
Multivariate Statistics

Jeff Powell

WESTERN SYDNEY
UNIVERSITY



Hawkesbury Institute
for the Environment

April 26, 2019

Contents

1	Multivariate statistics	2
1.1	Using multivariate statistics: when and why?	2
1.1.1	Caveats and considerations	2
1.2	Ordination: Exploring your data	3
1.2.1	Principal Components Analysis (PCA)	4
1.2.2	Correspondence analysis (CA)	8
1.2.3	Principal coordinates analysis (PCoA)	10
1.2.4	Non-metric multidimensional scaling (NMDS)	12
1.3	Two (or more)-table analysis: Unveiling drivers of data structure	14
1.3.1	Matrix correlations	15
1.3.2	Canonical analysis	15
1.3.3	More than two tables - variation partitioning	21
1.3.4	Using categorical variables in canonical analyses	29
1.3.5	'Experimental' frameworks: more working with factors	34
1.4	Functions used in this chapter	45
1.5	Exercises	46
1.5.1	Ordination	46
1.5.2	Analysis of Structure 1: two-table analysis	48
1.5.3	Analysis of Structure 2: variation partitioning	48
1.5.4	Analysis of Structure 3: 'experimental' systems	49

Chapter 1

Multivariate statistics

1.1 Using multivariate statistics: when and why?

In many studies, several types of data are collected on the same unit. Analysing these in combination may reveal patterns (resemblance between objects) that are indicative of structure in these data (how objects are organised along gradients), which may be more informative than looking at each individual variable. You may want to do some exploratory data analysis with these data to see whether certain individuals can be classified into groups based on their relative similarities. This would be useful, for example, for identifying soil types to be used as blocking factors in an experiment or groups of individuals that expressed similar traits in response to an experimental manipulation. You may also want to determine whether two or more variables provide essentially the same information, which will allow you to more efficiently focus your efforts on one of these. Or you may have specific hypotheses that you are interested in testing and are worried that, by performing these tests on multiple response variables, you may be inflating error rates associated with rejecting a correct null hypothesis (Type I error).

Further reading There are several resources for performing multivariate statistical analysis in **R** and a summary of these can be found at the Environmetrics CRAN Task View (<http://cran.r-project.org/web/views/Environmetrics.html>) under the “Ordination”, “Dissimilarity coefficients”, and “Cluster analysis” headings and at the Multivariate CRAN Task View (<https://cran.r-project.org/web/views/Multivariate.html>). A few tools are in the ‘stats’ package, which comes with your **R** distribution, but environmental scientists tend to mainly use the *vegan*, *ade4*, or *labdsv* packages. There is a lot of overlap between these packages in terms of the types of analyses that can be done, but the structure of the output is very different. Both *vegan* and *ade4* contain a function named `cca()`, so be careful if you have both packages loaded at the same time since the most recently loaded version of `cca()` will be called. You can specify the package in the call to the function – for example, `vegan::cca()`. Many other packages used by ecologists and evolutionary biologists rely on these packages. The Spatial CRAN Task View (<http://cran.r-project.org/web/views/Spatial.html>) will also be useful for more complex analysis of univariate and multivariate data with explicit spatial structure.

1.1.1 Caveats and considerations

You should think carefully about using multivariate statistics when the number of units being studied is on the same order as the number of variables being measured. A rule of thumb is that there should

be five independent units for every variable being measured. You can proceed with the analysis if you have less than this; in fact, you can even proceed with the analysis when you have fewer observations than you have variables. However, you should think carefully about replication when deciding whether it would be worthwhile to measure additional response variables, especially when testing hypotheses.

You should also think very carefully about what it is that you actually want to do. For instance, do you want to know whether there are patterns in the data, what are the important drivers generating patterns in the data, or which variables best represent patterns in the data? Once you have answered this question, it will become much easier to narrow down the possible approaches and identify the specific approach that is best suited to your data.

The choice of analysis can be determined by answering a series of questions:

- Am I looking for patterns in the data or do I want to make predictions about particular variables? For the former, use ordination methods (Section 1.2). For the latter, use multivariate ANOVA or GLM approaches (Section 1.3.5.3).
- For ordination, are my response variables continuous and approximately linear? If so, approaches based on principal components analysis (PCA; Section 1.2.1) are appropriate. If not, PCA may still work but this depends on gradient length, meaning the degree of overlap in the response variables among samples (e.g., if species composition, or betadiversity, does not change very much among samples; Section 1.2.2.1).
- For nonlinear and/or discrete responses, what might be an appropriate way to calculate the degree that the responses among samples differ? Raw count data may be used to calculate Chi-squared distances and perform correspondence analysis (CA; Section 1.2.2). Other data types (e.g., proportions and relative abundances) can be used to calculate various distance indices for principal coordinates analysis (PCoA; Section 1.2.3) or converted to ranks for nonmetric multidimensional scaling (NMDS; Section 1.2.4).
- Do I have additional data that may provide insight into the degree of resemblance among samples? If so, use constrained ordination (e.g., redundancy analysis [RDA], canonical correspondence analysis [CCA]) to partition variation in the response variables to these additional predictor variables prior to determining how much variation remains unexplained (Section 1.3.2).
- Do I want to estimate and compare the importance of specific predictor variables and their interactions? If so, use permutation-based approaches such as permutational MANOVA (PerMANOVA) to test specific hypotheses (Section 1.3.5.1).

Further reading Mike Palmer maintains a website (<http://ordination.okstate.edu>) that is helpful to get an overview of the many approaches that are available and the definitions of terms related to these approaches. The “GUSTA ME” project (GUide to STatistical Analysis in Microbial Ecology; <http://mb3is.megx.net/gustame/home>) includes accessible descriptions of multivariate statistical methods and worked examples (including flow charts) from the microbial ecology literature. “The R Book” by Michael Crawley includes a chapter on these methods that is also accessible to environmental scientists. For detailed discussion, see “Numerical Ecology” by Legendre and Legendre.

1.2 Ordination: Exploring your data

Unconstrained ordination approaches are very useful for simplifying multivariate data to visualise patterns. These should always be the first step in the analysis of multivariate data, even when you are interested in testing specific hypotheses regarding the potential indicators and drivers of data structure.

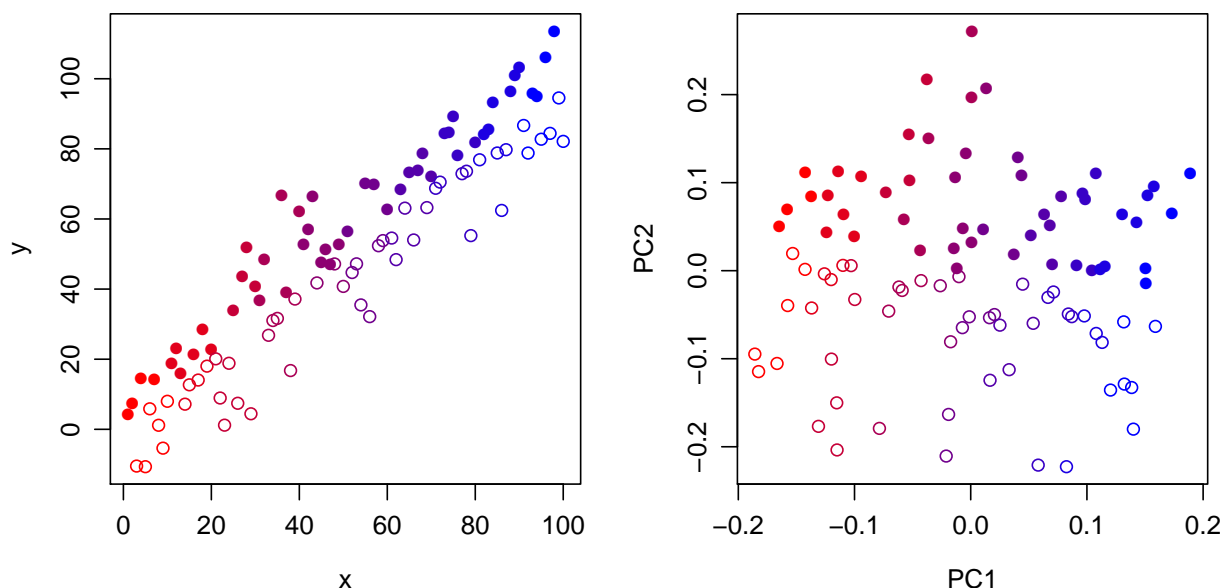


Figure 1.1: Example of two correlated variables (raw data shown on the left) and the resulting transformation following PCA (panel on right). The individual points are filled and coloured identically in both panels so that one can see the relationship between their position in each panel.

1.2.1 Principal Components Analysis (PCA)

PCA is ultimately the transformation of continuous multivariate data into new variables that maximise the amount of variance explained in the data, with each subsequent variable orthogonal (think perpendicular but in greater than two dimensions) to the previous variable and explaining decreasing amounts of variation. The new variables are linear combinations of the original variables.

Figure 1.1 shows an example of PCA, in which two correlated variables are transformed so that the first axis corresponds to the linear relationship between the two (represented by the shifting colours) and the second corresponds to the remaining variation (represented by the filled/unfilled circles). The individual points and the axes from the left panel are rotated in position to their locations on the right panel; their loadings reflect the degree of rotation.

