

# A Theory of Gradient Analysis

CAJO J.F. TER BRAAK AND I. COLIN PRENTICE

I.	Introduction	236
II.	Linear Methods	241
	A. Regression	241
	B. Calibration	242
	C. Ordination	243
	D. The Environmental Interpretation of Ordination Axes (Indirect Gradient Analysis)	244
	E. Constrained Ordination (Multivariate Direct Gradient Analysis)	245
III.	Non-Linear (Gaussian) Methods	246
	A. Unimodal Response Models	246
	B. Regression	249
	C. Calibration	250
	D. Ordination	250
	E. Constrained Ordination	251
IV.	Weighted Averaging Methods	251
	A. Regression	252
	B. Calibration	252
	C. Ordination	254
	D. Constrained Ordination	258
V.	Ordination Diagrams and Their Interpretation	259
	A. Principal Components: Biplots	259
	B. Correspondence Analysis: Joint Plots	261
	C. Redundancy Analysis	263
	D. Canonical Correspondence Analysis	265
VI.	Choosing the Methods	267
	A. Which Response Model?	267
	B. Direct or Indirect?	268
	C. Direct Gradient Analysis: Regression or Constrained Ordination?	269
VII.	Conclusions	271
	Acknowledgements	273
	References	273
	Appendix	278

## I. INTRODUCTION

All species occur in a characteristic, limited range of habitats; and within their range, they tend to be most abundant around their particular environmental optimum. The composition of biotic communities thus changes along environmental gradients. Successive species replacements occur as a function of variation in the environment, or (analogously) with successional time (Pickett, 1980; Peet and Loucks, 1977). The concept of niche space partitioning also implies the separation of species along “resource gradients” (Tilman, 1982). Gradients do not necessarily have physical reality as continua in either space or time, but are a useful abstraction for explaining the distributions of organisms in space and time (Austin, 1985). Austin’s review explores the interrelationships between niche theory and the concepts of ecological continua and gradients.

Our review concerns data analysis techniques that assist the interpretation of community composition in terms of species’ responses to environmental gradients in the broadest sense. Gradient analysis *sensu lato* includes direct gradient analysis, in which each species’ abundance (or probability of occurrence) is described as a function of measured environmental variables; the converse of direct gradient analysis, whereby environmental values are inferred from the species composition of the community; and indirect gradient analysis, *sensu Whittaker (1967)*, in which community samples are displayed along axes of variation in composition that can subsequently be interpreted in terms of environmental gradients. There are close relationships among these three types of analysis. Direct gradient analysis is a *regression* problem—fitting curves or surfaces to the relation between each species’ abundance or probability of occurrence (the response variable) and one or more environmental variables (the predictor variable(s)) (Austin, 1971). Inferring environmental values from species composition when these relationships are known is a *calibration* problem. Indirect gradient analysis is an *ordination* problem, in which axes of variation are derived from the total community data. Ordination axes can be considered as latent variables, or hypothetical environmental variables, constructed in such a way as to optimize the fit of the species data to a particular (linear or unimodal) statistical model of how species abundance varies along gradients (Ter Braak, 1985, 1987a). These latent variables are constructed without reference to environmental measurements, but they can subsequently be compared with actual environmental data if available. To these three well-known types of gradient analysis we add a fourth, *constrained ordination*, which has its roots in the psychometric literature on multidimensional scaling (Bloxom, 1978; De Leeuw and Heiser, 1980; Heiser, 1981). Constrained ordination also constructs axes of variation

in overall community composition, but does so in such a way as to explicitly optimize the fit to supplied environmental data (Ter Braak, 1986; Jongman *et al.*, 1987). Constrained ordination is thus a multivariate generalization of direct gradient analysis, combining aspects of regression, calibration and ordination. Table 1 gives an arbitrary selection of literature references, chosen simply to illustrate the wide range of ecological problems to which each of the four types of gradient analysis has been applied; the reader is also referred to Gauch (1982), who includes an extensive bibliography, and to Gittins (1985).

Standard statistical methods that assume linear relationships among variables exist for all four types of problems (regression, calibration, ordination and constrained ordination), but have found only limited application in ecology because of the generally non-linear, non-monotone response of species to environmental variables. Ecologists have independently developed a variety of alternative techniques. Many of these techniques are essentially heuristic, and have a less secure theoretical basis. These heuristic techniques can nevertheless give useful results, and can be understood as approximate solutions to statistical problems similar to those solved by standard methods, but formulated in terms of a *unimodal* (Gaussian or similar) response model instead of a linear one. We present here a theory of gradient analysis, in which the heuristic techniques are integrated with regression, calibration, ordination and constrained ordination as distinct, well-defined statistical problems.

The various techniques used for each type of problem are classified into families according to their implicit response model and the method used to estimate parameters of the model. We consider three such families (Table 2). First we treat the family of standard statistical techniques based on the linear response model, because these are conceptually the simplest and provide a basis for what follows, even though their ecological application is restricted. Second, we outline a family of somewhat more complex statistical techniques which are formal extensions of the standard linear techniques and incorporate unimodal (Gaussian-like) response models explicitly. Finally, we consider the family of heuristic techniques based on weighted averaging. These are not more complex than the standard linear techniques, but implicitly fit a simple unimodal response model rather than a linear one. Our treatment thus unites such apparently disparate data analysis techniques as linear regression, principal components analysis, redundancy analysis, Gaussian ordination, weighted averaging, reciprocal averaging, detrended correspondence analysis and canonical correspondence analysis in a single theoretical framework.

**Table 1** Selected applications of gradient analysis

Type of problem	Taxa	Environmental variables	Purpose of study
<i>Regression</i>			
Alderdice (1972)	Marine fish	Salinity, temperature	Defining ranges
Peet (1978)	Trees	Elevation, moisture, latitude	Biogeography
Wiens and Rotenberry (1981)	Birds	Vegetation structure	Niche characterization
Austin <i>et al.</i> (1984)	<i>Eucalyptus</i> spp.	Climatic indices	Habitat characterization
Bartlein <i>et al.</i> (1986)	Plant pollen types	Temperature, precipitation	Quaternary palaeoecology
<i>Calibration</i>			
Chandler (1970)	Benthic macro-invertebrates	Water pollution	Water quality management
Imbrie and Kipp (1971)	Foraminifera	Sea surface temperature	Palaeoclimatic reconstruction
Sládeček (1973)	Freshwater algae	Organic pollution	Ecological monitoring
Balloch <i>et al.</i> (1976)	Benthic macro-invertebrates	Water pollution	Ecological monitoring
Ellenberg (1979)	Terrestrial plants	Soil moisture, N, pH	Bioassay from vegetation
van Dam <i>et al.</i> (1981)	Diatoms	pH	Acid rain effects
Böcker <i>et al.</i> (1983)	Terrestrial plants	Soil moisture, N, pH	Bioassay from vegetation
Bartlein <i>et al.</i> (1984)	Plant pollen types	Temperature, precipitation	Palaeoclimatic reconstruction
Battarbee (1984)	Diatoms	pH	Acid rain effects
Charles (1985)	Diatoms	pH	Acid rain effects
Atkinson <i>et al.</i> (1986)	Beetles	Summer temperature, annual range	Palaeoclimatic reconstruction

*Ordination<sup>a</sup>*

van der Aart and Smeenk-Enserink (1975)	Spiders	Microenvironmental features	Habitat characterization
Kooijman and Hengeveld (1979)	Beetles	Lutum content, elevation	Habitat characterization
Wiens and Rotenberry (1981)	Birds	Vegetation structure	Niche characterization
Prodon and Lebreton (1981)	Birds	Vegetation structure	Niche characterization
Kalkhoven and Opdam (1984)	Birds	Habitat and ladscape features	Habitat characterization
Macdonald and Ritchie (1986)	Plant pollen types	Vegetation regions	Quaternary palaeoecology
<i>Constrained ordination</i>			
Webb and Bryson (1972)	Plant pollen types	Climate variables, airmass frequencies	Palaeoclimatic reconstruction
Gasse and Tekaia (1983)	Diatoms	pH classes	Palaeolimnology
Ås (1985)	Beetles	Vegetation types	Niche theory
Cramer and Hytteborn (1987)	Terrestrial plants	Time, elevation	Land uplift effects
Purata (1986)	Tropical trees	Successional boundary conditions	Study of secondary succession
Fångström and Willén (1987)	Phytoplankton	Physical/chemical variables	Environmental monitoring

<sup>a</sup> Excluding vegetation studies, where ordination is used routinely: see [Gauch \(1982\)](#) for a review.

**Table 2** Classification of gradient analysis techniques by type of problem, response model and method of estimation

Type of problem	Linear response model	Unimodal response model	
	Least-squares, estimation	Maximum likelihood estimation	Weighted averaging estimation
Regression	Multiple regression	Gaussian regression	Weighted averaging of site scores (WA)
Calibration	Linear calibration; "inverse regression"	Gaussian calibration	Weighted averaging of species' scores (WA)
Ordination	Principal components analysis (PCA)	Gaussian ordination	Correspondence analysis (CA); detrended correspondence analysis (DCA)
Constrained ordination <sup>a</sup>	Redundancy analysis (RDA) <sup>d</sup>	Gaussian canonical ordination	Canonical correspondence analysis (CCA); detrended CCA
Partial ordination <sup>b</sup>	Partial components analysis	Partial Gaussian ordination	Partial correspondence analysis; partial DCA
Partial constrained ordination <sup>c</sup>	Partial redundancy analysis	Partial Gaussian canonical ordination	Partial canonical correspondence analysis; partial detrended CCA

<sup>a</sup> Constrained multivariate regression.

<sup>b</sup> Ordination after regression on covariables.

<sup>c</sup> Constrained ordination after regression on covariables = constrained partial multivariate regression.

<sup>d</sup> "Reduced-rank regression" = "PCA of  $y$  with respect to  $x$ ".

## II. LINEAR METHODS

Species abundances may seem to change linearly through *short* sections of environmental gradients, so a linear response model may be a reasonable basis for analysing quantitative abundance data spanning a narrow range of environmental variation.

### A. Regression

If a plot of the abundance ( $y$ ) of a species against an environmental variable ( $x$ ) looks linear, or can easily be transformed to linearity, then it is appropriate to fit a straight line by linear regression. The formula  $y = a + bx$  describes the linear relation, with  $a$  the intercept of the line on the  $y$ -axis and  $b$  the slope of the line, or regression coefficient. Separate regressions can be carried out for each of  $m$  species.

We are usually most interested in how the abundance of each species changes with a change in the environmental variable, i.e. in the slopes  $b_k$  (the index  $k$  refers to species  $k$ ). If we first centre the data—by subtracting the mean of each species' abundances from the species data and the mean of the environmental values from the environmental data—the intercept disappears. Then if  $y_{ki}$  denotes the centred abundance of species  $k$  in the  $i$ th out of  $n$  sites, and  $x_i$  the centred environmental value for that site, the response model for fitting the straight lines becomes

$$y_{ki} = b_k x_i + e_{ki} \quad (1)$$

where  $e_{ki}$  is an error component with zero mean and variance  $v_{ki}$ . The standard estimator for the slope in Eq. (1) is

$$\tilde{b}_k = \sum_{i=1}^n y_{ki} x_i / s_x^2 \quad (2)$$

where  $s_x^2 = \sum_{i=1}^n x_i^2$ . This is the least-squares estimator, which is the best linear unbiased estimator when errors are uncorrelated and homogeneous across sites ( $v_{ki} = v_k$ ). It is also the maximum likelihood (ML) estimator when the errors are normally distributed. The fitted lines can be used to predict the abundances of species in a site with a known value of the environmental variable simply by reading off the graph.

Species experience the effect of more than one environmental variable simultaneously, so more than one variable may be required to account for variation in species abundances. The joint effect of two or more

environmental variables on a species can be analysed by multiple regression (see e.g. [Montgomery and Peck, 1982](#)). Standard computer packages are available to obtain least-squares (ML) estimates for the regression coefficients. Only when the environmental variables are uncorrelated will the partial regression coefficients be identical to the coefficients estimated by separate regressions using [Eq. \(1\)](#).

## B. Calibration

We now turn to the inverse problem, calibration. When the relationship between the abundances of species and the environmental variable we are interested in is known, we can infer values of that environmental variable for new sites from the observed species abundances. If we took into account the abundance of only a single species, we could simply read off the graph, starting from a value on the vertical axis. However, another species may well give a different estimate. We therefore need a good and unambiguous estimator that combines the information from all  $m$  species. In terms of [Eq. \(1\)](#), the  $b_k$  are now assumed to be known and  $x_i$  is unknown. The role of the  $b_k$  and  $x_i$  have been interchanged. By interchanging their roles in [Eq. \(2\)](#) as well, we obtain

$$\tilde{x}_i = \sum_{k=1}^m y_{ki} b_k / s_b^2 \quad (3)$$

where  $s_b^2 = \sum_{k=1}^m b_k^2$ . This is the least-squares estimator (and the ML-estimator) when the errors follow a normal distribution and are independent and homogeneous across species ( $v_{ki} = v_i$ ).

A problem with [Eq. \(3\)](#) is that these conditions are likely to be unrealistic, because effects of other environmental variables can cause correlation between the abundances of different species even after the effects of the environmental variable of interest have been removed. Further, the residual variance  $v_{ki}$  may be different for different species. If this occurs, we also need to take the residual correlations and variances into account. In practice, the residual correlations and variances are estimated from the residuals of the regressions used for estimating the  $b_k$ 's. Searching for the maximum of the likelihood with respect to  $x_i$  then leads to a general weighted least-squares problem ([Brown, 1979](#); [Brown, 1982](#)) that can be solved by using standard algorithms.

Inferring values of more than one environmental variable simultaneously has been given surprisingly little attention in the literature. However, [Williams \(1959\)](#) and [Brown \(1982\)](#) derived the necessary formulae from the ML-principle ([Cox and Hinkley, 1974](#)).



### C. Ordination

After having fitted a particular environmental variable to the species data by regression, we might ask whether another environmental variable would provide a better fit. For some species one variable may fit better, and for other species another variable. To get an overall impression we might judge the goodness-of-fit (explanatory power) of an environmental variable by the total regression sum of squares (Jongman *et al.*, 1987). The question then arises: what is the best possible fit that is theoretically obtainable with the straight line model of Eq. (1)?

