in a substantial loss of explanatory power (Carnes and Slade 1988, James and McCulloch 1990) as well as inferential problems in choosing which variables should remain (Mitchell-Olds and Shaw 1987). Another method is to avoid using marginal statistics during variable selection by pre-determining model composition (a priori modeling). This circumvents the problem of choosing which collinear variables should be excluded. In the absence of a reasonable a priori model, marginal statistics can also be avoided by using an “all possible subsets” method of analysis (Furnival 1971). F statistics and coefficients of determination are calculated for all possible combinations (subsets) of variables, and the subset with the greatest fit is identified as “best” using adjusted R^2 (or Akaike's Information Criteria, Mallows's C_p , PRESS, MSE, etc.; Neter et al. 1996). Since distinctions are not made between unique and shared contributions, all possible subsets analyses can help to identify reliably the final model that explains the most variability in the response, although the number of potential subsets can become analytically untreatable as the number of variables increases. An alternative to avoiding marginal statistics is to stabilize them using ridge regression, in which a constant is applied to the elements of the correlation matrix so that it is displaced from singularity,

TABLE 1. Simple linear regressions, and final models from standard multiple regression, sequential regression, and principal components regression (final models are after removal of insignificant explanatory variables; P values ≥ 0.15).

Method and variable	β	SE	t	$F_{2,35}$	P	r^2
Simple						
Wave orbital displacement	0.194	0.028	...	43.58	<0.001	0.55
Wave breaking depth	0.072	0.017	...	17.41	<0.001	0.33
Wind velocity	0.018	0.003	...	29.15	<0.001	0.45
Tidal height	-0.358	0.191	...	3.53	0.068	0.09
Standard						
Wave orbital displacement	0.139	0.038	3.62	...	0.001	0.55
Wind velocity	0.008	0.004	2.10	...	0.043	0.45
Total	26.06	<0.001	0.60
Sequential						
Wave orbit. displ. (1st prior.)	0.194	0.028	6.91	...	<0.001	0.55
Wind velocity (2nd prior.)	0.008	0.004	2.09	...	0.043	0.05
Total	26.06	<0.001	0.60
Principal components						
Principal component 1	0.157	0.024	6.66	...	<0.001	0.54
Principal component 4	-0.039	0.024	-1.69	...	0.100	0.03
Total	23.62	<0.001	0.57

Notes: Model intercepts were significant for all models ($P < 0.0001$) and are not given. In standard and principal components regressions, r^2 values represent total contributions, whereas in the sequential regression, r^2 values represent either unique or unique plus shared contributions as determined by assigned priorities.

increasing the precision of the coefficients (Birkes and Dodge 1993). A problem with all of these methods is that they still require the use of marginal statistics to estimate regression coefficients or determine the relative importance of individual explanatory variables, and thus offer no refuge from associated biases due to multicollinearity.

A more purposeful approach to solving the problems due to multicollinearity is to explore the functional nature of the collinearities, rather than avoid them. This requires methods for identifying and parameterizing the unique and shared contributions of explanatory variables to a response. Here, I used the kelp forest example data to illustrate how three such methods (residual/sequential regression, principle components regression, and structural equation modeling) can improve the reliability and interpretation of ecological multiple regression in the presence of multicollinearity.

Residual and sequential regression

When multicollinearity is limited to pairs of explanatory variables, the easiest way to disentangle unique from shared contributions is simply to assume that one variable is functionally more important than the other, assign the more important variable priority over the shared contribution, and ignore the shared contribution when analyzing the less important variable. This can be done by regressing the less important variable against the other, and replacing the less important variable with the residuals from the regression (see for example, Graham 1997). Priorities can be based on a researcher's own instincts and intuition, previously collected data, data currently under analysis, or the results