Performing these is simple using `prcomp`, but we will use `rda` in `vegan` for the sake of consistency (particularly valuable when showing plotting methods). We'll use the 'varechem' data from the `vegan` package, which contains observations of soil variables associated with 24 sites grazed by reindeer.

```
## library(vegan) # loads the 'vegan' library

data(varechem) # read data included in the package into your workspace
str(varechem)  # observe the object structure

## 'data.frame': 24 obs. of 14 variables:
## $ N      : num  19.8 13.4 20.2 20.6 23.8 22.8 26.6 24.2 29.8 28.1 ...
## $ P      : num  42.1 39.1 67.7 60.8 54.5 40.9 36.7 31 73.5 40.5 ...
## $ K      : num  140 167 207 234 181 ...
## $ Ca     : num  519 357 973 834 777 ...
## $ Mg     : num  90 70.7 209.1 127.2 125.8 ...
## $ S      : num  32.3 35.2 58.1 40.7 39.5 40.8 33.8 27.1 42.5 60.2 ...
## $ Al     : num  39 88.1 138 15.4 24.2 ...
## $ Fe     : num  40.9 39 35.4 4.4 3 ...
```

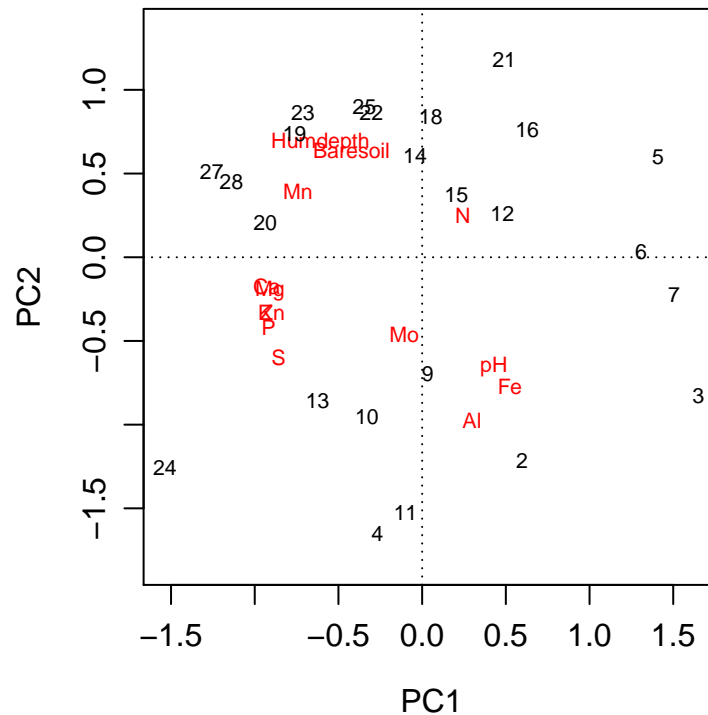


Figure 1.2: Biplot showing ordination of sites based on PCA of soil variables.

```
## $ Mn      : num  58.1 52.4 32.1 132 50.1 ...
## $ Zn      : num   4.5  5.4 16.8 10.7  6.6  9.1  7.4  5.2  9.3  9.1 ...
## $ Mo      : num   0.3  0.3  0.8  0.2  0.3  0.4  0.3  0.3  0.3  0.5 ...
## $ Baresoil: num  43.9 23.6 21.2 18.7  46 40.5 23 29.8 17.6 29.9 ...
## $ Humdepth: num   2.2  2.2  2  2.9  3  3.8  2.8  2  3  2.2 ...
## $ pH      : num   2.7  2.8  3  2.8  2.7  2.7  2.8  2.8  2.8  2.8 ...

# the variables differ in the scales of their variances
var(varechem$S) # variance in sulphur concentration
## [1] 136.1382

var(varechem$Ca) # variance in calcium concentration
## [1] 59332.17
```

The values associated with the different gradients are very different; for example, sulphur varies on a much smaller scale than calcium and their variances have to be standardised otherwise the large absolute variance associated with calcium will have much greater weight than that for sulphur. This is done with the `scale` argument (note that `scale=FALSE` is the default).

```
chem.pca <- rda(varechem, scale=TRUE) # use 'scale=TRUE' to standardise variances
plot(chem.pca) # plot the resulting object
```

The plot (Fig. 1.2) shows how the samples (indicated by number) are separated based on the first two principal components, the new variables resulting from the transformations of the original variables based on their ability to account for variation in the multivariate data. The samples are labelled according to the row names of the input dataframe or matrix. The black numbers indicate the loading of the individual samples ('sites') and the red labels indicate the loadings associated with the original variables ('species'; it is convention to refer to the columns in the multivariate response table as 'species',

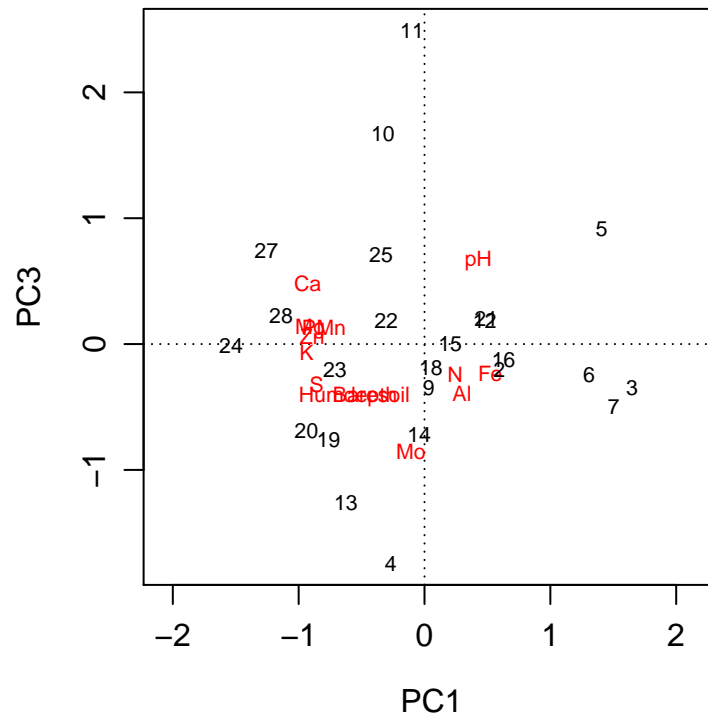


Figure 1.3: Biplot showing ordination (first and third principal components) of sites based on PCA of soil variables.

which becomes important when extracting these loadings). From this, we can see that the calcium and magnesium overlap entirely and, therefore, are positively correlated along both of the first two principal components. Because they are almost horizontal, their loadings associated with the first axis are much greater than that for the second axis. Sulfur and nitrogen are negatively correlated along both axes. Sample numbers positioned adjacent to variables along each axis indicates samples in which the values for those variable are high; from this we see that sample 24 is relatively high in calcium, magnesium, zinc, potassium, phosphorus, and sulfur while samples 3, 5, 6, and 7 are low in these elements.

```
varechem[c('5', '6', '7', '3', '24'), c('Ca', 'Mg', 'Zn', 'K', 'P', 'S')]
```

##	Ca	Mg	Zn	K	P	S
## 5	240.3	25.7	8.4	43.6	22.7	14.9
## 6	259.1	37.0	4.0	61.1	26.4	21.4
## 7	188.5	55.5	3.0	78.7	24.6	25.3
## 3	219.0	52.5	2.6	73.7	32.3	25.5
## 24	973.3	209.1	16.8	207.1	67.7	58.1

Using the `summary` function on the `rda` object provides a lot of information, including estimates of how much variation ('inertia') is associated with each PC axis. We can use this to evaluate how much of the variation in the entire dataset is represented by the PC axes shown in the PCA biplot.

```
summary(chem.pca, display=NULL) # use 'display=NULL' to suppress loadings tables
# (output not shown - run this yourself)
```

If we inspect the summary table we can see that these components explain 37 and 23 % of the variance in these data. The third component explains 12 % of variation, and we may be interested to see how samples and variables are distributed along this axis. To do this, we use the `choice` argument.