This question defines an ordination problem, i.e. to construct the single “hypothetical environmental variable” that gives the best fit to the species data according to Eq. (1). This hypothetical environmental variable is termed the *latent variable*, or simply the (first) ordination axis. Principal components analysis (PCA) provides the solution to this ordination problem. In Eq. (1),  $x_i$  is then the score of site  $i$  on the latent variable,  $b_k$  is the slope for species  $k$  with respect to the latent variable (also called the species loading or species score) and the eigenvalue of the first PCA axis is equal to the goodness-of-fit, i.e. the total sum of squares of the regressions of the species abundances on the latent variable. PCA provides the least-squares estimates of the site and species scores: these estimates are also ML estimates if the errors are independently and normally distributed with constant variance ( $v_{ki} = v$ ).

PCA is usually performed using a standard computer package, but several different algorithms can be used to do the same job. The following algorithm, known as the power method (Gourlay and Watson, 1973), makes the relationship between PCA and regression and calibration clear in a way that the usual textbook treatment, in terms of singular value decomposition of inner product matrices, does not; it also facilitates comparison with correspondence analysis, which we discuss later. The power method shows that PCA can be obtained by an alternating sequence of linear regressions and calibrations:

- Step 1* Start with some (arbitrary) initial site scores  $\{x_i\}$  with zero mean.
- Step 2* Calculate new species scores  $\{b_k\}$  by linear regression (Eq. (2)).
- Step 3* Calculate new site scores  $\{x_i\}$  by linear calibration (Eq. (3)).
- Step 4* Remove the arbitrariness in scale by standardizing the site scores as follows: new  $x_i = \text{old } x_i / n / s_x$ , with  $s_x$  as defined beneath Eq. (2).
- Step 5* Stop on convergence, i.e. when the newly obtained site scores are close to the site scores of the previous cycle of iteration, else go to Step 2.

The final scores do not depend on the initial scores.

The ordination problem for a two-dimensional linear model turns out to be relatively simple, compared with the regression and calibration problems. The solution does not need an alternating sequence of *multiple* regressions and calibrations, because the latent variables can always be chosen in such a way that they are uncorrelated; and if the latent variables are uncorrelated, then the multiple regressions and calibrations reduce to a series of separate linear regressions and calibrations. PCA provides the solution to the linear ordination problem in any number of dimensions; one latent variable is derived first, as in the one-dimensional case of Eq. (1), and the second latent variable can be obtained next by applying the same algorithm again but with one extra step—after Step 3, the trial scores are made uncorrelated with the first latent variable. On denoting the scores of the first axis by  $x_{i1}$ , this orthogonalization is computed by

*Step 3b* Calculate  $f = \sum_i x_i x_{i1} / n$ ,  
Calculate new  $x_i = \text{old } x_i - f x_{i1}$ .

Further latent variables (ordination axes) may be derived analogously. As in the one-dimensional case, PCA provides the ML-solution to the multi-dimensional linear ordination problem if the errors are independently and normally distributed with constant variance across species and sites. Jolliffe (1986) reviews the theory and applications of PCA.

#### D. The Environmental Interpretation of Ordination Axes (Indirect Gradient Analysis)

In indirect gradient analysis the species data are first subjected to ordination, e.g. using PCA, to find a few major axes of variation (latent variables) with a good fit to the species data. These axes are then interpreted in terms of known variation in the environment, often by using graphical methods (Gauch, 1982). A more formal method for the latter step would be to calculate correlation coefficients between environmental variables and each of the ordination axes. This analysis is similar to performing a multiple regression of each separate environmental variable on the axes (Dargie, 1984), because the axes are uncorrelated. A joint analysis of all environmental variables can be carried out by multiple regression of each ordination axis on the environmental variables:

$$x_i = c_0 + \sum_{j=1}^q c_j z_{ji} \quad (4)$$

in which  $x_i$  is the score of site  $i$  on that one ordination axis,  $z_{ij}$  denotes the value at site  $i$  of the  $j$ th out of  $q$  actual environmental variables, and  $c_j$  is the corresponding regression coefficient. For later reference, the error term in Eq. (4) is not shown. The multiple correlation coefficient  $R$  measures how well the environmental variables explain the ordination axis.

### E. Constrained Ordination (Multivariate Direct Gradient Analysis)

Indirect gradient analysis, as outlined above, is a *two-step* approach to relate species data to environmental variables. A few ordination axes that summarize the overall community variation are extracted in the first step; then in the second step one may calculate weighted sums of the environmental variables that most closely fit each of these ordination axes. However, the environmental variables that have been studied may turn out to be poorly related to the first few ordination axes, yet may be strongly related to other, “residual” directions of variation in species composition. Unless the first few ordination axes explain a very high proportion of the variation, this residual variation can be substantial, and strong relationships between species and environment can potentially be missed.

In constrained ordination this approach is made more powerful by combining the two steps into one. The idea of constrained ordination is to search for a few weighted sums of environmental variables that fit the data of all species best, i.e. that give the maximum total regression sum of squares. The resulting technique, redundancy analysis (Rao, 1964; van den Wollenberg, 1977), is an ordination of the species data in which the axes are constrained to be linear combinations of the environmental variables. These axes can be found by extending the algorithm of PCA described above with one extra step, to be performed directly after Step 3 (Jongman *et al.*, 1987):

*Step 3a* Calculate a multiple regression of the site scores  $\{x_i\}$  on the environmental variables (Eq. (4)), and take as new site scores the fitted values of this regression.

The regression is thus carried out within the iteration algorithm, instead of afterwards. On convergence, the coefficients  $\{c_j\}$  are termed canonical coefficients and the multiple correlation coefficient in Step 3a can be called the species-environment correlation.

Redundancy analysis is also known as reduced-rank regression (Davies and Tso, 1982), PCA of  $y$  with respect to  $x$  (Robert and Escoufier, 1976) and two-block mode C partial least-squares (Wold, 1982). It is intermediate between PCA and separate multiple regressions for each of the species: it is a constrained ordination, but it is also a constrained form of (multivariate)

multiple regression (Davies and Tso, 1982; Israëls, 1984). By inserting Eq. (4) into Eq. (1), it can be shown that the “regression” coefficient of species  $k$  with respect to environmental variable  $j$  takes the simple form  $b_k c_j$ . With two ordination axes this form would be, in obvious notation,  $b_{k1} c_{j1} + b_{k2} c_{j2}$ . With two ordination axes, redundancy analysis thus uses  $2(q+m)+m$  parameters to describe the species data, whereas the multiple regressions use  $m(q+1)$  parameters. One of the attractive features of redundancy analysis is that it leads to an ordination diagram that simultaneously displays (i) the main pattern of community variation as far as this variation can be explained by the environmental variables, and (ii) the main pattern in the correlation coefficients between the species and each of the environmental variables. We give an example of such a diagram later on.

Redundancy analysis is much less well known than canonical correlation analysis (Gittins, 1985; Tso, 1981), which is the standard linear multivariate technique for relating two sets of variables (in our case, the set of species and the set of environmental variables). Canonical correlation analysis is very similar to redundancy analysis, but differs from it in the assumptions about the error component: uncorrelated errors with equal variance in redundancy analysis and correlated normal errors in canonical correlation analysis (Tso, 1981; Jongman *et al.*, 1987). The most important practical difference is that redundancy analysis can analyse any number of species whereas in canonical correlation analysis the number of species ( $m$ ) must be less than  $n-q$  (Griffins, 1985: 24); this restriction is often a nuisance.

Canonical variates analysis, or multiple discriminant analysis, is simply the special case of canonical correlation analysis in which the “environmental” variables are a series of dummy variables reflecting a single-factor classification of the samples. A similar restriction on the number of species thus also applies to canonical variates analysis. Redundancy analysis with dummy variables provides an alternative to canonical variates analysis, evading this restriction.

### III. NON-LINEAR (GAUSSIAN) METHODS

#### A. Unimodal Response Models

Linear methods are appropriate to community analysis only when the species data are quantitative abundances (with few zeroes) and the range of environmental variation in the sample set is narrow. Alternative analytical methods can be derived from unimodal models.

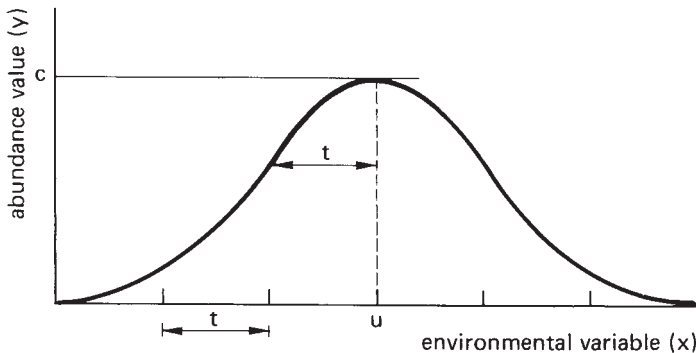
A unimodal response model for one environmental variable can be obtained by adding a quadratic term ( $x_i^2$ ) to the linear model, changing the response curve from a straight line into a parabola. But this quadratic model can predict large negative values, whereas species abundances are always zero or positive. A simple remedy for the problem of negative values is provided by the Gaussian response curve (Gauch and Whittaker, 1972) in which the *logarithm* of species abundance is a quadratic in the environmental variable:

$$\begin{aligned}\log y &= b_0 + b_1x + b_2x^2 \\ &= a - \frac{1}{2}(x - u)^2/t^2\end{aligned}\quad (5a)$$

where  $b_2 < 0$  (otherwise the curve would have a minimum instead of a mode). The coefficients  $b_0$ ,  $b_1$ , and  $b_2$  are most easily interpreted by transformation to  $u$ ,  $t$  and  $a$  (Figure 1),  $u$  being the species' optimum (the value of  $x$  at the peak),  $t$  being its tolerance (a measure of response breadth or ecological amplitude), and  $a$  being a coefficient related to the height of the peak (Ter Braak and Looman, 1986).

A closely related model can describe species data in presence-absence form. In analysing presence-absence data, we want to relate probability of occurrence ( $p$ ) to environment. Probabilities are never greater than 1, so rather than using Eq. (5a) we use the Gaussian logit model,

$$\log\left(\frac{p}{1-p}\right) = b_0 + b_1x + b_2x^2 \quad (5b)$$



**Figure 1** A Gaussian curve displays a unimodal relation between the abundance value ( $y$ ) of a species and an environmental variable ( $x$ ). ( $u$  = optimum or mode;  $t$  = tolerance;  $c$  = maximum =  $\exp(a)$ ).

which is very similar to the Gaussian model unless the peak probability is high ( $> 0.5$ ); then Eq. (5b) gives a curve that is somewhat flatter on top. The coefficients  $b_0$ ,  $b_1$ , and  $b_2$  can be transformed as before into coefficients representing the species' optimum, tolerance and maximum probability value.

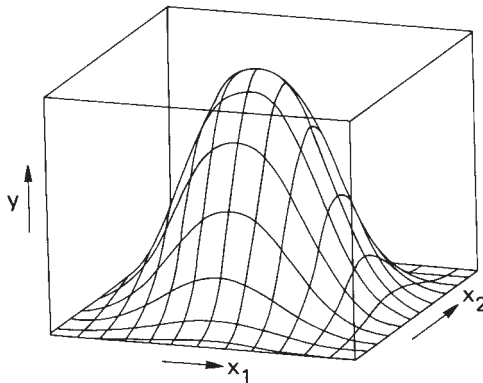
Although real ecological response curves are still more complex than implied by the Gaussian and Gaussian logit models, these models are nevertheless useful in developing statistical descriptive techniques for data showing mostly unimodal responses, just as linear models are useful in statistical analysis of data that are only approximately linear.

With two environmental variables, Eqs. (5a) and (5b) become full quadratic forms with both square and product terms (Alderdice, 1972). For example, the Gaussian model becomes

$$\log y = b_0 + b_1x_1 + b_2x_1^2 + b_3x_2 + b_4x_2^2 + b_5x_1x_2 \quad (6)$$

If  $b_2 + b_4 < 0$ , and  $4b_2b_4 - b_5^2 > 0$  then Eq. (6) describes a unimodal surface with ellipsoidal contours (Figure 2). If one of these conditions is not satisfied then Eq. (6) describes a surface with a minimum, or with a saddle point (e.g. Davison, 1983). Provided the surface is unimodal, its optimum ( $u_1, u_2$ ) can be calculated from the coefficients in Eq. (6) by

$$\left. \begin{aligned} u_1 &= (b_5b_3 - 2b_1b_4)/d \\ u_2 &= (b_5b_1 - 2b_3b_2)/d \end{aligned} \right\} \quad (7)$$



**Figure 2** A Gaussian surface displays a unimodal relation between the abundance value ( $y$ ) of a species and two environmental variables ( $x_1$  and  $x_2$ ).

where  $d = 4b_2b_4 - b_5^2$ . When  $b_5 \neq 0$ , the optimum with respect to  $x_1$ , depends on the value of  $x_2$ ; the environmental variables are then said to show interaction in their effect on the species. In contrast, when  $b_5 = 0$  the optimum with respect to  $x_1$ , does not depend on the value of  $x_2$  (no interaction) and Eq. (7) simplifies considerably (Ter Braak and Looman, 1986).

The unknown parameters of non-linear response models in the context of regression, calibration or ordination can (at least in theory) be estimated by the maximum likelihood principle, however difficult this may be in a particular situation. Usually iterative methods are required, and initial parameter values must be specified. The likelihood function may have local maxima, so that different sets of initial parameter values may result in different final estimates. It cannot be guaranteed that the global maximum has been found. Furthermore, all kinds of numerical problems may occur. However, the special cases of Gaussian and Gaussian logit response models do allow reasonably practical solutions, which we consider now.

## B. Regression

The regression problems of fitting Gaussian or Gaussian logit curves or surfaces are relatively straightforward, since these models can be fitted by Generalized Linear Modelling (GLM: McCullagh and Nelder, 1983; Dobson, 1983). An elementary introduction to GLM directed at ecologists is provided by Jongman *et al.* (1987). GLM is more flexible than ordinary multiple regression because one can specify “link functions” and error distributions other than the normal distribution. For example, the Gaussian models of Eqs. (5a) and (6) can be fitted with GLM to abundance data (which may include zeroes) by specifying the link function to be logarithmic and the error distribution to be Poissonian. The corresponding Gaussian logit models can be fitted with GLM to presence–absence data by specifying the link function to be logistic and the error distribution to be binomial-with-total-1. Alternatively, any statistical package that will do logit (= logistic) regression can be used to fit the Gaussian logic model. No initial estimates are needed and local maxima do not arise, so these techniques are quite practical for direct gradient analysis. For examples of the use of GLM in ecology see Austin and Cunningham (1981) and Austin *et al.* (1984).