of prior experiments that estimated the relative importance of one factor over another. Subsequent multiple regression analyses (residual regressions) will be unbiased since the explanatory variables are no longer statistically collinear. As multicollinearity among explanatory variables becomes more complicated, a modification of sequential regression (or hierarchical regression) can be used. Here it is also assumed that some variables are functionally more important than others, but fixed priorities are assigned to shared contributions for all variables in the model (Tabachnick and Fidell 1996). Marginal statistics are computed for variables in order of highest to lowest priority, with any given variable's marginal statistics ignoring variability already explained by higher priority variables. As such, the rank (order) of marginal statistics remains constant as variables are added or removed from the model, and the decision as to whether a particular variable should remain in the model does not depend on the presence of other variables. Furthermore, both unique and shared contributions are represented in the final parameterized model by the regression coefficients and coefficients of determination. The major concern when using these methods is whether the assigned priorities are relevant to the true functional importance of the variables, and thus, it is vital that researchers are critical of the criteria used to assign priorities.

The final model from a sequential regression analysis of the example data is presented in Table 1, where priorities were based on the unique contributions of each explanatory variable. Regression coefficients and the rank of marginal statistics were constant for each variable selection step (Appendix B) and confirmed

that, by assigning fixed priorities, the decision as to whether a particular variable should remain in the model does not depend on the presence of other variables and model composition is not affected by the use of marginal statistics. It was concluded from this analysis that the unique contribution of wave orbital displacement, plus its shared contribution with winds, was the most important predictor of the response, but that the unique contribution of winds was also important (Graham 1997). Note that, although the standard and sequential multiple regressions yielded the same final models, with sequential regression analyses both unique and shared contributions are represented by the regression coefficients and coefficients of determination, and the individual r^2 values summed to R^2 .

Principal components regression

Alternatively, in principal components regression, it is not generally believed that multicollinearity can be understood best by a hierarchical assignment of priorities, but that collinearities indicate the presence of underlying (latent) variables that are responsible for the shared contributions (Tabachnick and Fidell 1996). A principal components analysis is done on the explanatory variables that identify vectors (i.e., the linear combinations of variables) that account, successively, for the greatest variation in the observations of the explanatory variables; the principal components analysis is done in complete disregard of observed variability in the response. Scores of the orthogonal principal components are used as explanatory variables in a subsequent multiple regression analysis (Philippi 1993, Tabachnick and Fidell 1996, Legendre and Legendre 1998). Since principal components are orthogonal, partial regression coefficients and the rank of marginal statistics do not fluctuate as variables are added or removed and the results of principal components regression will be stable regardless of the severity of multicollinearity. Given that variable selection is unbiased in principal components regression, all principal components can and should be included during variable selection, avoiding the concerns of Mitchell-Olds and Shaw (1987) that explanatory power may be lost by limiting analyses to only those variables with high eigenvalues. The primary limitation of principal components regression lies in the biological interpretation of the principal components.

A principal components analysis was performed on the example data (Appendix C). PC1 accounted for 64% ($\lambda = 2.57$) of the variability among the variables, with wave orbital displacement, wave breaking depth, and wind velocity loading heavily and positively on this PC (all loadings ≥ 0.86); mean tidal height loaded moderately and negatively (loading = -0.54). PC1 thus represented high wave intensity, high wind velocity, and low tide height, or the occurrence of storms during low tides (see Graham [1997] for a detailed biological interpretation of these data). PC2 explained

only 20% of the variability ($\lambda = 0.81$) and appeared to represent mostly tides (loading = 0.84; all others ≤ 0.26). PC3 explained less than 10% of the variability ($\lambda = 0.37$) and primarily represented wind activity (loading = 0.49; all others ≤ 0.19). PC4 explained $\sim 6\%$ of the variability ($\lambda = 0.26$) and represented differences in the two estimates of wave intensity (OD and BD loaded -0.39 and 0.29 respectively; all others ≤ 0.13). The subsequent principal components regression confirmed the stability of regression coefficients and marginal statistics (Appendix C) and that individual r^2 values also summed to the total R^2 for the final model (Table 1). Not surprisingly, the PC that represented the occurrence of storms (PC1) explained the greatest amount of variation in the response. The importance of winds (PC3), however, was not emphasized in the principal components regression. Instead, PC4 was retained suggesting the importance of distinguishing between different aspects of wave intensity, despite the fact that PC4 explained only $\sim 6\%$ of the variability among explanatory variables. That the sequential and principal components regression analyses yielded different results when applied to identical data highlights the importance of determining whether latent variables are likely driving variability in the measured explanatory variables.