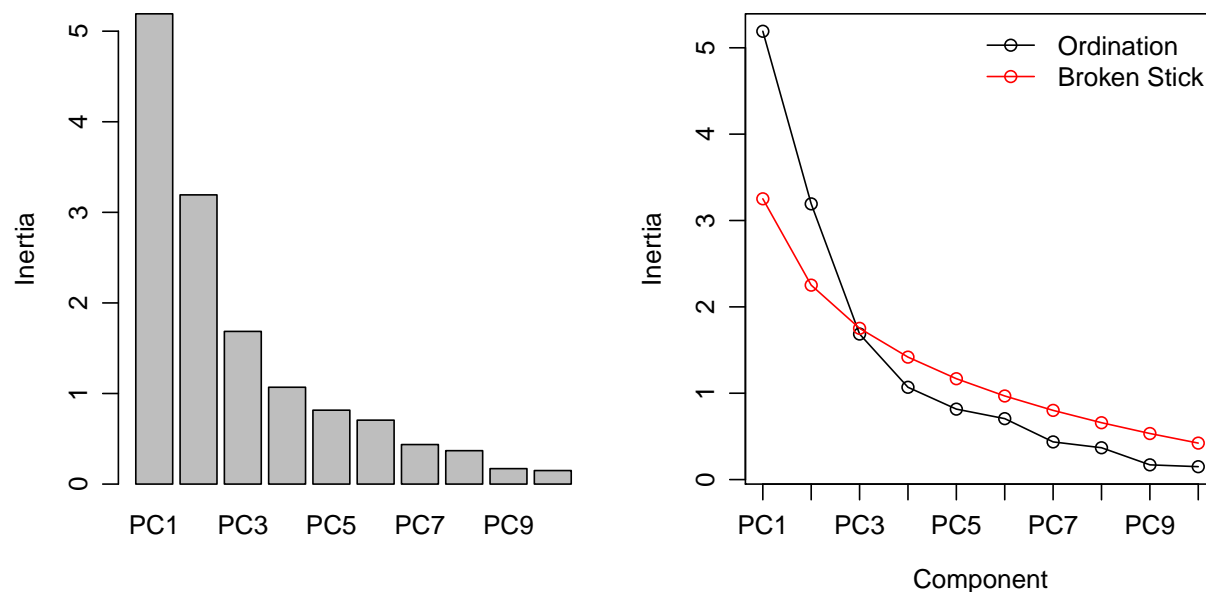


Figure 1.4: Screeplots for PCA of soil variables measured on plots grazed by raindeer. The panel on the left shows the variation (inertia) associated with each PC axis. The panel on the right compares observed variation associated with each axis in the ordination to that predicted based on a random expectation (broken-stick model).

The first principal component is still plotted on the x-axis of Fig. 1.3, but now the third component is plotted on the y-axis. We can see that calcium and magnesium are separated along the third axis (variation in calcium is explained by this axis, but not magnesium, which has loading close to zero). We also see that nitrogen and sulfur are positively correlated along the third axis (both are positioned on the negative side of this axis). The sample and variable scores for the first component are the same as before.

Most of the variation is accounted for in the first two principal components and > 90% is accounted for by the first six components. Any patterns observed associated with subsequent components will not contain much information. It is also possible to visualise the variation associated with each axis in a 'screeplot', as shown in Fig. 1.4 (left panel). In cases where most of the variation is explained in the first few axes, there should be a steep decline followed by a long, shallow tail as we view the screeplot from left to right.

```
# produce screeplot ('main=NULL' for a tidier plot)
screeplot(chem.pca, main=NULL)

# produce screeplot, including random expectation based on 'broken stick' model
screeplot(chem.pca, main=NULL, bstick=T, type='l')
```

The panel on the right of Fig. 1.4 shows the same data as a lineplot instead of barplot and comparing the result to that predicted based on a random, null expectation (for more information, see help page for the (bstick) function). The result suggests that only the first two PC axes should be used to interpret patterns since these are the only axes for which the observed inertia is higher than that of the null expectation.

Try this yourself Try to understand the importance of scaling variances by setting `scale = F` and looking at the result.

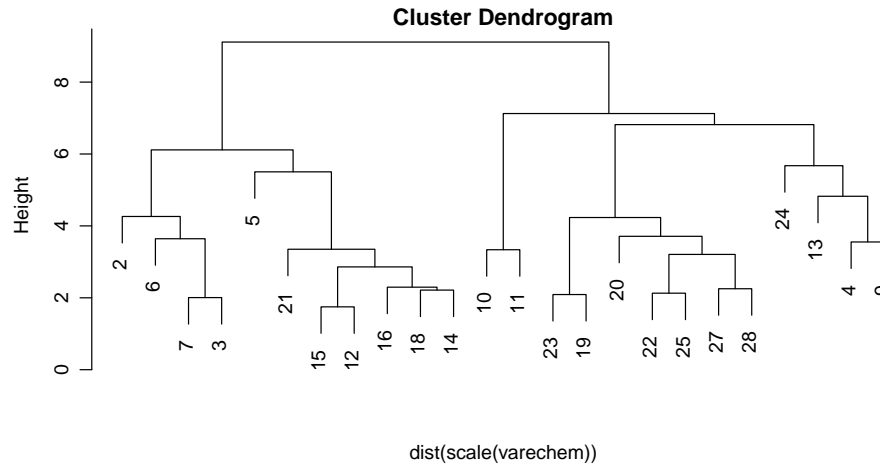


Figure 1.5: Dendrogram depicting the outcome of clustering based on variation in soil variables.

1.2.1.1 Cluster analysis

An alternative way to inspect the resemblance among samples is through plotting the result of hierarchical clustering, which clusters samples into groups that demonstrate increasingly more distant multivariate distances. The result of clustering can then be plotted as a dendrogram, in which nodes connecting short branches indicate groups of samples that exhibit similar responses. This has the advantage that patterns represent relationships across all variables, as opposed to only those with loadings associated with the first two or three PC axes. However, it is more difficult to see how each of the variables contributes to the patterns using this approach, compared with ordination using PCA.

Hierarchical clustering (using the `hclust` function) of the 24 grazed plots based on soil variables is shown in Fig. 1.5. Compare this result to the ordination in Fig. 1.2 and Fig. 1.3.

```
# input for 'hclust' is a distance matrix
# 'dist' produces a matrix of euclidean distances between each pair of samples
# 'scale' standardises the data before distance calculation
plot(hclust(dist(scale(varechem))))
```

1.2.2 Correspondence analysis (CA)

PCA uses euclidean distances, geometrical distances in multidimensional space, to estimate the divergence between samples. This is appropriate for continuous variables exhibiting normal distributions, like the above example, but not for data that resemble frequencies or counts. CA uses a different approach to estimate divergence among samples, based on calculation of χ^2 distances as for contingency tables. This property makes CA a better option than PCA for analysing tables of, for example, species counts in different environments.

CA is performed using `cca()` in `vegan` by providing only the response matrix as an argument to the function. We will use the plant community data collected from the same locations as the soil data analysed above; these data are stored in `varespec` and consist of observations relating to the percent cover associated with 44 plant species.

```
# read data included in the package into your workspace
data(varespec)

# performs CA on single input matrix
spec.ca <- cca(varespec)
```

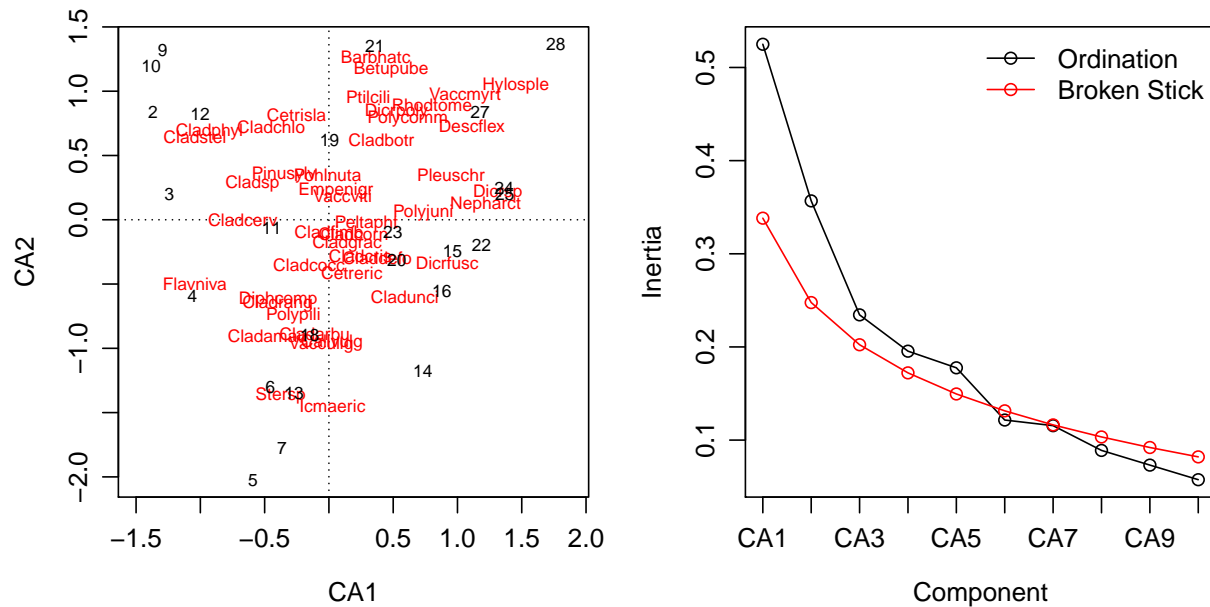


Figure 1.6: Biplot (left) and screeplot (right) showing ordination of sites based on plant communities following CA.

```
# plot resulting object and screeplot
plot(spec.ca)
screeplot(spec.ca, main=NULL, bstick=T, type='l')
```

The plots in Fig. 1.6 are interpreted in the same way as the PCA biplot and screeplot: the numbers represent the scores of the rows (samples) on the first two correspondence axes, the labels in red represent the loadings associated with the variables (species). The summary output (not shown here) is obtained and interpreted in the same way as the output from PCA.