The most common complications arise when the optimum for a species is estimated well outside the sampled range of environments, or if the fitted curve shows a minimum rather than a peak. These conditions suggest that the regression is ill-determined and that it might be better to fit a monotone curve by setting  $b_2 = 0$  in Eq. (5); a statistical test can be used to determine

whether this simplification is acceptable (Jongman *et al.*, 1987). Such cases are bound to arise in practice because any given set of samples will include some species that are near the edge of their range.

### C. Calibration

The calibration problem of inferring environmental values at sites from species data and known Gaussian (logit) curves by ML is feasible by numerical optimization, but no computer programs are available at present that are easy to use (Jongman *et al.*, 1987). Local maxima may occur in the likelihood, when the tolerances of the species are unequal, and one needs to specify an initial estimate. The assumption of independence of species responses is required, but might not be tenable in practice; it remains to be studied how important this assumption is. Dependency among species could most obviously be caused by the effects of additional, unconsidered environmental variables, in which case the best remedy would be to identify these variables and include them in the analysis. Inferring the values of more than one environmental variable simultaneously on the basis of several Gaussian (logit) response surfaces is also possible in principle, but has not been done as far as we know.

### D. Ordination

Ordination based on Gaussian (logit) curves aims to construct a latent variable such that these curves optimally fit the species data. This problem involves the ML estimation of site scores  $\{x_k\}$  and the species' optima  $\{u_k\}$ , tolerances  $\{t_k\}$  and maxima  $\{a_k\}$ , usually by an alternating sequence of Gaussian (logit) regressions and calibrations. This kind of ordination has been investigated by Gauch *et al.* (1974), Kooijman (1977), Kooijman and Hengeveld (1979), Goodall and Johnson (1982) and Ihm and Van Groenewoud (1975, 1984). The numerical methods required are computationally demanding, and in the general case, when the tolerances of the species are allowed to differ, the likelihood function typically contains many local maxima.

Kooijman (1977) and Goodall and Johnson (1982) reported numerical problems in their attempts to perform ML ordination using two-dimensional Gaussian-like models. A simple model with circular contours ( $b_2 = b_4$  and  $b_5 = 0$ ) may be amenable in practice, especially if  $b_2$  is not allowed to vary among species (Kooijman, 1977). This model is equivalent to the "unfolding model" used by psychologists to analyse preference data (Coombs, 1964; Heiser, 1981; Davison, 1983; DeSarbo and Rao, 1984).



But with more than two latent variables the Gaussian (logit) model with a second-degree polynomial as linear predictor contains so many parameters that it is likely to be difficult to get reliable estimates of them, even if all the interaction terms are dropped.

### E. Constrained Ordination

The constrained ordination problem for Gaussian-like response models is to construct ordination axes that are also linear combinations of the environmental variables, such that Gaussian (logit) surfaces with respect to these axes optimally fit the data. As in redundancy analysis (Section II.E), the joint effects of the environmental variables on the species are “channelled” through a few ordination axes which can be considered as composite environmental gradients influencing species composition. Ter Braak (1986) refers to this approach as Gaussian canonical ordination, the word canonical being chosen by analogy with canonical correlation analysis. The estimation problem is actually simpler than in unconstrained Gaussian ordination, and is more easily soluble in practice because the number of parameters to be estimated is smaller: instead of  $n$  site scores one has to estimate  $q$  canonical coefficients. Meulman and Heiser (1984) have applied similar ideas in the context of non-metric multidimensional scaling. Gaussian canonical ordination can also be viewed as multivariate Gaussian regression with constraints on the coefficients of the polynomial (Ter Braak, 1988). In multivariate Gaussian regression each species has its own optimum in the  $q$ -dimensional space formed by the environmental variables; the constraints imposed in Gaussian canonical ordination amount to a requirement that these optima lie in a low-dimensional subspace. If the optima lie close to a plane then the most important species–environment relationships can be depicted graphically in an ordination diagram.

## IV. WEIGHTED AVERAGING METHODS

Ecologists have developed alternative, heuristic methods that are simpler but have essentially the same aims as the methods of the previous section based on Gaussian-type models. Each method in the Gaussian family has a counterpart in the family of heuristic methods based on weighted averaging (WA). These methods have been used extensively, and even re-invented in different branches of ecology.

### A. Regression

WA can be used to estimate species' optima with respect to known environmental variables. When a species shows a unimodal relationship with environmental variables, the species' presences will be concentrated around the peak of this function. One intuitively reasonable estimate of the optimum is the average of the values of the environmental variable over those sites in which the species is present. With abundance data, WA applies weights proportional to species abundance; absences still carry zero weight. The estimate of the optimum for species  $k$  is thus

$$\tilde{u}_k = \sum_{i=1}^n y_{ki}x_i/y_{k+} \quad (8)$$

where  $y_{ki}$  is from now onwards the abundance (*not* centred) or presence/absence (1/0) of species  $k$  at site  $i$ ,  $y_{k+}$  is the species total ( $y_{k+} = \sum_i y_{ki}$ ) and  $x_i$  is the value of the environmental variable at site  $i$ . As a follow-up to an investigation of the theoretical properties of this estimator (Ter Braak and Barendregt, 1986), Ter Braak and Looman (1986) showed by simulation of presence-absence data that WA estimates the optimum of a Gaussian logit curve as efficiently as the ML technique of Gaussian logic regression provided:

*Condition 1a* The site scores  $\{x_i\}$  are equally spaced over the whole range of occurrence of the species along the environmental variable.

WA also proved to be only a little less efficient whenever the distribution of the environmental variable among the sites was reasonably homogeneous (rather than strictly equally spaced) over the whole range of species occurrences, or more generally for species with narrow ecological amplitudes. But the estimate of the optimum of a rare species may be imprecise, because the standard error of the estimate is inversely proportional to the square root of the number of occurrences. So for efficiency, we also need

*Condition 1b* The site scores  $\{x_i\}$  are closely spaced in comparison with the species' tolerance.

### B. Calibration

WA is also used in calibration, to estimate environmental values at sites from species' optima—which in this context are often called indicator values

(“Zeigerwerte”, [Ellenberg, 1979](#)) or scores ([Whittaker, 1956](#)). When species replace one another along the environmental variable of interest, i.e. have unimodal response functions with optima spread out along that variable, then species with optima close to the environmental value of a site will naturally tend to be represented at that site. Intuitively, to estimate the environmental value at a site, one can average the optima of the species that are present. With abundance data, the corresponding intuitive estimate is the weighted average,

$$\tilde{x}_i = \sum_{k=1}^m y_{ki} u_k / y_{+i} \quad (9)$$

where  $y_{+i}$  is the site total ( $y_{+i} = \sum_k y_{ki}$ ).

[Ter Braak and Barendregt \(1986\)](#) showed that WA estimates the value  $x_i$  of a site as well as the corresponding ML techniques if the species show Gaussian curves and Poisson-distributed abundance values (or, for presence-absence data, show Gaussian logit curves), and provided:

- Condition 2a* The species' optima are equally spaced along the environmental variable over an interval that extends for a sufficient distance in both directions from the true value  $x_i$ ;  
*Condition 3* The species have equal tolerances;  
*Condition 4* The species have equal maximum values.

These conditions amount to a “species packing model” wherein the species have equal response breadth and equal spacing ([Whittaker et al., 1973](#)). The conditions may be relaxed somewhat ([Ter Braak and Barendregt, 1986](#)) without seriously affecting the efficiency of the WA-estimate. When the optima are uniformly distributed instead of being equally spaced, the efficiency is still high if the maximum probabilities of occurrence are small ( $< 0.5$ ). The species' maximum values may differ, but they must not show a trend along the environmental variable (for instance, leading to species-rich samples at one end of the gradient and species-poor samples at the other end). The efficiency of WA is less good if the tolerances substantially differ among species; a tolerance weighted version of WA, as suggested by [Zelinka and Marvan \(1961\)](#) and [Goff and Cottam \(1967\)](#), would be more efficient since it would give greater weight to species of narrower tolerance, which are more informative about the environment.

Under Conditions 2a–4 above, the standard error of the estimate of  $\tilde{x}_i$  is approximately  $t/\sqrt{y_{+i}}$ , where  $t$  is the (common) species tolerance. For the weighted average to be practically useful, the number of species encountered

in a site should therefore not be too small (not less than five). We therefore need the extra condition (cf. Section 5 in [Ter Braak and Barendregt, 1986](#)):

*Condition 2b* The species' optima must be closely spaced in comparison with their tolerances.

An alternative heuristic method of calibration is by "inverse regression". This is simply multiple linear regression of the environmental variable on the species abundances ([Brown, 1982](#)): the environmental variable is treated as if it were the response variable and the species abundances, possibly transformed, as predictor variables. The regression coefficients can be estimated from the training set of species abundances and environmental data, the resulting equations being applied directly to infer environmental values from further species abundance data. When applied to data on percentage composition, e.g. pollen spectra or diatom assemblages ([Bartlein et al., 1984](#); [Charles, 1985](#)), the method differs from WA calibration only in the way in which the species optima are estimated, since the linear combination of percentage values used to estimate the environmental value is by definition a weighted average of the regression coefficients.

### C. Ordination

[Hill \(1973\)](#) turned weighted averaging into an ordination technique by applying alternating WA regressions and calibrations to a species-by-site data table. The algorithm of this technique of "reciprocal averaging" is similar to that given earlier for PCA:

- Step 1* Start with arbitrary, but unequal, initial site scores  $\{x_i\}$ .
- Step 2* Calculate new species scores  $\{u_k\}$  by WA (Eq. (8)).
- Step 3* Calculate new site scores  $\{x_i\}$  by WA (Eq. (9)).
- Step 4* Remove the arbitrariness in scale by standardizing the site scores by new  $x_i = \{\text{old } x_i - z\} / s$  where  $z = \sum_i y_{+i} x_i / \sum_i y_{+i}$  and

$$s^2 = \frac{\sum_i y_{+i} (x_i - z)^2}{\sum_i y_{+i}} \quad (10)$$

- Step 5* Stop on convergence, else go to Step 2.

As in PCA, the resulting site and species scores do not depend on the initial scores. The final scores produced by this reciprocal averaging algorithm form the first eigenvector or ordination axis of correspondence analysis (CA), an eigenvector technique that is widely used especially in the French-language literature ([Laurec et al., 1979](#); [Hill, 1974](#)). As with the power algorithm for PCA, the reciprocal averaging algorithm makes clear the

relationship between CA and regression and calibration—this time, with WA regression and calibration. The method of standardization chosen in Step 4 is arbitrary, but chosen for later reference. On convergence,  $s$  in Step 4 is equal to the eigenvalue of the first axis, and lies between 0 and 1.

Correspondence analysis has many applications outside ecology. Nishisato (1980), Greenacre (1984) and Gifi (1981) provide a variety of different rationales for correspondence analysis, each adapted to a particular type of application. Heiser (1987) and Ter Braak (1985, 1987a) develop rationales for correspondence analysis that are particularly relevant to ecological applications.

Ter Braak (1985) showed that CA approximates ML Gaussian (logit) ordination under Conditions 1–4 listed above, i.e. under just these conditions for which WA is as good as ML-regression and ML-calibration. In practice CA can never be exactly equivalent to ML ordination, because Condition 1a implies that the range of site scores is broad enough to include the ranges of all of the species, whereas Condition 2a implies that there must be species with their optima situated beyond the edge of the range of site scores. These conditions cannot both be satisfied if the range of site scores is finite. As a result, CA shows an edge effect: the site scores near the ends of the axes become compressed relative to those in the middle (Gauch, 1982). This effect becomes less strong, however, as the range of site scores becomes wider and the spacing of the site scores and species scores becomes closer relative to the average species' tolerance.

Conditions 1–4 also disallow “deviant” sites and rare species. CA is sensitive to both (Hill, 1974; Feoli and Feoli Chiapella, 1980; Oksanen, 1983). This sensitivity may be useful in some applications, but is a nuisance if the aim is to detect major gradients. Deviant sites (and, possibly, the rarest species) should therefore ideally be removed from the data before analysis by CA.

As in PCA, further ordination axes can be extracted in CA by adding an extra step after Step 3, making the trial scores on the second axis uncorrelated with the (final) scores on the first axis. (In the calculation of  $f$  in Step 3b (see Section II.C) the sites are weighted proportional to the site total  $y_{+i}$ . This weighting is implicitly applied from now on.) However, there is a problem with the second and higher axes in CA. The problem is the well-known but hitherto not well-understood “arch effect” (Hill, 1974). If the species data come from an underlying one-dimensional Gaussian model the scores on the second ordination axis show a parabolic (“arch”) relation with those of the first axis; if the species data come from a two-dimensional Gaussian model in which the true site and species scores are located homogeneously in a rectangular region in two-dimensional space (the extension to two dimensions of Conditions 1a and 2a), the scores of the second ordination axis lie not in a rectangle but in an arched band (Hill and

Gauch, 1980). The arch effect arises because the axes are extracted sequentially in order of decreasing “variance”. Suppose CA has succeeded in constructing a first axis, such that species appear one after the other along that axis as in a species packing model. Then a possible second axis is obtained by folding the first axis in the middle and bringing the ends together. This axis is a superposition of two species packing models, each with half the gradient length of the first axis. It is a candidate for becoming the second axis, because it has *no linear correlation* with the first CA-axis yet has as much as half the gradient length of the first axis (Jongman *et al.*, 1987). The folded axis by itself thus “explains” a part of the variation in the species data, even though when taken jointly with the first axis it contributes nothing. Even if there is a strong second gradient, CA will not associate it with the second axis if it separates the species less than a folded first axis. As a result of the arch effect, the two-dimensional CA-solution is generally not a good approximation to the ML-solution (two-dimensional Gaussian ordination).