Structural equation modeling

Like residual/sequential regression and principal components regression, in structural equation modeling (SEM), it is generally assumed that the best functional multiple regression model is one that can account for both unique and shared contributions. Moreover, like a priori modeling, SEM does not simply explore data to search for relationships between the response and explanatory variables, but rather sets out to test and parameterize hypothesized relationships among the variables. As such, SEM can be used to develop accurate and meaningful final multiple regression models when collinearities among explanatory variables are thought to be present (Hayduk 1987, Loehlin 1987, Bollen 1989, Bentler 1995, Ullman 1996, Shipley 1999). Hypothetical causal links among variables (both unique and shared contributions) are specified and structural equations (models) are developed that represent each potential combination of links. Regression coefficients are then parameterized simultaneously for each link of each model (Bentler 1995, Ullman 1996) and the overall fit of the models are compared as with "all possible subsets" techniques. In its generalized form, SEM directly incorporates latent variables into its models that can represent shared contributions (Ullman 1996; for ecological examples see Brown and Weis 1995, Bishop and Schemske 1998, Gough and Grace 1999), and thus avoids many of the problems identified by Petratis et al. (1996) for path analysis. Still, the successful application of SEM to ecological data is vulnerable to inferential errors made during model development and selection (Ullman 1996,

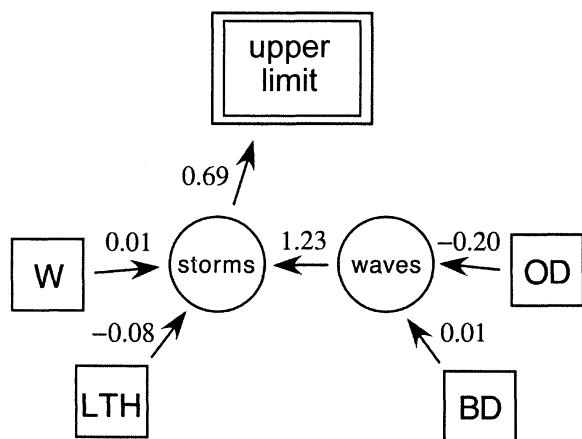


FIG. 2. A structural equation model representing the relationships among four measured explanatory variables (wave orbital displacement [OD], wave breaking depth [BD], mean tidal height [LTH], and wind velocity [W]), two latent variables (storm intensity and wave activity), and the response variable (giant kelp shallow limit). Arrows depict the proposed links between variables. Parameterized regression coefficients are associated with each link. The coefficients were parameterized using iterative normal-theory maximum likelihood available with EQS 6 for Windows (Multivariate Software, Encino, California, USA). The latent variables were developed using the covariance matrix and varimax rotation, and initiated using adjusted principal components according to Bentler (1995).

Shiple 1999): for example, alternate models may exist that differ greatly in the form of their hypothetical causal links, yet may explain similar amounts of variability in the response.

An SEM was developed for the example data, representing one potential relationship between the four predictor variables (wave orbital displacement, wave breaking depth, mean tidal height and wind velocity) and the response (giant kelp shallow limit) (Fig. 2; Appendix D). It was hypothesized that two latent variables were important in driving variability in the response. The first structural equation specified that the latent variable wave intensity could be estimated by a linear combination of wave orbital displacement and wave breaking depth. The second structural equation specified that the latent variable storm intensity could be estimated by a linear combination of wind velocity, mean tidal height, and the latent variable wave intensity. The final structural equation simply specified the linear relationship between the latent variable storm intensity and the response ($R^2 = 0.59$). Again, the results of the parameterized SEM support the conclusions of the sequential and principal components regressions, identifying the underlying importance of storm activity during low tides in driving variability in giant kelp upper limits. Furthermore, by including latent variables into the model, various unique and shared contributions among explanatory variables were explicitly parameterized. However, although R^2 was almost identical among the various

methods (i.e., 0.59–0.60), the adjusted R^2 was in fact lower for the SEM (0.51) than the sequential (0.57) and principal components (0.55) regressions, due to the greater number of SEM regression coefficients that needed to be parameterized. Thus, although the incorporation of latent variables adds flexibility during model development, SEM may not provide the greatest explanatory power for all data analyses.