There was no scaling of the variables in our example but this may be desirable if, for example, the species vary greatly in total abundance along the whole gradient. The `decostand` function will transform the input matrix into a standardised form with the desired characteristics. There are several methods available in `decostand` but the most commonly used are `standardize` (analogous to `scale=T`, but results in an error when using `cca` due to negative values), `normalize` (results in row sums of squares equal to one), `total` (returns proportional relative abundances within each sample), and `hellinger` (as for `total` but then returns the square root of the result).

In Fig. 1.7, the plot on the left shows the result of CA on the untransformed species abundances. The figure on the right shows the result when the data have been normalized (using `cca(decostand(varespec, method='normalize'))`). The spread of the site scores is slightly less triangular on the right than on the left, suggesting that the transformed data are not as distorted and the analysis is more likely to result in meaningful patterns.

1.2.2.1 PCA or CA?

Multivariate count or frequency data may be analysed using either PCA or CA. PCA (and RDA, see Section 1.3.2) may be used for variables that exhibit a roughly normal distribution and with few absences (zeros). These data may be referred to as exhibiting low turnover or a short gradient. CA (and CCA, see Section 1.3.2) is generally used for longer gradients, for example, when variables are recorded as present in some samples but not in others. One can determine which is the case by estimating gradient

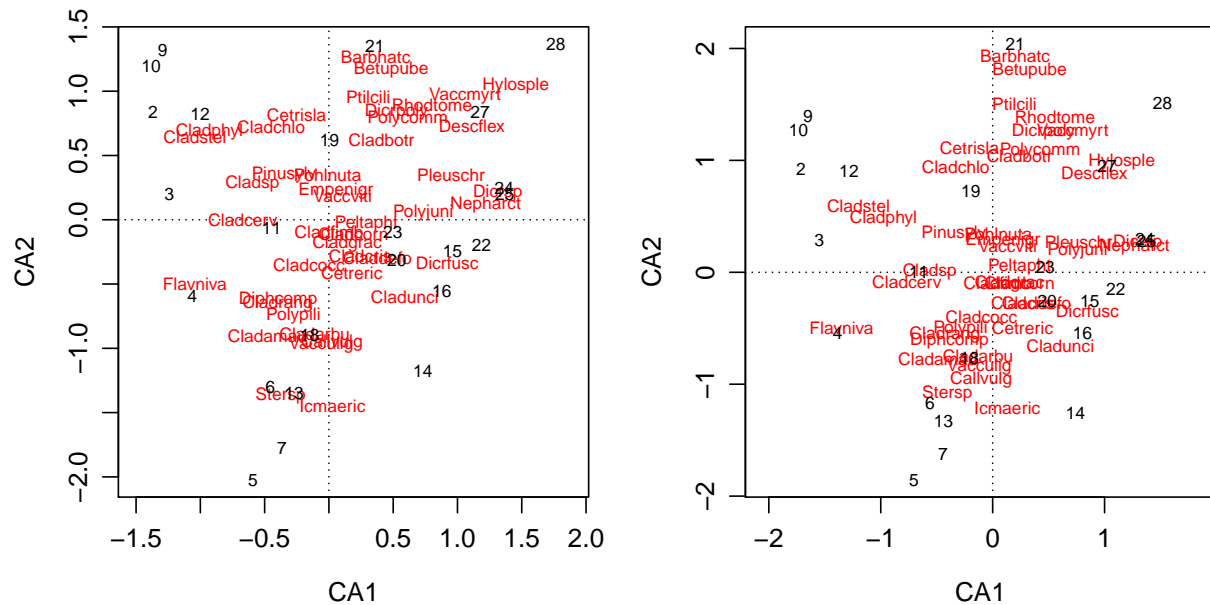


Figure 1.7: Correspondence analysis of community matrix based on raw data (left) and normalized data (right).

length using the `decorana()` function.

```
spec.dca <- decorana(varespec)
summary(spec.dca, display='none') # note different argument for 'display'

##
## Call:
## decorana(veg = varespec)
##
## Detrended correspondence analysis with 26 segments.
## Rescaling of axes with 4 iterations.
##
##          DCA1   DCA2   DCA3   DCA4
## Eigenvalues  0.5235 0.3253 0.20010 0.19176
## Decorana values 0.5249 0.1572 0.09669 0.06075
## Axis lengths  2.8161 2.2054 1.54650 1.64864
```

The gradient length along the first axis is 2.8161. A rule of thumb is that, since this value is less than 3, PCA would be an appropriate method (use CA when this value is greater than 3). However, it is not much smaller than 3; we could also try transforming the data using an alternative distance index prior to analysing them (next section).

1.2.3 Principal coordinates analysis (PCoA)

Just as euclidean distances are not appropriate for estimating the divergence between samples when analysing frequency data, some data types require an alternative to χ^2 distances. This may be the case for data in which there is a high frequency of absences (zeros) or for which binary (presence/absence) variables were recorded. PCoA allows for the analysis of these data using distance matrices estimated from various multidimensional scaling indices. Many common metrics for estimating community dissimilarity are available in the `vegdist()` function (see the help page for these indices and their formulae). The distance matrix is used as the input for `wcmdscale()` in `vegan`, which performs the

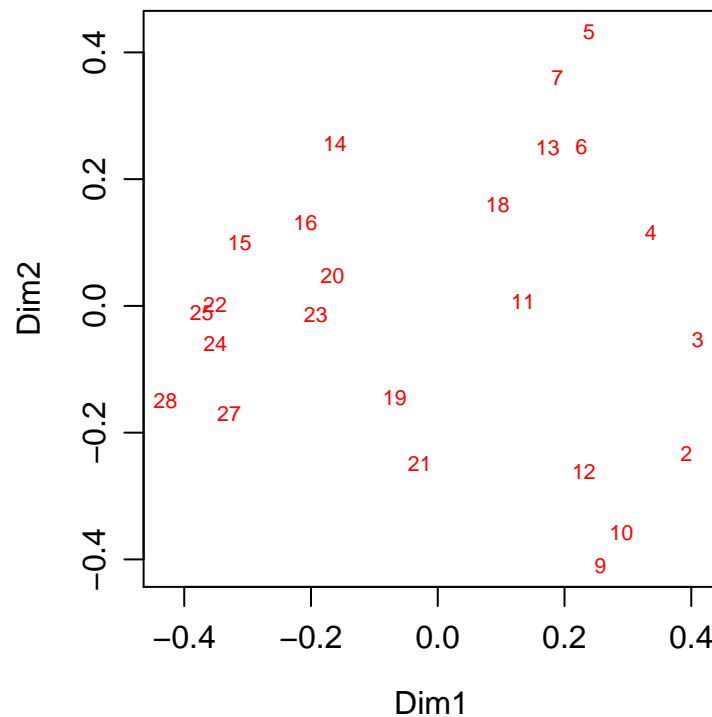


Figure 1.8: Principal coordinates analysis of the vegetation community matrix, using Bray-Curtis dissimilarities.

PCoA.

```
# use the default index ('bray') and 'eig=TRUE' to save the calculated eigenvalues
spec.pco <- wcmdscale(vegdist(varespec), eig=TRUE)
plot(spec.pco)
```

Notice in Fig. 1.8 that the ordinated data exhibit a roughly even spread along both the x- and y-axes, suggesting that we could interpret the patterns in these data in a meaningful way (although they do look a bit like the outline of Australia!). Also notice that there are no labels indicating the species loadings; this is because this information is not retained in the distance matrix that was provided as input. There is no built-in `summary()` function for returning interpretable output from PCoA analysis, but we can look at the eigenvalues to see how much variation is accounted for by the different principal coordinate axes (these values are available since we specified `eig=TRUE`).

```
# Return the eigenvalues
spec.pco$eig
## [1] 1.7552165397 1.1334455380 0.4429018480 0.3698054310 0.2453531540
## [6] 0.1960920773 0.1751130911 0.1284466728 0.0971594360 0.0759600747
## [11] 0.0637177905 0.0583225124 0.0394933793 0.0172699235 0.0051011077
## [16] -0.0004131222 -0.0064653552 -0.0133147491 -0.0253943546 -0.0375104890
## [21] -0.0480068852 -0.0537145779 -0.0741390257

# Return the proportion of variance explained (non-negative eigenvalues)
eigens <- spec.pco$eig[spec.pco$eig >= 0]
# proportion explained by each axis
eigens / sum(eigens)
## [1] 0.365411388 0.235967413 0.092205933 0.076988288 0.051079075
```

```
## [6] 0.040823612 0.036456082 0.026740790 0.020227228 0.015813819
## [11] 0.013265147 0.012141926 0.008221966 0.003595355 0.001061979

# cumulative proportion explained
cumsum(eigens / sum(eigens))

## [1] 0.3654114 0.6013788 0.6935847 0.7705730 0.8216521 0.8624757 0.8989318
## [8] 0.9256726 0.9458998 0.9617136 0.9749788 0.9871207 0.9953427 0.9989380
## [15] 1.0000000
```

The first two values (1.75, 1.13) are much larger than the rest, suggesting that most of the variation is accounted for by these two axes. There are also negative eigenvalues; to estimate the relative proportions of variance explained, we only use the nonnegative eigenvalues (indexing can be used to select these, as shown above).