Hill and Gauch (1980) developed detrended correspondence analysis (DCA) as a heuristic modification of CA designed to remedy both the edge effect and the arch effect. The edge effect is removed in DCA by non-linear rescaling of the axis. Assuming a species packing model with randomly distributed species’ optima, Hill and Gauch (1980) noted that the variance of the optima of the species present at a site (the “within-site variance”) is an estimate of the average response curve breadth of those species (they used the standard deviation as a measure of breadth, which is about equal to tolerance as we define it). Because of the edge effect, the species’ curves before rescaling are narrower near the ends of the axis than in the middle, and the within-site variance is correspondingly smaller in sites near the ends of the axis than in sites in the middle. The rescaling therefore attempts to equalize the within-site variance at all points along the ordination axis by dividing the axis into small segments, expanding the segments with sites with small within-site variance, and contracting the segments with sites with large within-site variance. The site scores are then calculated as weighted averages of the species scores and the scores are standardized such that the within-site variance is equal to 1.

Hill and Gauch (1980) defined the length of the ordination axis to be the range of the site scores. This length is expressed in “standard-deviation units” (SD). The tolerance of the species’ curves along the rescaled axis are close to 1, and each curve therefore rises and falls over about 4 SD. Sites that differ by 4 SD can thus be expected to have no species in common. Even if non-linear resealing is not used, one can still set the average within-site variance of the species scores along a CA-axis equal to 1 by linear rescaling (Hill, 1979; Ter Braak, 1987b), so as to ensure that this useful interpretation of the length of the axis still approximately holds.

The arch effect, a more serious problem in CA, is removed in DCA by the heuristic method of “detrending-by-segments”. This method ensures that at any point along the first ordination axis, the mean value of the site scores on subsequent axes is approximately zero. In order to achieve this, the first axis is divided into a number of segments and the trial site scores are adjusted within each segment by subtracting their mean after some smoothing across segments. Detrending-by-segments is built into the reciprocal averaging algorithm, and replaces Step 3b. Subsequent axes are derived similarly by detrending with respect to each of the existing axes.

DCA often works remarkably well in practice (Hill and Gauch, 1980; Gauch *et al.*, 1981). It has been critically evaluated in several recent simulation studies. Ter Braak (1985) showed that DCA gave a much closer approximation to ML Gaussian ordination than CA did, when applied to simulated data based on a two-dimensional species packing model in which species have identically shaped Gaussian surfaces and the optima and site scores are uniformly distributed in a rectangle. This improvement was shown to be mainly due to the detrending, not to the non-linear rescaling of axes. Kenkel and Orlóci (1986) found that DCA performed substantially better than CA when the two major gradients differed in length, but also noted that DCA sometimes “collapsed and distorted” CA results when there were: (a) few species per site, and (b) the gradients were long (we believe (a) to be the real cause of the collapse). Minchin (1987) further found that DCA can flatten out some of the variation associated with one of the underlying gradients. He ascribed this loss of information to an instability in the detrending-by-segments method. Pielou (1984, p. 197) warned that DCA is “overzealous” in correcting the “defects” in CA, and “may sometimes lead to the unwitting destruction of ecologically meaningful information”. Minchin’s (1987) results indicate some of the conditions under which such loss of information can occur.

DCA is popular among practical field ecologists, presumably because it provides an effective approximate solution to the ordination problem for a unimodal response model in two or more dimensions—given that the data are reasonably representative of sections of the major underlying environmental gradients. Two modifications might increase its robustness with respect to the problems identified by Minchin (1987). First, non-linear, rescaling aggravates these problems; since the edge effect is not too serious, we advise against the routine use of non-linear rescaling. Second, the arch effect needs to be removed (as Heiser, 1987, also noted), but this can be done by a more stable, less “zealous” method of detrending which was also briefly mentioned by Hill and Gauch (1980): namely detrending-by-polynomials. Under the one-dimensional Gaussian model, it can be shown that the second CA-axis is a quadratic function of the first axis, the third axis is a cubic function of the first axis, and so on (Hill, 1974; Iwatsubo, 1984).



Detrending-by-polynomials can be incorporated into the reciprocal averaging algorithm by extending Step 3b such that the trial scores are not only made uncorrelated with the previous axes, but are also made uncorrelated with polynomials of the previous axes. The limited experience so far suggests that detrending up to fourth-order polynomials should be adequate. In contrast with detrending-by-segments, the method of detrending-by-polynomials removes only specific defects of CA that are now theoretically understood.

#### D. Constrained Ordination

Just as CA/DCA is an approximation to ML Gaussian ordination, so is canonical correspondence analysis (CCA) an approximation to ML Gaussian canonical ordination (Ter Braak, 1986). CCA is a modification of CA in which the ordination axes are restricted to be weighted sums of the environmental variables, as in Eq. (4). CCA can be obtained from CA as redundancy analysis was obtained from PCA. An algorithm can be obtained by adding to the CA algorithm an extra multiple regression step. The only difference from Step 3a of redundancy analysis (see Section II.E) is that the sites must be weighted in the regression proportional to their site total  $y_{+i}$  (Ter Braak, 1986). CCA can also be obtained as the solution of an eigenvalue problem (Ter Braak, 1986). It is closely related to “redundancy analysis for qualitative variables” (Israëls, 1984) but has a different rationale and is applied to a different type of data.

In constrained ordination the constraints always become less strict as more environmental variables are included. If  $q \geq n - 1$ , then there are no real constraints, and CA and CCA become equivalent. As in CA, the edge effect in CCA is a minor problem that is best left untreated. Detrending may sometimes be required to remove the arch effect, i.e. to prevent CCA from selecting weighted sums of environmental variables that are approximately polynomials of previous axes. Detrending-by-segments does not work very well here for technical reasons; detrending-by-polynomials is better founded and more appropriate (see Appendix and Ter Braak, 1987b). However, the arch effect in CCA can be eliminated much more elegantly, simply by dropping superfluous environmental variables (Ter Braak, 1987a). Variables that are highly correlated with the “arched” axis (often the second axis) are the most likely to be superfluous. If the number of environmental variables is small enough for the relationship of individual variables to the ordination axes to be significant, the arch effect is not likely to occur at all.

CCA can be sensitive to deviant sites, but only when they are outliers with regard to both species composition and environment. When



realistically few environmental variables are included, CCA is thus more robust than CA in this respect too.

CCA leads to an ordination diagram that simultaneously displays (a) the main patterns of community variations, as far as these reflect environmental variation, and (b) the main pattern in the weighted averages (not correlations as in redundancy analysis) of each of the species with respect to the environmental variables (Ter Braak, 1986, 1987a). CCA is thus intermediate between CA and separate WA calculations for each species. Geometrically, the separate WA calculations give each species a point in the  $q$ -dimensional space of the environmental variables, which indicates the centre of the species' distribution. CCA attempts to provide a low-dimensional representation of these centres; CCA is thus also a constrained form of WA, in which the weighted averages are restricted to lie in a low-dimensional subspace.

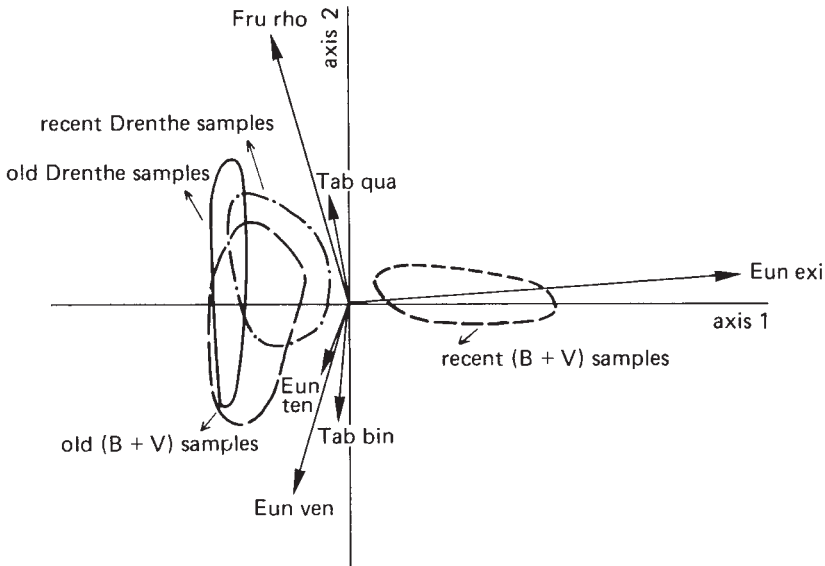
Like redundancy analysis, CCA can be used with dummy "environmental" variables to provide an ordination constrained to show maximum separation among pre-defined groups of samples. This special case of CCA is described, for example, by Feoli and Orłóci (1979) under the name of "analysis of concentration", by Greenacre (1984, Section 7.1) and by Gasse and Tekaiia (1983).

## V. ORDINATION DIAGRAMS AND THEIR INTERPRETATION

The linear ordination techniques (PCA and redundancy analysis) and, the ordination techniques based on WA (CA/DCA and CCA) represent community data in substantially different ways. We focus on two-dimensional ordination diagrams, as these are the easiest to construct and to inspect, and illustrate the interpretation of each type of diagram with an example.

### A. Principal Components: Biplots

PCA fits planes to each species' abundances in the space defined by the ordination axes. The species' point  $(b_{k1}, b_{k2})$  may be connected with the origin  $(0,0)$  to give an arrow (Figure 3). Such a diagram, in which sites are marked by points and species by arrows is called a "biplot" (Gabriel, 1971). There is a useful symbolism in this use of arrows: the arrow points in the direction of maximum variation in the species' abundance, and its length is proportional to this maximum rate of change. Consequently, species on the edge of the diagram (far from the origin) are the most important for



**Figure 3** Biplot based on principal components analysis of diatom assemblages from Dutch moorland pools (schematic after van Dam *et al.*, 1981). The arrows for the six most frequent species and the regions where different categories of samples lie jointly display the approximate community composition in each of the regions (old, c. 1920; recent, 1978; B+V, from the province of Brabant and the Veluwe). Abbreviations: Eun exi, *Eunotia exigua*; Eun ten, *Eunotia tenella*; Eun ven, *Eunotia veneris*; Fru rho, *Frustulia rhomboides* var. *saxonica*; Tab bin, *Tabellaria binalis*; Tab qua, *Tabellaria quadriseptata*.

indicating site differences; species near the centre are of minor importance. Ter Braak (1983) provides more detailed, quantitative rules for interpreting PCA ordination diagrams.

van Dam *et al.* (1981) applied PCA to data consisting of diatom assemblages from 16 Dutch moorland pools, sampled in the 1920s and again in 1978, to investigate the impact of acidification on these shallow water bodies. Ten clearwater (non-humic) pools were situated in the province of Brabant and on the Veluwe and six brownwater (humic) pools in the province of Drenthe. Figure 3 displays the major variation in the data. The arrow of *Eunotia exigua* indicates that this species increases strongly along the first principal component: *E. exigua* is abundant in the recent Brabant and Veluwe samples, which lie on the right-hand side of the diagram, and rare in the remaining samples, which lie more to the left. The second axis accounts for some of the difference among the old and recent samples from Drenthe. These groups differ in the abundances of *Frustulia*

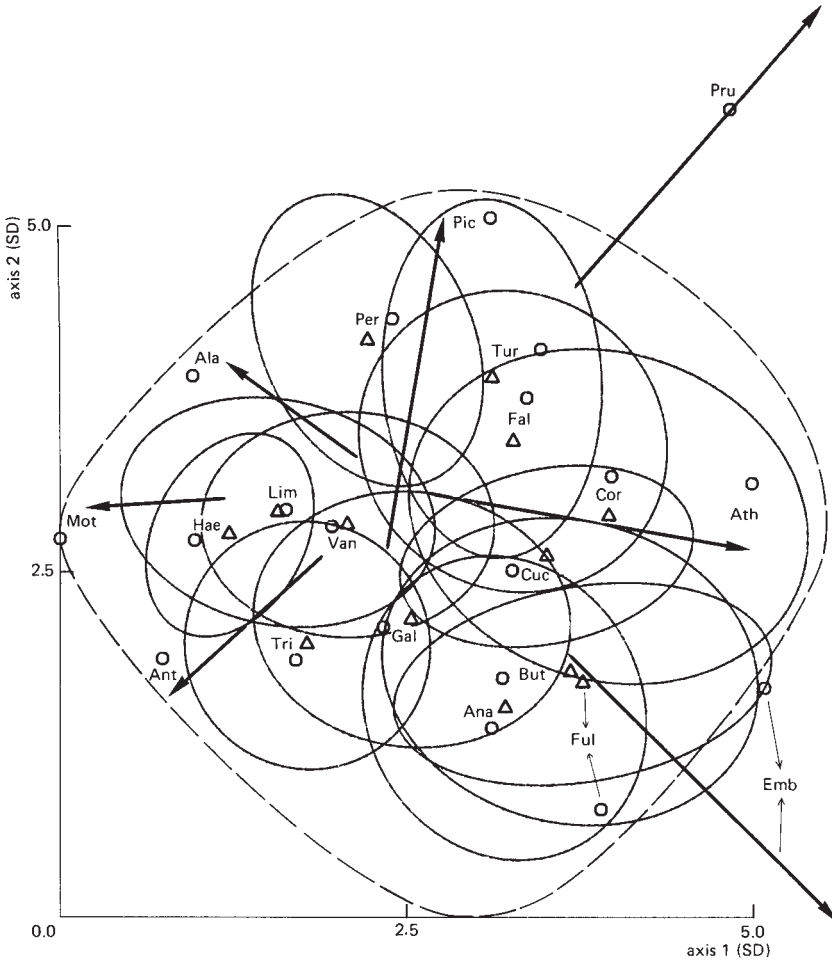
*rhomboides* var. *saxonica*, *Tabellaria quadrisepitata*, *Eunotia tenella*, *Tabellaria binalis*, and *Eunotia veneris*, as shown by the directions of the arrows for these species in Figure 3. As *E. exigua* is acidobiontic and the first principal component is strongly correlated with the sulphate concentration of the 1978 samples, this component clearly depicts the impact of acidification of the moorland pools in Brabant and the Veluwe (and to a smaller extent also in Drenthe). Thus van Dam *et al.* (1981) used PCA to summarize the changes in diatom composition between the 1920s and 1978. PCA helped them to detect that the nature of the change differed among provinces, hence stressing the importance for diatoms of the distinction between clearwater and brownwater pools.