Post-analysis

The application of one of the above techniques should not be considered the final step in analysis of collinear data. First, each technique demands the standard set of parametric assumptions: normality, constant variance, and independence of error terms. As such, thorough analysis of model residuals should always follow the application of multiple regression techniques. Some techniques (e.g., principal components analysis) additionally require (1) nonsingular matrices of the correlation-covariance among explanatory variables, and (2) that the number of observations of the response greatly exceeds the number of explanatory variables (Tabachnick and Fidell 1996). Second, the generality of estimated regression coefficients should be validated against data that are collected independently of those used during model parameterization. Such validation procedures may also be useful for assessing whether a given multiple regression technique offers the greatest explanatory power. Finally, structural equation modeling and residual, sequential, and principal components regression all deal with shared vs. unique variance contributions differently, and therefore provide diverse perspectives as to the nature of the underlying multicollinearity. As such, ecologists will likely find it most useful to explore multicollinear data with a combination of techniques.

CONCLUSION

This study has quantitatively shown that statistical and inferential problems created by multicollinearity can be extremely severe under realistic ecological conditions. Although straightforward techniques exist for diagnosing and remediating the effects of multicollinearity in multiple regression, they are not commonly utilized in ecology. Still, most of these procedures only help to stabilize the statistical analyses, making them less biased, less subjective, and more repeatable, but only the statistical collinearity will have been removed from the data. The explanatory variables are still, by nature and in nature, correlated, whether or not functionally. Aside from designing manipulative experiments to break correlations among explanatory variables, no technique exists that allows researchers to infer the different functional relationships between the response and explanatory variables. Experiments, however, cannot be applied under all field situations and are especially difficult during the exploratory stage of data collection and model development. It is then that

the determination of relative importance of individual explanatory variables via sampling, and thus a distinction between unique and shared variance contributions, becomes important. The suite of techniques described herein compliment each other and offer ecologists useful alternatives to standard multiple regression for identifying ecologically relevant patterns in collinear data. Each comes with its own set of benefits and limitations, yet together they allow ecologists to directly address the nature of shared variance contributions in ecological data.

ACKNOWLEDGMENTS

M. Edwards, S. Thrush, P. Wainwright, G. Leonard, L. Ferry-Graham, D. Strong, A. Ellison, and two anonymous reviewers provided useful comments on the manuscript and participated in various discussions of multicollinearity. Special thanks to B. Hughes for assistance with the literature searches.