How does one choose one of the available indices for calculating dissimilarity? One way to do so is to pick the index that provides the best rank-order similarity (i.e., nonparametric correlation based on two rank-transformed variables) to the gradient under study, using `rankindex()`. This gradient might be represented in an environmental matrix (here we will use the soil chemistry data from above) or some other indicator of a putative gradient (e.g., sampling coordinates).

```
# For rankindex(), the first argument is the gradient,
# the second is the community matrix. The output is a named vector of rank-
# -order similarities, each representing a dissimilarity index.

# First calculate rank-order similarity on the raw community matrix
rankindex(scale(varechem), varespec)

##          euc          man          gow          bra          kul
## 0.2396330 0.2735087 0.2288358 0.2837910 0.2839834

# Then calculate rank-order similarity after applying the
# hellinger transformation (see ?decostand)
rankindex(scale(varechem), decostand(varespec, method='hellinger'))

##          euc          man          gow          bra          kul
## 0.2842265 0.2562909 0.2570694 0.3183687 0.3168585
```

This function can also be used to determine which standardisation approach should be applied to the data. Comparing the values returned by the two calls to `rankindex` immediately above, we see that the Bray-Curtis index ('bra') provides a higher rank-order similarity than the other tested indices when the hellinger transformation is applied to the community matrix.

Try this yourself Try using different dissimilarity indices to see the effect this has on the PCoA result.

1.2.4 Non-metric multidimensional scaling (NMDS)

NMDS maximises the differences between samples on few dimensions, which can be particularly useful for visual representation of the dissimilarities between samples. It is “nonmetric” because the data undergo rank-order transformation and their positions are moved during the procedure to minimise stress. The `metaMDS` function in `vegan` performs NMDS on either a table of community data or a distance matrix calculated from this table, as above. Since the purpose of this approach is to visualise the spread of the data in reduced dimensions, it makes sense to plot the data (Fig. 1.9).

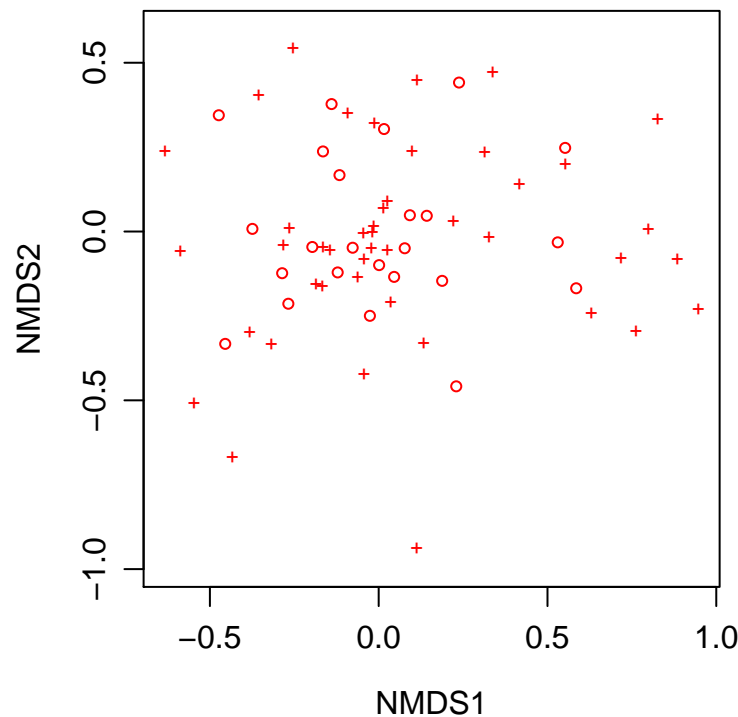


Figure 1.9: NMDS ordination of sites based on variation in plant communities

```
spec.nmds <- metaMDS(varespec, trymax=40) # (output not shown)
plot(spec.nmds)
```

The procedure is iterative, and technical data are output to the screen as it runs. *** Solution reached should appear at the end of this output (note that the output above is abbreviated). If not, redo the analysis and increase the value associated with the `trymax` argument (the default is 20). The stress in the final run is 0.183, which is acceptable; stress values greater than 0.3 are unsatisfactory suggesting that the dissimilarities are not effectively captured by these two dimensions. Increase the number of dimensions using the `k` argument (the default is 2).

1.2.4.1 Manipulating graphics from *vegan* objects

Figure 1.9 shows the loadings for the sites (circles) and species (crosses) from the above NMDS analysis. This is a good point to show how graphics can be manipulated as no labels are included in the plot. While the plotting functions in *vegan* are useful to visualise output from *rda* and *cca* and for interpreting this output, they are very basic and the results are not necessarily of publication quality. However, we can use our output and *R*'s default plotting functions to modify the graphics as we did in chapter ?? . The information needed to generate a plot for publication can be found using the `scores` function, which returns 'site' and 'species' loadings based on which is specified using the `display` argument.

```
# Set up plot window
plot(scores(spec.nmds, display='species'), type='n')

# Sites loadings
text(scores(spec.nmds, display='sites'),
      labels=rownames(scores(spec.nmds, display='sites')),
      cex=0.8, col="blue3")
```

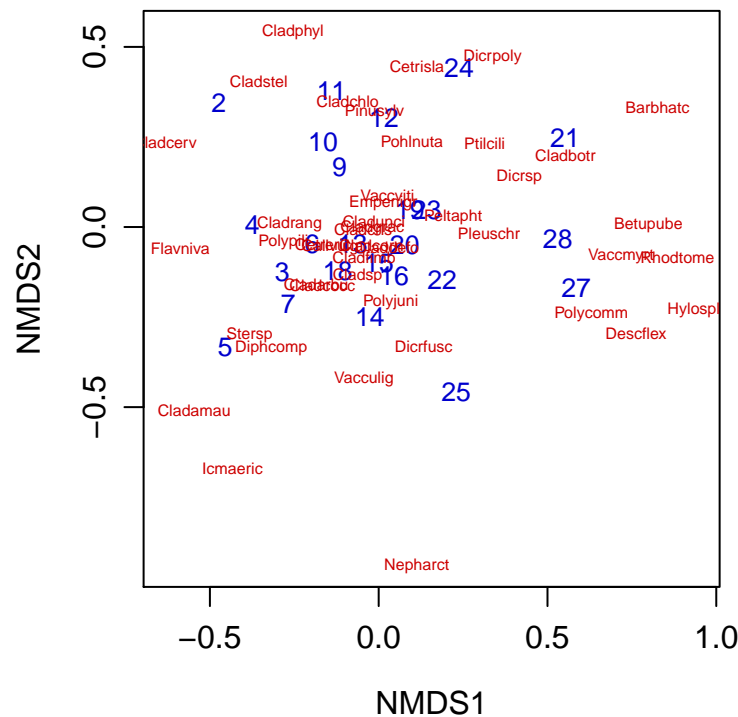


Figure 1.10: NMDS ordination of sites based on variation in plant communities.

```
# Species loadings
text(scores(spec.nmds, display='species'),
      labels=rownames(scores(spec.nmds, display='species')),
      cex=0.5, col="red3")
```

This was a multi-step process. First we set up the plot window, including the figure axes, but suppressed the plotting of the data points using `type='n'`. We used the `species` loadings to set up the plot window in this case as these expand further along both axes than the `sites` loadings; had we used the `sites` loadings, some of the labels belonging to species would fall outside the plot area. Then we used `text()` to plot the row names for the table of sites loadings in place of the data points. The second call to `text()` plots the row names for the table of species loadings; we used the `cex` argument in each line to reduce the text size.

Try this yourself Using the examples from PCA, CA, and PCoA, plot the results and manipulate. Note that for PCoA, no species loadings are provided as the analysis is performed on a distance matrix, not species abundances.

1.3 Two (or more)-table analysis: Unveiling drivers of data structure

Up to this point, we have used ordination to look for structure in multivariate data. Once that structure is observed, you may have ideas about what factors are driving that structure and would like to test the hypotheses that these factors are actually important. To do this is to apply a constraint to the data,

partitioning variation in the data to specific factors, or linear combinations of these specific factors (as above), and then comparing this partitioned variation to any remaining variation in the data.

1.3.1 Matrix correlations

Given two tables, each containing multivariate data about two aspects of the system under study (for example, species abundances and environmental characteristics), it may be of value to determine whether the responses across the two tables are related. In ecology, it is common to have collected data on species abundances in a variety of environments and to try to explain variation in those species' distributions due to characteristics of their environments. Similarly, physiologists and evolutionary biologists often take multiple types of measurements on a number of organisms and may wish to determine whether aspects of these organisms or their environment are important predictors of their responses.