### B. Correspondence Analysis: Joint Plots

In CA and DCA both sites and species are represented by points, and each site is located at the centre of gravity of the species that occur there. One may therefore get an idea of the species composition at a particular site by looking at “nearby” species points. Also, in so far as DCA approximates the fitting of Gaussian (logit) surfaces (Figure 2), the species points are approximately the optima of these surfaces; hence the abundance or probability of occurrence of a species tends to decrease with distance from its location in the diagram.

Figure 4 illustrates this interpretation of the species' points as optima in ordination space. DCA was applied to presence–absence data on 51 bird species in 526 contiguous, 100 m × 100 m grid-cells in an area with pastures and scattered woodlots in the Rhine valley near Amerongen, the Netherlands (Opdam *et al.*, 1984). Figure 4 shows the DCA scores of the 20 most frequent species by small circles, and the outline (dashed) of the region in which the scores for the grid-cells fall (the individual grid-cells are not shown, to avoid crowding). Opdam *et al.* (1984) interpreted the first axis, of length 5.6 SD, as a gradient from open to closed landscape and the second axis, of length 5.3 SD, as a gradient from wet to drier habitats.

In order to test the interpretation of species' scores as optima, we fitted a response surface for each species by logit regression using Eq. (6) with the first and the second DCA-axes as the predictor variables  $x_1$  and  $x_2$ . For 13 of the 20 bird species, the fitted surface had a maximum. The optimum was calculated for each of these species by Eq. (7) and plotted as a triangle in Figure 4. The fitted optima lie close to the DCA scores. The regression analysis also allowed us to estimate species' tolerances in ordination space: these are indicated in Figure 4 by ellipses representing the region within which each species occurs with at least half of its maximum probability, according to the fitted surface.



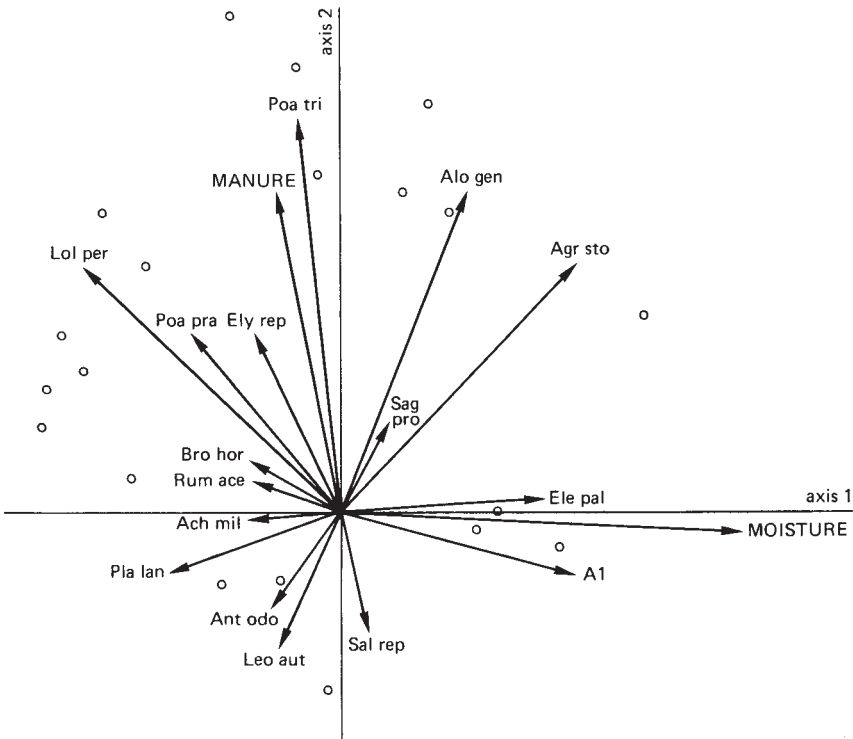
**Figure 4** Joint plot based on detrended correspondence analysis (DCA) of bird species communities in the Rhine valley near Amerongen, the Netherlands (data from *Opdam et al., 1984*), displaying the major variation in bird species composition across the landscape. This plot shows the DCA-scores (○) of the 20 most frequent species and the region in which the samples fall (— — —). Also shown are optima (△) and lines of equal probability for the 13 species whose probability surfaces had clear maxima (as fitted by Gaussian logit regression), and arrows representing directions of increase for the seven species whose probability surfaces were monotonic. Abbreviations: Ala, *Alauda arvensis*; Ana, *Anas platyrhynchos*; Ant, *Anthus pratensis*; Ath, *Athene noctua*; But, *Buteo buteo*; Cor, *Corvus corone*; Cuc, *Cuculus canorus*; Emb, *Emberiza schoeniclus*; Fal, *Fatco tinnunculus*; Ful, *Fulica atra*; Gal, *Gallinago gallinago*; Hae, *Haematopus ostralegus*; Lim, *Limosa limosa*; Mot, *Motacilla flava flava*; Per, *Perdix perdix*; Pic, *Pica pica*; Pru, *Prunella modularis*; Tri, *Tringa totanus*; Tur, *Turdus merula*; Van, *Vanellus vanellus*.

The fitted surfaces for the remaining seven species had a minimum or saddle point, suggesting that their optima are located well outside the sampled range. For these species we fitted a “linear” logit surface by setting  $b_2$ ,  $b_4$  and  $b_5$  in Eq. (6) to zero. The direction of steepest increase of each of the fitted surfaces is indicated in Figure 4 by an arrow through the centroid of the site points; the beginning and end points of each arrow correspond to fitted probabilities of 0.1 and 0.9 respectively. As expected from our interpretation of DCA, these arrows point more or less in the same direction as the DCA scores of the corresponding species (Figure 4).

In contrast to the PCA-diagram, the species points on the edge of the CA- or DCA-diagram are often rare species, lying there either because they prefer extreme (environmental) conditions, or (very often) because their few occurrences by chance happen to fall in sites with extreme conditions; one cannot decide between these possibilities without additional data. Such peripheral species have little influence on the analysis and it is often convenient not to display them at all. Furthermore, species near the centre of the diagram may be ubiquitous, unrelated to the ordination axes, bimodal, or in some other way not fitting a unimodal response model—or they may be genuinely specific with a habitat-optimum near the centre of the sampled range of habitats. The correct interpretation may be found by the kind of secondary analysis shown in Figure 4, or more straightforwardly just by plotting the species’ abundances in ordination space.

### C. Redundancy Analysis

In redundancy analysis sites are indicated by points, and both species and environmental variables are indicated by arrows whose interpretation is similar to that of the arrows in the PCA biplot. The pattern of abundance of each species among the sites can be inferred in exactly the same way as in a PCA biplot, and so may the direction of variation of each environmental variable. One may also get an idea of the correlations between species’ abundances and environmental variables. Arrows pointing in roughly the same direction indicate a high positive correlation, arrows crossing at right angles indicate near-zero correlation, and arrows pointing in opposite directions indicate high negative correlation. Species and environmental variables with long arrows are the most important in the analysis; the longer the arrows, the more confident one can be about the inferred correlation. (It is assumed here that for the purpose of the ordination diagram the environmental variables have been standardized to zero mean and unit variance, so as to make the lengths of arrows comparable.) Jongman *et al.* (1987) provide more quantitative rules for interpreting the ordination diagrams derived from redundancy analysis.



**Figure 5** Biplot based on redundancy analysis of vegetation with respect to three environmental variables (quantity of manure, soil moisture and thickness of the A1 horizon) in dune meadows (○) on the island of Terschelling, The Netherlands. The arrows for plant species and environmental variables display the approximate (linear) correlation coefficients between plant species and the environmental variables. Abbreviations: Ach mil, *Achillea millefolium*; Agr sto, *Agrostis stolonifera*; Alo gen, *Alopecurus geniculatus*; Ant odo, *Anthoxanthum odoratum*; Bro hor, *Bromus hordaceus*; Ele pal, *Eleocharis palustris*; Ely rep, *Elymus repens*; Leo aut, *Leontodon autumnalis*; Loi per, *Lolium perenne*; Pla lan, *Plantago lanceolata*; Poa pra, *Poa pratensis*; Poa tri, *Poa trivialis*; Rum ace, *Rumex acetosa*; Sag pro, *Sagina procumbens*; Sal rep, *Salix repens*.

The data we use to illustrate redundancy analysis were collected to study the relation between the vegetation and management of dune meadows on the island of Terschelling, The Netherlands (M. Batterink and G. Wijffels, unpublished). **Figure 5** displays the main variation in the vegetation in relation to three environmental variables (thickness of the A1 horizon, moisture content of the soil and quantity of manuring). The arrows for *Poa trivialis* and *Elymus repens* make small angles with the arrow for manuring;

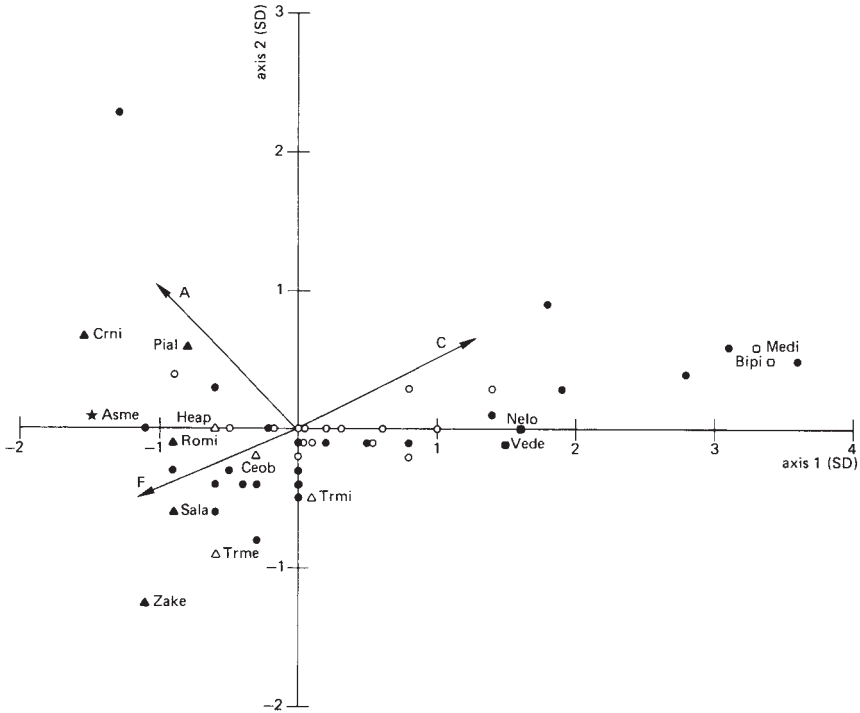
these species are inferred to be positively correlated with manuring. *Salix repens* and *Leontodon autumnalis* have arrows pointing in directions roughly opposite to that of manuring, and are inferred to be negatively correlated with manuring. The former species are thus most abundant in the heavily manured meadows of standard farms (positioned at the top of the diagram), whereas the latter species are most abundant in the unmanured meadows (owned by the nature conservancy and positioned at the bottom of the diagram). The relationships of the species with moisture and thickness of the A1 horizon can be inferred in a similar way. The short arrows for *Bromus hordaceus* and *Sagina procumbens*, for example, indicate that their abundance is not so much affected by moisture, manure and thickness of the A1 horizon. Redundancy analysis can summarize the species–environment relationships in such an informative way, because the gradients are short ( $\approx 2SD$ : Ter Braak, 1987b).

#### D. Canonical Correspondence Analysis

In CCA, since species are assumed to have unimodal response surfaces with respect to linear combinations of the environmental variables, the species are logically represented by points (corresponding to their approximate optima in the two-dimensional environmental subspace), and the environmental variables by arrows indicating their direction and rate of change through the subspace.

Purata (1986, and unpublished results) applied CCA to plant species abundance data from 40 abandoned cultivation sites within Mexican tropical rain forest. Data were available for 24 of these sites on the regrowth age (A), the length of the cropping period in the past (C), and the proportion of the perimeter that had remained forested (F). These three variables were used as environmental variables in CCA. The remaining 16 sites were entered as “passive” sites, to be positioned with respect to the CCA axes according to their floristic composition in relation to the “active” sites.

Figure 6 illustrates the results. The first axis, with length 4.7 SD, was interpreted as an indicator of the general trend of secondary succession. The direction of the arrow for regrowth age shows that this trend runs broadly from right to left. The species' locations are consistent with their life-history characteristics: the trend of succession runs from ruderals (to the right), through pioneer shrubs and trees, to late-secondary canopy dominants and shade-tolerant understorey species (to the left). The directions of the other two arrows in relation to axis 1 show that a long cropping period delays succession, while an extensive forested perimeter accelerates succession. Axis 2 (3.0 SD) may (more speculatively) differentiate species whose



**Figure 6** Ordination diagram based on canonical correspondence analysis of successional plant communities with respect to three environmental variables (regrowth age A, length of cropping period C, and extent of forested perimeter F) on abandoned cultivation sites within Mexican tropical rain forest (Purata, 1986 and unpublished). ●, sites with environmental data; ○, sites added “passively” on the basis of floristic composition. The species shown are a selection among the 285 included in the analysis. □, denotes ruderals; ■, pioneer shrubs; △, pioneer trees; ▲, late-secondary canopy trees; and ★, an understory palm. Abbreviations: Asme, *Astrocaryum mexicanum*; Bipi, *Bidens pilosa*; Ceob, *Cecropia obtusifolia*; Crni, *Croton nitens*; Heap, *Heliocarpus appendicalatus*; Medi, *Melampodium divaricatum*; Nelo, *Neurolaena lobata*; Pial, *Piper amalago*; Romi, *Robinsonella mirandae*; Sala, *Sapium lateriflorum*; Trme, *Trichospermum mexicanum*; Trmi, *Trema micrantha*; Vede, *Vernonia deppeana*; Zake, *Zanthoxylum kellermanii*.

establishment is favoured by the presence of mature forest around the site from those that simply require a long time to grow.

CCA also allows the computation of unconstrained, “residual” axes summarizing floristic variation that remains after the effect of the environmental variables has been taken out. In Purata’s study, the successive eigenvalues of the first three (constrained) CCA axes were 0.49, 0.34 and 0.18. (There cannot be more constrained axes than there are environmental



variables). The first residual axis gave an eigenvalue of 0.74, showing that at least as much floristic variation was *not* explained by the environmental variables. In our experience, terrestrial community data commonly give a residual eigenvalue as large as the first constrained eigenvalue, however carefully the environmental variables are chosen. Thus DCA and CCA tend to give different ordinations, and CCA is more powerful in detecting relationships between species composition and environment.