LITERATURE CITED

- Bentler, P. M. 1995. EQS: structural equations program manual. Multivariate Software, Encino, California, USA.
- Birkes, D., and Y. Dodge. 1993. Alternative methods of regression. John Wiley and Sons, New York, New York, USA.
- Bishop, J. G., and D. W. Schemske. 1998. Variation in flowering phenology and its consequences for lupines colonizing Mount St. Helens. *Ecology* **79**:534–546.
- Bollen, K. A. 1989. Structural equations with latent variables. John Wiley and Sons, New York, New York, USA.
- Brown, D. G., and A. E. Weis. 1995. Direct and indirect effects of prior grazing of goldenrod upon the performance of a leaf beetle. *Ecology* **76**:426–436.
- Carnes, B. A., and N. A. Slade. 1988. The use of regression for detecting competition with multicollinear data. *Ecology* **69**:1266–1274.
- Chatterjee, S., A. S. Hadi, and B. Price. 2000. Regression analysis by example. John Wiley and Sons, New York, New York, USA.
- Cohen, J., and P. Cohen. 1983. Applied multiple regression/correlation analysis for the behavioral sciences. Lawrence Erlbaum, Mahwah, New Jersey, USA.
- Draper, N. R., and H. Smith. 1998. Applied regression analysis. John Wiley and Sons, New York, New York, USA.
- Furnival, G. M. 1971. All possible regressions with less computations. *Technometrics* **13**:403–408.
- Gough, L., and J. B. Grace. 1999. Effects of environmental change on plant species density: comparing predictions with experiments. *Ecology* **80**:882–890.
- Graham, M. H. 1997. Factors determining the upper limit of giant kelp, *Macrocystis pyrifera* Agardh, along the Monterey Peninsula, central California, USA. *Journal of Experimental Marine Biology and Ecology* **218**:127–149.
- Hayduk, L. A. 1987. Structural equation modeling with LISREL: essentials and advances. Johns Hopkins University Press, Baltimore, Maryland, USA.
- Hocking, R. R. 1996. Methods and applications of linear models: regression and the analysis of variance. John Wiley and Sons, New York, New York, USA.
- James, F. C., and C. E. McCulloch. 1990. Multivariate analysis in ecology and systematics: panacea or Pandora's box? *Annual Review of Ecology and Systematics* **21**:129–166.
- Legendre, P., and L. Legendre. 1998. Numerical ecology. Elsevier, Amsterdam, The Netherlands.
- Loehlin, J. C. 1987. Latent variable models. Lawrence Erlbaum, Mahwah, New Jersey, USA.
- Mitchell-Olds, T., and R. G. Shaw. 1987. Regression analysis of natural selection: statistical inference and biological interpretation. *Evolution* **41**:1149–1161.
- Neter, J., M. H. Kutner, C. J. Nachtsheim, and W. Wasserman. 1996. Applied linear statistical models. Irwin, Chicago, Illinois, USA.
- Petraitis, P. S., A. E. Dunham, and P. H. Niewiarowski. 1996. Inferring multiple causality: the limitations of path analysis. *Functional Ecology* **10**:421–431.
- Philippi, T. E. 1993. Multiple regression: herbivory. Pages 183–210 in S. M. Scheiner and J. Gurevitch, editors. Design and analysis of ecological experiments. Chapman and Hall, New York, New York, USA.
- Scheffe, H. 1959. The analysis of variance. John Wiley and Sons, New York, New York, USA.
- Shipley, B. 1999. Testing causal explanations in organismal biology: causation, correlation and structural equation modeling. *Oikos* **86**:374–382.
- Tabachnick, B. G., and L. S. Fidell. 1996. Using multivariate statistics. HarperCollins, New York, New York, USA.
- Ullman, J. B. 1996. Structural equation modeling. Pages 709–812 in B. G. Tabachnick and L. S. Fidell, editors. Using multivariate statistics. HarperCollins College Publishers, New York, New York, USA.

APPENDIX A

Tables of SYSTAT output for backwards stepwise multiple regression for the model are available in ESA's Electronic Data Archive: *Ecological Archives* E084-073-A1.

APPENDIX B

Residual transformation equations and SYSTAT output for backwards stepwise sequential regression are available in ESA's Electronic Data Archive: *Ecological Archives* E084-073-A2.

APPENDIX C

Results of principal components analysis and SYSTAT output for backwards stepwise principal components regression are available in ESA's Electronic Data Archive: *Ecological Archives* E084-073-A3.

APPENDIX D

EQS protocol and output for structural equation modeling of the original explanatory variables are available in ESA's Electronic Data Archive: *Ecological Archives* E084-073-A4.

SUPPLEMENT

Data used in standard (Appendix A), sequential (Appendix B), and principal components regressions (Appendix C), and structural equation models (Appendix D) of the effect of various environmental factors on the distribution of giant kelp are available in ESA's Electronic Data Archive: *Ecological Archives* E084-073-S1.