Correlations are estimated on dissimilarity matrices, not on the raw data, so it is necessary to transform each table of data using `vegdist()` as for PCoA above.

```
mantel(vegdist(varespec, 'bray'), vegdist(scale(varechem), 'euclidian'))
```

```
##
## Mantel statistic based on Pearson's product-moment correlation
##
## Call:
## mantel(xdis = vegdist(varespec, "bray"), ydis = vegdist(scale(varechem),      "euclidian"))
##
## Mantel statistic r: 0.3047
##      Significance: 0.001
##
## Upper quantiles of permutations (null model):
##   90%   95%  97.5%   99%
## 0.112 0.144 0.166 0.202
## Permutation: free
## Number of permutations: 999
```

Note that we used different indices for the two tables. Bray-Curtis dissimilarities were estimated for the species-sample table, while euclidean distances were estimated for the table containing soil variables (after standardising these data using the `scale()` function). The function provides a correlation coefficient and an outcome of the hypothesis test that the coefficient does not differ from zero, by permutation. Here we see that the two matrices are correlated, which may mean that the variables in one table are important drivers of the values in the other or that both types of variables share a common driver or set of drivers.

1.3.2 Canonical analysis

By combining ordination techniques and multiple linear regression, we can use regression approaches to explain variation in one multivariate data table by variation in another multivariate data table from observations of the same objects. Prior to the availability of suitable software programs, principal components were extracted from the ordination of the explanatory matrix and then related to the ordination of the response matrix; this approach is named indirect gradient analysis. This approach has gradually been replaced by canonical analysis, or direct gradient analysis, in which the explanatory matrix is directly involved in the ordination of the response matrix.

1.3.2.1 Redundancy analysis (RDA)

RDA is an extension of PCA and, as such, is appropriate for estimating the importance of constraining variables along short gradients that display low species turnover (recall Section 1.2.2.1). Calculations are performed using the `rda()` function in `vegan`.

```
vare.rda <- rda(varespec, scale(varechem), scale=T)
vare.rda

## Call: rda(X = varespec, Y = scale(varechem), scale = T)
##
##              Inertia Proportion Rank
## Total          44.0000      1.0000
## Constrained    28.5273      0.6483   14
## Unconstrained  15.4727      0.3517    9
## Inertia is correlations
##
## Eigenvalues for constrained axes:
##  RDA1  RDA2  RDA3  RDA4  RDA5  RDA6  RDA7  RDA8  RDA9  RDA10  RDA11  RDA12
## 5.548 4.529 3.566 2.946 2.369 2.240 1.831 1.373 1.140 1.027 0.712 0.553
## RDA13 RDA14
## 0.389 0.303
##
## Eigenvalues for unconstrained axes:
##  PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8  PC9
## 4.965 2.582 2.059 1.740 1.446 0.930 0.743 0.548 0.458

anova(vare.rda)

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(X = varespec, Y = scale(varechem), scale = T)
##      Df Variance      F Pr(>F)
## Model   14   28.527 1.1852 0.159
## Residual   9   15.473

plot(vare.rda)
```

From the printed object, we see that 65 % of inertia is accounted for by the constraining variables. The plot shows the loadings associated with the species in the samples, and the variables in the explanatory table over the first two constrained axes (explaining 23 % of the total variation; you can calculate this yourself or find it stored in the summary object under `summary(vare.rda)[['cont']][['importance']]`). `vegan` includes an anova-like function that tests the significance of the constraints using permutation. Here we see that even though the variance explained is high, the model containing the constraint does not provide a significantly better fit. This may be due to the penalty associated with the large number of variables in the explanatory matrix; simply including all of the data without providing any thought as to which data may be useful is unlikely to generate an interpretable response.

There are a few options for selecting only a subset of variables that appear to be important predictors of changes in the community matrix. One approach is to use `envfit` to identify those variables whose vectors are significantly correlated with the site loadings.

```
vare.rda <- rda(varespec, scale(varechem), scale=T)
envfit(vare.rda, scale(varechem))

##
```

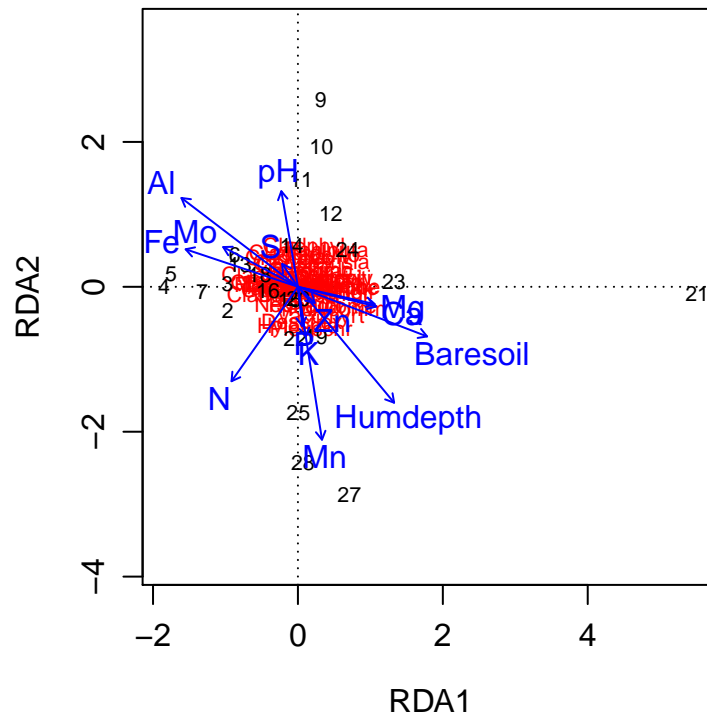


Figure 1.11: RDA of plant communities after constraining variation by soil variables.

```
## ***VECTORS
##
##          RDA1      RDA2      r2 Pr(>r)
## N          -0.46296 -0.88638 0.2575 0.042 *
## P           0.05485 -0.99849 0.0320 0.744
## K           0.08770 -0.99615 0.0517 0.592
## Ca          0.95338 -0.30177 0.0980 0.340
## Mg          0.96319 -0.26883 0.0929 0.331
## S          -0.43111  0.90230 0.0142 0.865
## Al          -0.68922  0.72455 0.3562 0.015 *
## Fe          -0.91722  0.39839 0.2135 0.085 .
## Mn          0.08864 -0.99606 0.4900 0.004 **
## Zn          0.59628 -0.80278 0.0094 0.898
## Mo          -0.81301  0.58225 0.1133 0.272
## Baresoil    0.89251 -0.45103 0.2942 0.020 *
## Humdepth    0.50530 -0.86295 0.4058 0.008 **
## pH          -0.10028  0.99496 0.1921 0.098 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Permutation: free
## Number of permutations: 999

varechem_subs <- varechem[,c('N','Al','Fe','Mn','Baresoil','Humdepth')]
vare.rda.envfit <- rda(varespec,
                      scale(varechem_subs),
                      scale=T)
anova(vare.rda.envfit)
```

```
## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(X = varespec, Y = scale(varechem_subs), scale = T)
##           Df Variance      F Pr(>F)
## Model      6   14.139 1.3416 0.032 *
## Residual  17   29.861
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

To test the significance of the relationship between the community and environmental matrices, we only included those variables whose vectors were significantly (or marginally-significantly) correlated with the sites loadings on the first two axes. We can see that this relationship is now significant, where it wasn't when including all environmental variables.

Try this yourself Note that the code provided here only identifies variables whose vectors are correlated with the loadings on the first two axes. This is appropriate in cases where these axes account for most of the variation in the community matrix. Use the `choices` argument in `envfit` to increase the number axes over which environmental vectors are correlated. See how this affects the significance of the relationship between the community and environmental matrices.

If the goal is to identify which variables are generally important, it would be better to use forward and/or reverse selection to select a subset of significant variables. We can do this using the `ordistep()` function in `vegan` after fitting two models: one containing all predictors and the other containing no predictors. The function prints output to the screen after each iteration and stops once adding/removing a variable to/from the model no longer significantly improves the model fit. Here we select potential drivers of plant community composition from a standardised matrix of environmental variables, using both forward and backward selection.

```
# select particular variables to proceed with
# here we use both forward and backward selection

# generate a new dataframe containing scaled predictor variables
# since difficult to do this in the formula
varechem.sclsd <- decostand(varechem, method='standardize')

# have to use the formula interface so generate a new 'rda' object
# including all predictors (use '.' after the '~')
vare.rda <- rda(varespec ~ ., data=varechem.sclsd, scale=T)

# set up the null case with no predictors (be sure to include the
# 'data' argument, even though no predictors)
vare.pca <- rda(varespec ~ 1, data=varechem.sclsd, scale=T)

# select variables in each predictor table (output not shown)
step.env <- ordistep(vare.pca, scope=formula(vare.rda))
```

We can then look at the result for when including only those variables that best explain variation between communities. We can also look at the significance level associated with each variable included in this model.

```
# show the object summary
step.env
```

```
## Call: rda(formula = varespec ~ Baresoil + Mn, data =
## varechem.scl, scale = T)
##
##              Inertia Proportion Rank
## Total          44.0000      1.0000
## Constrained    6.5676      0.1493    2
## Unconstrained 37.4324      0.8507   21
## Inertia is correlations
##
## Eigenvalues for constrained axes:
##  RDA1  RDA2
## 3.352 3.215
##
## Eigenvalues for unconstrained axes:
##  PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8
## 7.817 4.235 3.140 3.092 2.774 2.369 2.231 2.124
## (Showed only 8 of all 21 unconstrained eigenvalues)
# evaluate the statistical significance of the constraint
anova(step.env)

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(formula = varespec ~ Baresoil + Mn, data = varechem.scl, scale = T)
##           Df Variance      F Pr(>F)
## Model      2    6.568 1.8423 0.005 **
## Residual 21   37.432
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# presents results in an ANOVA-like table with the retained predictors
step.env$anova

##           Df      AIC      F Pr(>F)
## + Baresoil  1 91.931 1.7805 0.02 *
## + Mn        1 91.919 1.8363 0.03 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Try this yourself Use the `direction` argument to perform either 'forward' or 'backward' selection of predictor variables to include. Does this result in different variables being selected?