## VI. CHOOSING THE METHODS

### A. Which Response Model?

Regression methods can fit response models with a wide variety of shapes. The linear and Gaussian-like models are convenient starting points; more complex shapes can be fitted by adding further parameters, if the data are sufficiently detailed to support it. Other species may be used as additional explanatory variables if the specific aim is to detect species interactions (Fresco, 1982). The shapes of the response functions may be made even more general by applying Box–Cox transformations to the explanatory variables (Bartlein *et al.*, 1986) or still more general by fitting splines (Smith, 1979). Even with all these modifications, regression can still be done with standard packages for Generalized Linear Modelling.

After species response curves or surfaces have been fitted by regression, calibration based on the maximum likelihood principle can be used to make inferences about the environment from community data. If the surfaces fitted by regression have complex shapes, then calibration by *numerical* maximization of the likelihood may be problematic. But even then, if there are only a few environmental variables involved, the “most likely” combination of environmental values can be searched for on a grid across the environmental space (Atkinson *et al.*, 1986; Bartlein *et al.*, 1986). So the type of response model used in both regression and calibration should generally be guided by the characteristics and resolution of the data, and inspection of the data and the residuals after regression should show whether the model being used is adequate for the purpose.

In contrast to regression and calibration, the ordination problem requires the simultaneous estimation of large numbers of parameters and cannot be solved practically without some constraints on the structure one wants to fit. That these constraints may seem unduly restrictive simply shows that there are limits to what ordination can achieve. The number of ordination axes to be extracted must be small, and the type of response model must be restricted, in order to permit a solution. For example, it seems necessary

to disregard the possibility of bimodal species distributions (Hill, 1977). Certainly bimodal distributions sometimes occur, but ordination has to assume that species “on average” have simple distributions—otherwise, the problem would be insoluble; the utility of ordination techniques depends on them being robust with respect to departures from the simple models on which they are based. The Gaussian model seems to be of the right order of complexity for ordination of ecological data, but the full second-degree model of Eq. (6) is already difficult to fit (Kooijman, 1977; Goodall and Johnson, 1982). The Gaussian model with circular contour lines and equal species tolerances, i.e. the unfolding model, might provide a good compromise between practical solubility and realism in ordination. Promising algorithms for unfolding are developed by Heiser (1987) and DeSarbo and Rao (1984). DCA provides a reasonably robust approximation to ML Gaussian ordination and requires far less computing time. Similarly, ML Gaussian canonical ordination is technically feasible, but CCA provides a practical and robust approximation to it.

Non-linear methods are appropriate if a reasonable number of species have their optima located within the data set. If the gradient length is reduced to less than about 3 SD, the approximations involved in WA become worse and ultimately (if the gradient length is less than about 1.5 SD) the methods yield poor results because most species are behaving monotonically over the observed range. Thus if the community variation is within a narrow range, the linear ordination methods—PCA and redundancy analysis—are appropriate. If the community variation is over a wider range, non-linear ordination methods—including DCA and CCA—are appropriate.

## **B. Direct or Indirect?**

Direct gradient analysis allows one to study the part (large or small) of the variation in community composition that can be explained by a particular set of environmental variables. In indirect gradient analysis attention is first focused on the major pattern of variation in community composition; the environmental basis of this pattern is to be established later. If the relevant environmental data are to hand, the direct approach—either fitting separate response surfaces by regression for each major species, or analysing the overall patterns of the species–environment relationship by constrained ordination—is likely to be more effective than the traditional indirect approach. However, indirect gradient analysis does have the advantage that no prior hypothesis is needed about what environmental variables are relevant. One does not need to measure the environmental variables in advance, and one can use informal field knowledge to help interpret the

patterns that emerge—hence the emphasis in the literature on ordination as a technique for “hypothesis generation”, the implication being that experimental or more explicit statistical approaches can be used for subsequent hypothesis testing. This distinction is not hard and fast, but it does draw attention to the strengths and limitations of indirect gradient analysis.

In [Section V.D](#), we showed in passing how an indirect gradient analysis can be carried out *after* a direct gradient analysis in order to summarize the community variation that remains after known effects have been removed. When the known environmental variables are not the prime object of study, they are called concomitant variables ([Davies and Tso, 1982](#)) or covariables. It would be convenient to solve for the residual (unconstrained) axes without having to extract all the constrained axes first. Fortunately, this is straightforward. In the iterative algorithm for PCA and CA, one simply extends Step 3b such that the trial scores are not only made uncorrelated with any previous axis (if present) but are also made uncorrelated with all specified covariables (see Appendix for details). In this way the effects of the covariables are partialled out from the ordination; hence the name “partial ordination”. The theory of “partial components analysis” and “partial correspondence analysis”, as we call these extensions of PCA and CA, is given by [Gabriel \(1978, theorem 3\)](#) and [Ter Braak \(1988\)](#), respectively. [Swaine and Greig-Smith \(1980\)](#) used partial components analysis to obtain an ordination of within-plot vegetation change in permanent plots. Partial correspondence analysis, or its detrended form, would be more appropriate if the gradients were long.

### C. Direct Gradient Analysis: Regression or Constrained Ordination?

Whether to use constrained ordination (multivariate direct gradient analysis) instead of a series of separate regressions (the traditional type of direct gradient analysis) depends on whether or not there is any advantage in analysing all the species simultaneously. Both constrained and unconstrained ordination assume that the species react to the *same* composite gradients of environmental variables, while in regression a separate composite gradient is constructed for each species. Regression can therefore allow more detailed descriptions and more accurate prediction and calibration, if properly carried out (with due regard to its statistical assumptions) and if sufficient data are available. However, ecological data that are collected over a large range of habitat variation require non-linear models, and building good non-linear models by regression is demanding in time and computation. In CCA the composite gradients are linear combinations of environmental variables and the non-linearity enters through a unimodal

response model with respect to a few composite gradients, taken care of in CCA by the procedure of weighted averaging. Constrained ordination is thus easier to apply, and requires less data, than regression; it provides a summary of the species–environment relationship, and we find it most useful for the exploratory analysis of large data sets.

Constrained ordination can also be carried out *after* regression, in order to relate the residual variation to other environmental variables. This type of analysis, called “partial constrained ordination”, is useful when the explanatory (environmental) variables can be subdivided in two sets, a set of covariables—the effects of which are not the prime object of study—and a further set of environmental variables whose effects are of particular interest.

For example, in the illustration of [Section V.C](#), the study was initiated to investigate differences in vegetation among dune meadows that were exploited under different management regimes (standard farming, biodynamical farming, nature management, among others). Standard CCA showed systematic differences in vegetation among management regimes. A further question is then whether these differences can be fully accounted for by the environmental variables moisture, quantity of manure and thickness of the A1 horizon, whose effects are displayed in [Figure 5](#), or whether the variation that remains after fitting the three environmental variables (three constrained ordination axes) is systematically related to management regimes. This question can be tackled using partial constrained ordination, with the three environmental variables as covariables, and a series of dummy variables (for each of the management regimes) as the variables of interest.

Technically, partial constrained ordination can be carried out by any computer program for constrained ordination. The usual environmental variables are replaced by the residuals obtained by regressing each of the variables of interest on the covariables (see Appendix). [Davies and Tso \(1982\)](#) gave the theory behind partial redundancy analysis; [Ter Braak \(1988\)](#) derived partial canonical correspondence analysis as an approximation to “partial Gaussian canonical ordination”.

Partial constrained ordination has the same essential aim as [Carleton’s \(1984\)](#) residual ordination, i.e. to determine the variation in the species data that is uniquely attributable to a particular set of environmental variables, taking into account the effects of other (co-) variables; however, Carleton’s method is somewhat less powerful, being based on a pre-existing DCA which may already have removed some of the variation of interest. Partial constrained ordination is, by contrast, a true direct gradient analysis technique which seems promising, e.g. for the analysis of permanent plot data (effects of time, with location and/or environmental data as covariables), and a variety of other applications in which effects of particular environmental variables are to be sorted out from the “background” variation imposed by other variables.

## VII. CONCLUSIONS

Regression, calibration, ordination and constrained ordination are well-defined statistical problems with close interrelationships. Regression is the tool for investigating the nature of individual species' response to environment, and calibration is the tool for (later) inferring the environment from species composition at an individual site. Both tools come in various degrees of complexity. The simplest are linear and WA regression and calibration. The linear methods are applicable over short ranges of environment, where species' abundance appears to vary monotonically with variation in the environment. The WA methods are applicable over wider ranges of environment; WA regression is a crude method to estimate each species' optimum, and WA calibration just averages the optima of the species that are present. WA works with presence-absence data. If abundances are available, they provide the weights. These WA techniques can be shown to give approximate estimates of the species' optima and environmental values when the species' response surfaces (the relationships between the species' abundance, or probability of occurrence, and the environmental variables) are Gaussian (or for probabilities, Gaussian-logit) in form. Gaussian regression and calibration are also possible, but the WA techniques are simpler and are approximations to the Gaussian methods.

These simple tools are suitable when there are many species of interest and the exact form of the response surface is not critical, and they are very easy to use. If the form of the response surfaces *is* critical, more complex models can be fitted using Generalized Linear Modelling (for regression) and maximum likelihood techniques (for calibration). These more complex tools are becoming important in the theoretical study of species-environment relationships (Austin, 1985) and environmental dynamics (Bartlein *et al.*, 1986). Naturally, they require skilled users who are aware of their statistical assumptions, limitations and pitfalls.

Ordination and constrained ordination can be related to the simpler methods of regression and calibration. Ordination is the tool for exploratory analysis of community data with no prior information about the environment. Constrained ordination is the equivalent tool for the analysis of community variation in relation to environment. Both implicitly assume a common set of environmental variables and a common response model for all of the species. (Without these simplifying assumptions, they could not work; such major simplifications of data can only be achieved at the expense of some realism.) The basic ordination techniques are PCA and CA. PCA constructs axes that are as close as possible to a linear relationship with the species. These axes can be found by a converging sequence of alternating linear regressions and calibrations. Each axis after the first is obtained by

partialling out linear relationships with the previous axis. CA is mathematically related to PCA, but has a very different effect. CA axes can be found by a converging sequence of WA regressions and calibrations. In CA, axes after the first are obtained analogously with PCA; in DCA they are obtained by removing all trends, linear or non-linear, with respect to previous axes. CA suffers from the arch effect, which DCA eliminates. DCA is a reasonably robust approximation to Gaussian ordination, in which the axes are constructed so that the species response curves with respect to the axes are Gaussian in form. Gaussian ordination is feasible but not convenient. DCA is much more practical. But there are problems with the detrending, and the method can break down when the connections between sites are too tenuous. Some modifications—including an improved method of detrending—may improve DCA's robustness; alternatively, some forms of nonmetric multidimensional scaling may be more robust (Kenkel and Orłóci, 1986; Minchin, 1987).

Constrained ordination methods have the added constraint that the ordination axes must be linear combinations of environmental variables. This constraint can be implemented as an extra multiple regression step in the general iterative ordination algorithm. PCA then becomes redundancy analysis (a more practical alternative to canonical correlation), Gaussian ordination becomes Gaussian canonical ordination, and CA becomes CCA (Table 2). The constraint makes Gaussian canonical ordination somewhat more stable than its unconstrained equivalent, but still CCA provides a much more practical alternative. All these constrained methods are most powerful if the number of environmental variables is small compared to the number of sites. Then the constraints are much stronger than in normal ordination, and the common problems of ordination (such as the arch effect, the need for detrending and the sensitivity to deviant sites) disappear.

Often, community–environment relationships have been explored by “indirect gradient analysis”—ordination, followed by interpretation of the axes in terms of environmental variables. But if the environmental data are to hand, constrained ordination (“multivariate direct gradient analysis”) provides a more powerful means to the same end. Hybrid (direct/indirect) analyses are also possible. In partial ordination and partial constrained ordination, the analysis works on the variation that remains after the effects of particular environmental, spatial or temporal “covariables” have been removed.

The choice between linear and non-linear ordination methods is not a matter of personal preference. Where gradients are short, there are sound statistical reasons to use linear methods. Gaussian methods break down, and edge effects in CA and related techniques become serious; the representation of species as arrows becomes appropriate. As gradient lengths

increase, linear methods become ineffective (principally through the “horseshoe effect”, which scrambles the order of samples along the first axis as well as creates a meaningless second axis); Gaussian methods become feasible, and CA and related techniques become effective. The representation of species as points, representing their optima, becomes informative. The range 1.5–3 SD for the first axis represents a “window” over which both PCA and CA/DCA, or both redundancy analysis and CCA, can be used to good effect.

## ACKNOWLEDGEMENTS

We thank Dr M. P. Austin, Dr P. J. Bartlein, Professor L. C. A. Corsten, J. A. Hoekstra, Dr P. Opdam and Dr H. van Dam for comments on the manuscript. Our collaboration was supported by a Netherlands Science Research Council (ZWO) grant to I.C.P. and a Swedish Natural Science Research Council (NFR) grant to the project “Simulation Modelling of Natural Forest Dynamics”. We also thank Dr S. E. Purata V. for supplying unpublished results.

## REFERENCES

- Alderdice, D.F. (1972) Factor combinations: responses of marine poikilotherms to environmental factors acting in concert. In: *Marine Ecology* (Ed. by O. Kinne), Vol. 1, Part 3, pp. 1659–1722. John Wiley, New York.
- Ås (1985) Biological Community patterns in insular environments. *Acta Unit. Ups.* **792**, 1–55.
- Atkinson, T.C., Briffa, K.R., Coope, G.R., Joachim, M.J. and Perry, D.W. (1986) Climatic calibration of coleopteran data. In: *Handbook of Holocene Palaeoecology and Palaeohydrology* (Ed. by B.E. Berglund), pp. 851–858. John Wiley, Chichester.
- Austin, M.P. (1971) Role of regression analysis in plant ecology. *Proc. Ecol. Soc. Austr.* **6**, 63–75.
- Austin, M.P. (1985) Continuum concept, ordination methods, and niche theory. *Ann. Rev. Ecol. Syst.* **16**, 39–61.
- Austin, M.P. and Cunningham, R.B. (1981) Observational analysis of environmental gradients. *Proc. Ecol. Soc. Austr.* **11**, 109–119.
- Austin, M.P., Cunningham, R.B. and Fleming, P.M. (1984) New approaches to direct gradient analysis using environmental scalars and statistical curve-fitting procedures. *Vegetatio* **55**, 11–27.
- Balloch, D., Davies, C.E. and Jones, F.H. (1976) Biological assessment of water quality in three British rivers: the North Esk (Scotland), the Ivel (England) and the Taf (Wales). *Wat. Pollut. Control* **75**, 92–114.



- Bartlein, P.J., Webb, T. III and Fleri, E. (1984) Holocene climatic changes in the Northern Midwest: pollen-derived estimates. *Quaff. Res.* **22**, 361–374.
- Bartlein, P.J., Prentice, I.C. and Webb, T. III (1986) Climatic response surfaces from pollen data for some eastern North American taxa. *J. Biogeogr.* **13**, 35–57.
- Battarbee, R.W. (1984) Diatom analysis and the acidification of lakes. *Phil. Trans. Roy. Soc. London Ser. B* **305**, 451–477.
- Bloxom, B. (1978) Constrained multidimensional scaling in N spaces. *Psychometrika* **43**, 397–408.
- Böcker, R., Kowarik, I. and Bornkamm, R. (1983) Untersuchungen zur Anwendung der Zeigerwerte nach Ellenberg. *Verh. Ges. Oekol.* **11**, 35–56.
- Brown, G.H. (1979) An optimization criterion for linear inverse estimation. *Technometrics* **21**, 575–579.
- Brown, P.J. (1982) Multivariate calibration. *J. Roy. Statist. Soc. B* **44**, 287–321.
- Carleton, T.J. (1984) Residual ordination analysis: a method for exploring vegetation environment relationships. *Ecology* **65**, 469–477.
- Chandler, J.R. (1970) A biological approach to water quality management. *Wat. Pollut. Control* **69**, 415–421.
- Charles, D.F. (1985) Relationships between surface sediment diatom assemblages and lakewater characteristics in Adirondack lakes. *Ecology* **66**, 994–1011.
- Coombs, C.H. (1964) *A Theory of Data*. John Wiley, New York.
- Cox, D.R. and Hinkley, D.V. (1974) *Theoretical Statistics*. Chapman and Hall, London.
- Cramer, W. and Hytteborn, H. (1987) The separation of fluctuation and long-term change in vegetation dynamics of a rising sea-shore. *Vegetatio* **69**, 155–167.
- Dargie, T.C.D. (1984) On the integrated interpretation of indirect site ordinations: a case study using semi-arid vegetation in southeastern Spain. *Vegetatio* **55**, 37–55.
- Davies, P.T. and Tso, M.K.-S. (1982) Procedures for reduced-rank regression. *Appl. Statist.* **31**, 244–255.
- Davison, M.L. (1983) *Multidimensional Scaling*. John Wiley, New York.
- De Leeuw, J. and Heiser, W. (1980) Multidimensional scaling with restrictions on the configuration. In: *Multivariate Analysis-V* (Ed. by P.R. Krishnaiah), pp. 501–522. North-Holland, Amsterdam.
- DeSarbo, W.S. and Rao, V.R. (1984) GENFOLD2: a set of models and algorithms for the general unfolding analysis of preference/dominance data. *J. Class.* **1**, 147–186.
- Dobson, A.J. (1983) *Introduction to Statistical Modelling*. Chapman and Hall, London.
- Ellenberg, H. (1979) Zeigerwerte der Gefäßpflanzen Mitteleuropas. *Scripta Geobotanica* **9**, 1–121.
- Fängström, I. and Willén, E. (1987) Clustering and canonical correspondence analysis of phytoplankton and environment variables in Swedish lakes. *Vegetatio* **71**, 87–95.
- Feoli, E. and Feoli Chiapella, L. (1980) Evaluation of ordination methods through simulated coenoclines: some comments. *Vegetatio* **42**, 35–41.
- Feoli, E. and Orlóci, L. (1979) Analysis of concentration and detection of underlying factors in structured tables. *Vegetatio* **40**, 49–54.



- Fresco, L.F.M. (1982) An analysis of species response curves and of competition from field data: some results from heath vegetation. *Vegetatio* **48**, 175–185.
- Gabriel, K.R. (1971) The biplot graphic display of matrices with application to principal component analysis. *Biometrika* **58**, 453–467.
- Gabriel, K.R. (1978) Least squares approximation of matrices by additive and multiplicative models. *J. Roy. Statist. Soc. B* **40**, 186–196.
- Gasse, F. and Tekaia, F. (1983) Transfer functions for estimating paleoecological conditions (pH) from East African diatoms. *Hydrobiologia* **103**, 85–90.
- Gauch, H.G. (1982) *Multivariate Analysis in Community Ecology*. Cambridge Univ. Press, Cambridge.
- Gauch, H.G. and Whittaker, R.H. (1972) Coenocline simulation. *Ecology* **53**, 446–451.
- Gauch, H.G., Chase, G.B. and Whittaker, R.H. (1974) Ordination of vegetation samples by Gaussian species distributions. *Ecology* **55**, 1382–1390.
- Gauch, H.G., Whittaker, R.H. and Singer, S.B. (1981) A comparative study of nonmetric ordinations. *J. Ecol.* **69**, 135–152.
- Gifi, A. (1981) *Nonlinear Multivariate Analysis*. Department of Data Theory, University of Leiden, Leiden.
- Gittins, R. (1985) *Canonical Analysis. A Review With Applications in Ecology*. Springer-Verlag, Berlin.
- Goff, F.G. and Cottam, G. (1967) Gradient analysis: The use of species and synthetic indices. *Ecology* **48**, 793–806.
- Goodall, D.W. and Johnson, R.W. (1982) Non-linear ordination in several dimensions: a maximum likelihood approach. *Vegetatio* **48**, 197–208.
- Gourlay, A.R. and Watson, G.A. (1973) *Computational Methods for Matrix Eigen Problems*. John Wiley, New York.
- Greenacre, M.J. (1984) *Theory and Applications of Correspondence Analysis*. Academic Press, London.
- Heiser, W.J. (1981) *Unfolding Analysis of Proximity Data*. Thesis, University of Leiden, Leiden.
- Heiser, W.J. (1987) Joint ordination of species and sites: the unfolding technique. In: *Developments in Numerical Ecology* (Ed. by P. Legendre and L. Legendre), pp. 189–221. Springer-Verlag, Berlin.
- Hill, M.O. (1973) Reciprocal averaging: an eigenvector method of ordination. *J. Ecol.* **61**, 237–249.
- Hill, M.O. (1974) Correspondence analysis: a neglected multivariate method. *Appl. Statist.* **23**, 340–354.
- Hill, M.O. (1977) Use of simple discriminant functions to classify quantitative phytosociological data. In: *First International Symposium on Data Analysis and Informatics* (Ed. by E. Diday, L. Lebart, J.P. Pages and R. Tomassone), Vol. 1, pp. 181–199. INRIA, Chesnay.
- Hill, M.O. (1979) *DECORANA—A FORTRAN Program for Detrended Correspondence Analysis and Reciprocal Averaging*. Section of Ecology and Systematics, Cornell University, Ithaca, New York.
- Hill, M.O. and Gauch, H.G. (1980) Detrended correspondence analysis: an improved ordination technique. *Vegetatio* **42**, 47–58.

- Ihm, P. and Van Groenewoud, H. (1975) A multivariate ordering of vegetation data based on Gaussian type gradient response curves. *J. Ecol.* **63**, 767–777.
- Ihm, P. and Van Groenewoud, H. (1984) Correspondence analysis and Gaussian ordination. *COMPSTAT Lectures* **3**, 5–60.
- Imbrie, J. and Kipp, N.G. (1971) A new micropaleontological method for quantitative paleoclimatology: application to a late Pleistocene Caribbean core. In: *The Late Cenozoic Glacial Ages* (Ed. by K.K. Turekian), pp. 71–181. Yale University Press, New Haven, CT.
- Israëls, A.Z. (1984) Redundancy analysis for qualitative variables. *Psychometrika* **49**, 331–346.
- Iwatsubo, S. (1984) The analytical solutions of eigenvalue problem in the case of applying optimal scoring method to some types in data. In: *Data Analysis and Informatics 3* (Ed. by E. Diday), pp. 31–40. North-Holland, Amsterdam.
- Jolliffe, I.T. (1986) *Principal Component Analysis*. Springer-Verlag, Berlin.
- Jongman, R.H.G., Ter Braak, C.J.F. and Van Tongeren, O.F.R. (1987) *Data Analysis in Community and Landscape Ecology*. Pudoc, Wageningen.
- Kalkhoven, J. and Opdam, P. (1984) Classification and ordination of breeding bird data and landscape attributes. In: *Methodology in Landscape Ecological Research and Planning* (Ed. by J. Brandt and P. Agger), Vol. 3, pp. 15–26. Roskilde Universitetsforlag GeoRue, Theme 3, Roskilde.
- Kenkel, N.C. and Orlóci, L. (1986) Applying metric and nonmetric multidimensional scaling to ecological studies: some new results. *Ecology* **67**, 919–928.
- Kooijman, S.A.L.M. (1977) Species abundance with optimum relations to environmental factors. *Ann. Syst. Res.* **6**, 123–138.
- Kooijman, S.A.L.M. and Hengeveld, R. (1979) The description of a non-linear relationship between some carabid beetles and environmental factors. In: *Contemporary Quantitative Ecology and Related Econometrics* (Ed. by G.P. Patil and M.L. Rosenzweig), pp. 635–647. International Co-operative Publishing House, Fairland MD.
- Laurec, A., Chardy, P., de la Salle, P. and Rickaert, M. (1979) Use of dual structures in inertia analysis: ecological implications. In: *Multivariate Methods in Ecological Work* (Ed. by L. Orlóci, C.R. Rao and W.M. Stiteler), pp. 127–174. International Co-operative Publishing House, Fairland MD.
- McCullagh, P. and Nelder, J.A. (1983) *Generalized Linear Models*. Chapman and Hall, London.
- Macdonald, G.M. and Ritchie, J.C. (1986) Modern pollen spectra from the western interior of Canada and the interpretation of Late Quaternary vegetation development. *New Phytol.* **103**, 245–268.
- Meulman, J. and Heiser, W.J. (1984) Constrained multidimensional scaling: more directions than dimensions. In: *COMPSTAT 1984*, pp. 137–142. Physica-Verlag, Vienna.
- Minchin, P. (1987) An evaluation of the relative robustness of techniques for ecological ordination. *Vegetatio* **69**, 89–107.
- Montgomery, D.C. and Peck, E.A. (1982) *Introduction to Linear Regression Analysis*. John Wiley, New York.

- Nishisato, S. (1980) *Analysis of Categorical Data: Dual Scaling and Its Applications*. University of Toronto Press, Toronto.
- Oksanen, J. (1983) Ordination of boreal heath-like vegetation with principal component analysis, correspondence analysis and multidimensional scaling. *Vegetation* **52**, 181–189.
- Opdam, P.F.M., Kalkhoven, J.T.R. and Phillippona, J. (1984) *Verband tussen Broedvogelgemeenschappen en Begroeiing in een Landschap bij Amerongen*. Pudoc, Wageningen.
- Peet, R.K. (1978) Latitudinal variation in southern Rocky Mountain forests. *J. Biogeogr.* **5**, 275–289.
- Peet, R.K. and Loucks, O.L. (1977) A gradient analysis of southern Wisconsin forests. *Ecology* **58**, 485–499.
- Pickett, S.T.A. (1980) Non-equilibrium coexistence of plants. *Bull. Torrey bot. Club* **107**, 238–248.
- Pielou, E.C. (1984) *The Interpretation of Ecological Data*. John Wiley, New York.
- Prodon, R. and Lebreton, J.-D. (1981) Breeding avifauna of a Mediterranean succession: the holm oak and cork oak series in the eastern Pyrenees. 1. Analysis and modeling of the structure gradient. *Oikos* **37**, 21–38.
- Purata, S.E. (1986) Studies on secondary succession in Mexican tropical rain forest. *Acta Univ. Ups.* Comprehensive Summaries of Uppsala Dissertations from the Faculty of Science 19. Almqvist and Wiksell International, Stockholm.
- Rao, C.R. (1964) The use and interpretation of principal components analysis in applied research. *Sankhya A* **26**, 329–358.
- Robert, P. and Escoufier, Y. (1976) A unifying tool for linear multivariate statistical methods: the RV-coefficient. *Appl. Statist.* **25**, 257–265.
- Sládeček, V. (1973) System of water quality from the biological point of view. *Arch. Hydrobiol. Beiheft* **7**, 1–218.
- Smith, P.L. (1979) Splines as a useful and convenient statistical tool. *Am. Stat.* **33**, 57–62.
- Swaine, M.D. and Greig-Smith, P. (1980) An application of principal components analysis to vegetation change in permanent plots. *J. Ecol.* **68**, 33–41.
- Ter Braak, C.J.F. (1983) Principal components biplots and alpha and beta diversity. *Ecology* **64**, 454–462.
- Ter Braak, C.J.F. (1985) Correspondence analysis of incidence and abundance data: properties in terms of a unimodal response model. *Biometrics* **41**, 859–873.
- Ter Braak, C.J.F. (1986) Canonical correspondence analysis: a new eigenvector technique for multivariate direct gradient analysis. *Ecology* **67**, 1167–1179.
- Ter Braak, C.J.F. (1987a) The analysis of vegetation–environment relationships by canonical correspondence analysis. *Vegetatio* **69**, 69–77.
- Ter Braak, C.J.F. (1987b) *CANOCO-a FORTRAN Program for Canonical Community Ordination by [Partial] [Detrended] [Canonical] Correspondence Analysis, Principal Components Analysis and Redundancy Analysis (Version 2.1)*. Agriculture Mathematics Group, Wageningen.
- Ter Braak, C.J.F. (1988) Partial canonical correspondence analysis. In: *Classification Methods and Related Methods of Data Analysis* (Ed. by H.H. Bock), pp. 551–558. North-Holland, Amsterdam.