Variations to RDA include transformation of the response matrix prior to analysis (tb-RDA) using `decostand` (as in Section 1.2.2) or generating a response matrix from principle coordinates after PCoA (distance based [db]-RDA or constrained analysis of principal coordinates [CAP]; as in Section 1.2.3). We'll deal with these in subsequent sections.

1.3.2.2 Canonical correspondence analysis (CCA)

CCA is an extension of CA, but serves the same purpose as RDA in relation to PCA. CCA is performed using `cca` in `vegan`. Its use in `vegan` is very similar to `rda()`, so we will not deal with it further here. It may be preferred over RDA for long gradients, over which species distributions are limited to a subset

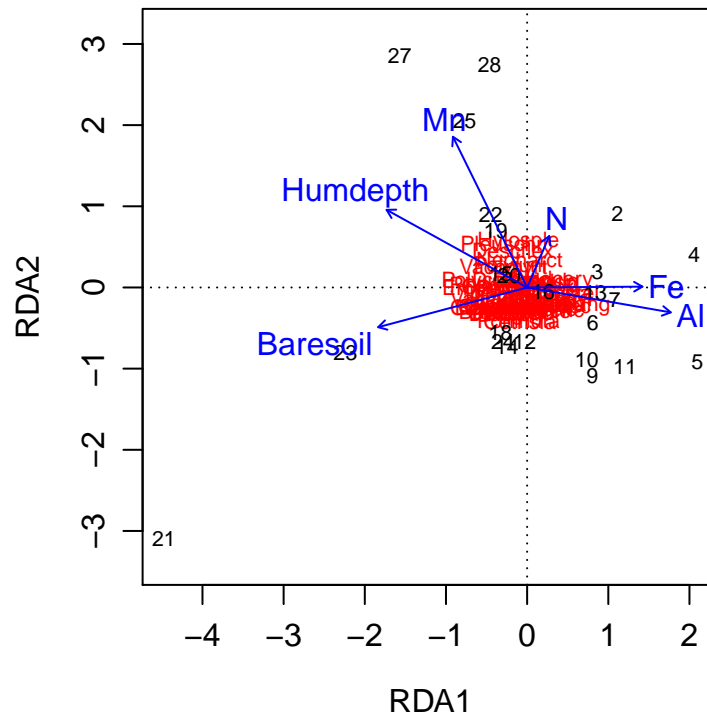


Figure 1.12: RDA ordination of sites based on variation in plant communities, constrained by significant soil properties

of the gradient and do not overlap for many species (recall Section 1.2.2.1), but db-RDA or NMDS-based approaches can also be used.

1.3.2.3 More manipulating graphics from *vegan* objects

The plots that are produced from using `plot` on *vegan* objects are useful for visualising patterns but not so nice for publication. For example, look at the result for the analysis performed in Section 1.3.2.1, shown in Fig. 1.12.

We can make much nicer plots (Fig. 1.13) showing the loadings associated with species and soil properties by extracting relevant information from the resulting object.

```
## get the relevant information from the RDA object

# 'sites' and 'species' loadings can be retrieved using 'scores' function
sites <- scores(vare.rda.envfit, display='sites')
spp <- scores(vare.rda.envfit, display='species')

# proportion of variation explained by each constrained axis (for plot axes)
# can't just use sum(eig) for denominator, this vector only includes constraint
eig <- vare.rda.envfit[['CCA']][['eig']]
eig / vare.rda.envfit[['tot.chi']]

##      RDA1      RDA2      RDA3      RDA4      RDA5      RDA6
## 0.10188359 0.08724602 0.06177211 0.03575530 0.02230310 0.01238436

## plot one figure on top of another
```

```

par(mfrow=c(2, 1), mar=c(5, 5, 1, 1))
## plot showing loadings for species (response variables)
plot(sites, pch=16, col='grey', xlim=c(-5, 5), ylim=c(-4, 4), main='Species',
      xlab='Axis 1 (10.2 %)', ylab='Axis 2 (8.7 %)')
# add lines to separate the plot regions
abline(v=0, h=0, lty='dashed')
# use a multiplier (here, 7) to spread the points away from the origin for readability
# also, only print the first four letters to reduce overlap
text(spp[, 1] * 7, spp[, 2] * 7, substr(rownames(spp), 1, 4), col='blue', cex=0.6)

## plot showing loadings for soil properties (constraining variables)
plot(sites, pch=16, col='grey', xlim=c(-5, 5), ylim=c(-4, 4), main='Constraints',
      xlab='Axis 1 (10.2 %)', ylab='Axis 2 (8.7 %)')
# add lines to separate the plot regions
abline(v=0, h=0, lty='dashed')
# add text and arrows linked to variable labels
text(vare.rda.envfit, display='bp', col='blue', cex=0.75)

```

Try this yourself In the above example, adjust the par settings to modify the plot to your liking.

1.3.3 More than two tables - variation partitioning

Canonical approaches calculate the correspondence between community composition and environmental properties, and allow for the estimation of each variable's explanatory power as a predictor of community shifts. These environmental properties may belong to different categories (e.g., chemical, physical, climatic) and/or may be measured along known spatial gradients. In these cases, we may want to partition variation in community composition to groups of variables in order to gain a general sense of how important each group of variables is for driving community shifts. We can do this using the `varpart` function in the `vegan` package.

In this example, we partition variation in the `varespec` plant data to two categories of environmental data: soil elemental chemistry and soil exposure (prevalence of bare soil and depth of the humus layer). We perform variation partitioning by including both types of predictor matrices as separate arguments, and only including the variables that appeared to be important for explaining variation (using the `envfit` results from above in this example).

```

# partition variation among two predictor tables:
# 1) soil elemental chemistry ('N', 'Al', 'Fe', and 'Mn')
# 2) soil exposure ('Baresoil' and 'Humdepth')
vare.var <- varpart(varespec,
                    ~ N + Al + Fe + Mn,
                    ~ Baresoil + Humdepth,
                    data=varechem.scl, scale=T)

# plot variation association with each partition ('bg' used to set colours)
plot(vare.var, bg=1:3, Xnames=c('chemistry', 'exposure'), id.size=0.75)

# show variation associated with each partition and across both partitions
vare.var

##
## Partition of variance in RDA

```

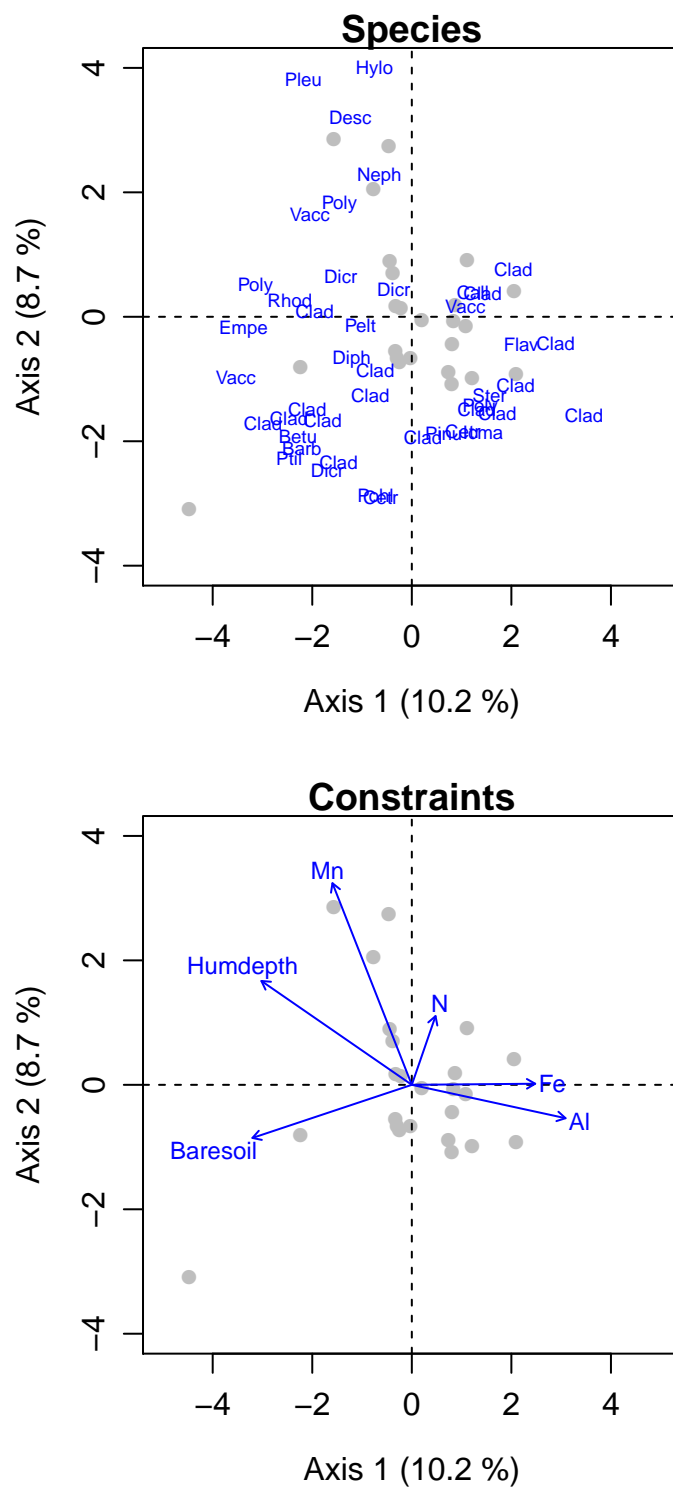


Figure 1.13: RDA ordination of sites based on variation in plant communities, showing plant species loadings (top) and loadings for constraining variables (bottom)