- Ter Braak, C.J.F. and Barendregt, L.G. (1986) Weighted averaging of species indicator values: its efficiency in environmental calibration. *Math. Biosci.* **78**, 57–72.
- Ter Braak, C.J.F. and Looman, C.W.N. (1986) Weighted averaging, logistic regression and the Gaussian response model. *Vegetatio* **65**, 3–11.
- Tilman, D. (1982) *Resource Competition and Community Structure*. Princeton University Press, Princeton.
- Tso, M.K.-S. (1981) Reduced-rank regression and canonical analysis. *J. Roy. Statist. Soc. B* **43**, 183–189.
- van der Aart, P.J.M. and Smeenk-Enserink, N. (1975) Correlations between distribution of hunting spiders (Lycosidae, Ctenidae) and environmental characteristics in a dune area. *Neth. J. Zool.* **25**, 1–45.
- van den Wollenberg, A.L. (1977) Redundancy analysis. An alternative for canonical correlation analysis. *Psychometrika* **42**, 207–219.
- van Dam, H., Suurmond, G. and Ter Braak, C.J.F. (1981) Impact of acidification on diatoms and chemistry of Dutch moorland pools. *Hydrobiologia* **83**, 425–459.
- Webb, T. III and Bryson, R.A. (1972) Late- and postglacial climatic change in the northern Midwest, USA: quantitative estimates derived from fossil spectra by multivariate statistical analysis. *Quat. Res.* **2**, 70–115.
- Whittaker, R.H. (1956) Vegetation of the Great Smoky Mountains. *Ecol. Monogr.* **26**, 1–80.
- Whittaker, R.H. (1967) Gradient analysis of vegetation. *Biol. Rev.* **49**, 207–264.
- Whittaker, R.H., Levin, S.A. and Ropt, R.B. (1973) Niche, habitat and ecotope. *Am. Natur.* **107**, 321–338.
- Wiens, J.A. and Rotenberry, J.T. (1981) Habitat associations and community structure of birds in shrubsteppe environments. *Ecol. Monogr.* **51**, 21–41.
- Williams, E.J. (1959) *Regression Analysis*. John Wiley, New York.
- Wold, H. (1982) Soft modeling. The basic design and some extensions. In: *Systems Under Indirect Observation. Causality–Structure–Prediction* (Ed. by K.G. Jöreskog and H. Wold), Vol. 2, pp. 1–54. North-Holland, Amsterdam.
- Zelinka, M. and Marvan, P. (1961) Zür Präzisierung der biologischen Klassifikation der Reinheit fließender Gewässer. *Arch. Hydrobiol.* **57**, 389–407.

## APPENDIX

A general iterative algorithm can be used to carry out the linear and weighted-averaging methods described in this review. The algorithm is essentially the one used in the computer program CANOCO (Ter Braak, 1987b). It operates on response variables, each recording the abundance or presence/absence of a species at various sites, and on two types of explanatory variables: environmental variables and covariables. By environmental variables we mean here explanatory variables of prime interest, in

contrast with covariables which are “concomitant” variables whose effect is to be removed. When all three types of variables are present, the algorithm describes how to obtain a *partial constrained ordination*. The other linear and WA techniques are all special cases, obtained by omitting various irrelevant steps.

Let  $Y = [y_{ki}]$  ( $k = 1, \dots, m; i = 1, \dots, n$ ) be a species-by-site matrix containing the observations of  $m$  species at  $n$  sites ( $y_{ki} \geq 0$ ) and let  $Z_1 = [z_{1ji}]$  ( $l = 0, \dots, p; i = 1, \dots, n$ ) and  $Z_2 = [z_{2ji}]$  ( $j = 1, \dots, q; i = 1, \dots, n$ ) be covariable-by-site and environmental variable-by-site matrices containing the observations of  $p$  covariables and  $q$  environmental variables at the same  $n$  sites, respectively. The first row of  $Z_1$ , with index  $l = 0$ , is a row of 1's which is included to account for the intercept in Eq. (4). Further, denote the species and site scores on the  $s$ th ordination by  $\mathbf{u} = [u_k]$  ( $k = 1, \dots, m$ ) and  $x = [x_i]$  ( $i = 1, \dots, n$ ), the canonical coefficients of the environmental variables by  $\mathbf{c} = [c_j]$  ( $j = 1, \dots, q$ ) and collect the site scores on the  $(s-1)$  previous ordination axes as rows of the matrix  $A$ . If detrending-by-polynomials is in force (Step A10), then the number of rows of  $A$ ,  $s_A$  say, is greater than  $s - 1$ . In the algorithm we use the assign statement “ $:=$ ”, for example  $a := b$  means “ $a$  is assigned the value  $b$ ”. If the left-hand side of the assignment is indexed by a subscript, it is assumed that the assignment is made for all permitted subscript values; the subscript  $k$  will refer to species ( $k = 1, \dots, m$ ), the subscript  $i$  to sites ( $i = 1, \dots, n$ ) and the subscript  $j$  to environmental variables ( $j = 1, \dots, q$ ).

**Preliminary Calculations**

*P1* Calculate species totals  $\{y_{k+}\}$ , site totals  $\{y_{+i}\}$  and the grand total  $y_{++}$ . If a linear method is required, set

$$r_k := 1, \quad w_i := 1, \quad w_i^* := \frac{1}{n} \tag{A.1}$$

and if a weighted averaging method is required, set

$$r_k := y_{k+}, \quad w_i := y_{+i}, \quad w_i^* := y_{+i}/y_{++} \tag{A.2}$$

*P2* Standardize the environmental variables to zero mean and unit variance. For environmental variable  $j$  calculate its mean  $\bar{z}$  and variance  $v$

$$\bar{z} := \sum_i w_i^* z_{2ji}, \quad v := \sum_i w_i^* (z_{2ji} - \bar{z})^2 \tag{A.3}$$

and set  $z_{2ji} := (z_{2ji} - \bar{z})/\sqrt{v}$ .

*P3* Calculate for each environmental variable  $j$  the residuals of the multiple regression of the environmental variables on the covariables, i.e.

$$\mathbf{c}_j^* := (\mathbf{Z}_1 \mathbf{W} \mathbf{Z}_1')^{-1} \mathbf{Z}_1 \mathbf{W} \mathbf{z}_{2j} \quad (\text{A.4})$$

$$\tilde{\mathbf{z}}_{2j} := \mathbf{z}_{2j} - \mathbf{Z}_1' \mathbf{c}_j^* \quad (\text{A.5})$$

where  $\mathbf{z}_{2j} = (z_{2j}, \dots, z_{2jm})'$ ,  $\mathbf{W} = \text{diag}(w_1, \dots, w_n)$  and  $\mathbf{c}_j^*$  is the  $(p+1)$ -vector of the coefficients of the regression of  $\mathbf{z}_{2j}$  on  $\mathbf{Z}_1$ . Now define  $\tilde{\mathbf{Z}}_2 = [\tilde{\mathbf{z}}_{2ji}]$  ( $j=1, \dots, q$ ,  $i=1, \dots, n$ ).

### Iteration Algorithm

*Step A0* Start with arbitrary, but unequal site scores  $\mathbf{x} = [x_i]$ . Set  $x_i^0 = x_i$

*Step A1* Derive new species scores from the site scores by

$$u_k := \sum_i y_{ki} x_i / r_k \quad (\text{A.6})$$

*Step A2* Derive new site scores  $\mathbf{x}^* = [x_i^*]$  from the species scores

$$x_i^* := \sum_k y_{ki} u_k / w_i \quad (\text{A.7})$$

*Step A3* Make  $\mathbf{x}^* = [x_i^*]$  uncorrelated with the covariables by calculating the residuals of the multiple regression of  $\mathbf{x}^*$  on  $\mathbf{Z}_1$ :

$$\mathbf{x}^* := \mathbf{x}^* - \mathbf{Z}_1' (\mathbf{Z}_1 \mathbf{W} \mathbf{Z}_1')^{-1} \mathbf{Z}_1 \mathbf{W} \mathbf{x}^* \quad (\text{A.8})$$

*Step A4* If  $q \leq s_A$ , set  $x_i := x_i^*$  and skip Step A5.

*Step A5* If  $q > s_A$ , calculate a multiple regression of  $\mathbf{x}^*$  on  $\tilde{\mathbf{Z}}_2$

$$\mathbf{c} := (\tilde{\mathbf{Z}}_2 \mathbf{W} \tilde{\mathbf{Z}}_2')^{-1} \tilde{\mathbf{Z}}_2 \mathbf{W} \mathbf{x}^* \quad (\text{A.9})$$

and take as new site scores the fitted values:

$$\mathbf{x} := \tilde{\mathbf{Z}}_2' \mathbf{c} \quad (\text{A.10})$$

*Step A6* If  $s > 1$ , make  $\mathbf{x} = [x_i]$  uncorrelated with previous axes by calculating the residuals of the multiple regression of  $\mathbf{x}$  on  $\mathbf{A}$ :

$$\mathbf{x} := \mathbf{x} - \mathbf{A}' (\mathbf{A} \mathbf{W} \mathbf{A}')^{-1} \mathbf{A} \mathbf{W} \mathbf{x} \quad (\text{A.11})$$

*Step A7* Standardized  $\mathbf{x} = [x_i]$  to zero mean and unit variance by

$$\begin{aligned} \tilde{\mathbf{x}} &:= \sum_i w_i^* x_i, \quad \sigma^2 := \sum_i w_i^* (x_i - \bar{x})^2 \\ x_i &:= (x_i - \bar{x})/\sigma \end{aligned} \tag{A.12}$$

*Step A8* Check convergence, i.e. if

$$\sum_i w_i^* (x_i^0 - x_i)^2 < 10^{-10} \tag{A.13}$$

goto Step A9, else set  $x_i^0 := -x_i$  and goto Step A1.

*Step A9* Set the eigenvalue  $\lambda$  equal to  $\sigma$  in (A.12) and add  $\mathbf{x} = [x_i]$  as a new row to the matrix  $A$ .

*Step A10* If detrending-by-polynomials is required, calculate polynomials of  $\mathbf{x}$  up to order 4 and first-order polynomials of  $\mathbf{x}$  with the previous ordination axes,

$$x_{2i} := x_i^2, \quad x_{3i} := x_i^3, \quad x_{4i} := x_i^4, \quad x_{(bi)} := -x_i a_{bi} \tag{A.14}$$

where  $a_{bi}$  are the site scores of a previous ordination axis ( $b = 1, \dots, s - 1$ ). Now perform for each of the  $(s + 2)$ -variables in (A.14) the Steps A3–A6 and add the resulting variables as new variables to the matrix  $A$ .

*Step A11* Set  $s := s + 1$  and goto Step A0 if required and if further ordination axes can be extracted, else stop.

At convergence, the algorithm gives the solution with the greatest real value of  $\lambda$  to the following transition formulae [where  $R = \text{diag}(r_1, \dots, r_m)$  and  $W = \text{diag}(w_1, \dots, w_n)$  and where the notation  $B^0$  is used to denote  $B'(BWB')^{-1}BW$ , the projection operator on the row space of a matrix  $B$  in the metric defined by the matrix  $[W]$ :

$$\mathbf{u} = R^{-1} Y \mathbf{x} \tag{A.15}$$

$$\mathbf{x}^* = (I - \tilde{Z}_1^0) W^{-1} Y' \mathbf{u} \tag{A.16}$$

$$\mathbf{c} = (\tilde{Z}_2 W \tilde{Z}_2')^{-1} \tilde{Z}_2 W \mathbf{x}^* \tag{A.17}$$

$$\lambda \mathbf{x} = (I - A^0) \tilde{Z}_2' \mathbf{c} \tag{A.18}$$

The tilde above  $Z_2$  is there as a reminder that the original environmental variables were replaced by residuals of a regression on  $Z_1$  in (A.5), i.e. in

terms of the original variables

$$\tilde{Z}'_2 = (I - Z_1^0)Z'_2 \quad (\text{A.19})$$

### Remarks

- (1) Note that  $u_k$  in the algorithm takes the place of  $b_k$  in [Section II](#).
- (2) Special cases of the algorithm are: constrained ordination:  $p=0$ ; partial ordination:  $q=0$ ; (unconstrained) ordination:  $p=0$ ,  $q=0$ ; linear calibration and weighted averaging:  $p=0$ ,  $q=1$ ; (partial) multiple regression:  $m=1$ . The corresponding transition formulae follow from [\(A.15\)–\(A.18\)](#) with the proviso that, if  $q=0$ ,  $Z_2$  in [\(A.19\)](#) must be replaced by the  $n \times n$  identity matrix and generalized matrix inverses are used. Note that, if  $p=0$ ,  $Z_1$  is a  $1 \times n$  matrix containing 1's;  $Z_1$  renders the centring of the species data in the linear methods in [Section II](#) redundant.
- (3) The standardization in P2 removes the arbitrariness in the units of measurement of the environmental variables, and makes the canonical coefficients comparable among each other, but does not influence the values of  $\lambda$ ,  $\mathbf{u}$  and  $\mathbf{x}$  to be obtained in the algorithm.
- (4) Step A6 simplifies to Step 3b of the main text if the rows of  $A$  are  $W$ -orthonormal. The steps A3–A6 form a single projection of  $\mathbf{x}^*$  on the column space of  $(I - A^0)\tilde{Z}'_2$  if and only if  $A$  defines a subspace of the row space of  $\tilde{Z}'_2$ . As each ordination axis defines such a subspace, this is trivially so without detrending. The method of detrending-by-polynomials as defined in Step A10, ensures that  $A$  defines also a subspace of  $\tilde{Z}'_2$  if detrending is in force. The transition formulae [\(A.15\)–\(A.18\)](#) define an eigenvalue equation of which all eigenvalues are real non-negative ([Ter Braak, 1987b](#)).
- (5) If a particular scaling of the biplot or the joint plot is wanted, the ordination axes may require linear rescaling. With linear methods one can choose between a Euclidean distance biplot and a covariance biplot, which focus on the approximate Euclidean distances between sites and correlations among species, respectively ([Ter Braak, 1983](#)). With weighted averaging methods it is customary to use the site scores  $\mathbf{x}^*$  [\(A.16\)](#) and the species scores  $\mathbf{u}$  [\(A.15\)](#) to prepare an ordination diagram after a linear rescaling so that the average within-site variance of the species scores is equal to 1 (cf. [Section IV.C](#)), as is done in DECORANA ([Hill, 1979](#)) and CANOCO ([Ter Braak, 1987b](#)).