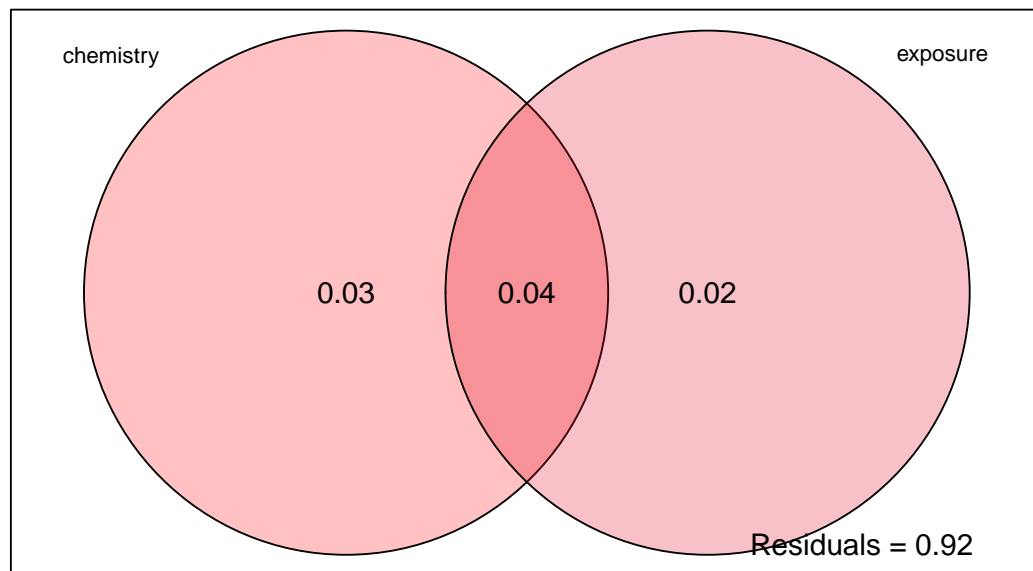


Figure 1.14: Venn diagram showing variation partitioned to variables associated individually with soil chemistry, with soil exposure, or across both.


```
##
## Call: varpart(Y = varespec, X = ~N + Al + Fe + Mn, ~Baresoil +
## Humdepth, data = varechem.scl, scale = T)
## Columns of Y were scaled to unit variance
##
## Explanatory tables:
## X1: ~N + Al + Fe + Mn
## X2: ~Baresoil + Humdepth
##
## No. of explanatory tables: 2
## Total variation (SS): 1012
##          Variance: 44
## No. of observations: 24
##
## Partition table:
##          Df R.squared Adj.R.squared Testable
## [a+b] = X1      4  0.22774      0.06516    TRUE
## [b+c] = X2      2  0.13468      0.05227    TRUE
## [a+b+c] = X1+X2  6  0.32134      0.08182    TRUE
## Individual fractions
## [a] = X1|X2      4          0.02955    TRUE
## [b]              0          0.03561   FALSE
## [c] = X2|X1      2          0.01666    TRUE
## [d] = Residuals          0.91818   FALSE
## ---
## Use function 'rda' to test significance of fractions of interest
```

The partition table displayed above has two sections. The first contains the total amount of variation associated with each partition and across both partitions. The second section of the partition table looks at the individual fractions, or the amount of variation that is attributed solely to each partition and the amount that cannot be attributed to an individual partition. The Venn diagram in Figure 1.14 summarises the results in the second section.

The statistical significance of the values in the first section of the partition table can be tested by using the `anova()` function on an `rda()` object.

```
# test significance of variation in partition 'X1' (chemistry)
anova(rda(varespec ~ N + Al + Fe + Mn, data=varechem.scl, scale=T))

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(formula = varespec ~ N + Al + Fe + Mn, data = varechem.scl, scale = T)
##          Df Variance      F Pr(>F)
## Model      4  10.021 1.4008 0.025 *
## Residual  19  33.979
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# test significance of variation in partition 'X2' (exposure)
anova(rda(varespec ~ Baresoil + Humdepth, data=varechem.scl, scale=T))

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
```

```
##
## Model: rda(formula = varespec ~ Baresoil + Humdepth, data = varechem.scl, scale = T)
##           Df Variance      F Pr(>F)
## Model      2      5.926 1.6343 0.008 **
## Residual  21     38.074
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# test significance of variation in both partitions
anova(rda(varespec ~ N + Al + Fe + Mn + Baresoil + Humdepth,
          data=varechem.scl, scale=T))

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(formula = varespec ~ N + Al + Fe + Mn + Baresoil + Humdepth, data = varechem.scl, scale = T)
##           Df Variance      F Pr(>F)
## Model      6     14.139 1.3416 0.025 *
## Residual  17     29.861
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

These are all significant, which is not surprising because we already demonstrated this in the section [1.3.2.1](#). The statistical significance of the individual fractions in the second section of the partition table can also be tested by using the `anova()` function on an `rda()` object that uses the `Condition()` function to remove variation associated with variables in the other partition.

```
# test significance of variation in partition 'X1' (chemistry) after
# accounting for variation in 'X2'
anova(rda(varespec ~ N + Al + Fe + Mn
          + Condition(Baresoil + Humdepth), data=varechem.scl, scale=T))

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(formula = varespec ~ N + Al + Fe + Mn + Condition(Baresoil + Humdepth), data = varechem.scl, scale = T)
##           Df Variance      F Pr(>F)
## Model      4      8.2132 1.169 0.179
## Residual  17     29.8608

# test significance of variation in partition 'X2' (exposure) after
# accounting for variation in 'X1'
anova(rda(varespec ~ Baresoil + Humdepth
          + Condition(N + Al + Fe + Mn), data=varechem.scl, scale=T))

## Permutation test for rda under reduced model
## Permutation: free
## Number of permutations: 999
##
## Model: rda(formula = varespec ~ Baresoil + Humdepth + Condition(N + Al + Fe + Mn), data = varechem.scl, scale = T)
##           Df Variance      F Pr(>F)
## Model      2      4.1185 1.1724 0.229
## Residual  17     29.8608
```

Neither of the individual fractions, associated only with the 'exposure' partition or the 'chemistry' par-

tion, is significant. This is, again, not surprising – the variation explained by each of these individual partitions is small and there is a larger fraction of variation that is explained but cannot be partitioned individually (the overlapping section of the Venn diagram in [1.14](#)).

1.3.3.1 Variation Partitioning - incorporating spatial processes

In some cases, we have information on the physical locations from which our samples were collected. We can include these spatial characteristics in our analyses because the patterns linked to these variables can represent processes associated with unmeasured (and spatially autocorrelated) environmental variation and/or species dispersal among each of the locations. To do this, we use the `pcnm()` function in the `vegan` library. This function calculates Principal Coordinates of Neighbour Matrices (PCNMs) to generate a dataframe containing variables that represent different spatial scales. We demonstrate this here using the 'mite' data that come with `vegan`, including count data for 35 mite species observed across 70 samples.

```
# load tables containing species, environmental variables, and geographic coordinates

# species-sample table - species in columns, samples in rows
data(mite)
dim(mite)

## [1] 70 35

# five environmental variables associated with each sample location
data(mite.env)
summary(mite.env)

##      SubsDens      WatrCont      Substrate      Shrub      Topo
## Min.   :21.17  Min.   :134.1  Sphagn1 :25  None:19  Blanket:44
## 1st Qu.:30.01  1st Qu.:314.1  Sphagn2 :11  Few :26  Hummock:26
## Median :36.38  Median :398.5  Sphagn3 : 1  Many:25
## Mean   :39.28  Mean   :410.6  Sphagn4 : 2
## 3rd Qu.:46.81  3rd Qu.:492.8  Litter  : 2
## Max.   :80.59  Max.   :827.0  Barepeat : 2
##                               Interface:27

# 'x' and 'y' coordinates associated with each sample location
data(mite.xy)
plot(mite.xy)

# calculate PCNMs from a Euclidean distance matrix of sample coordinates
# and extract scores associated with these new variables
# convert to data.frame for downstream steps
mite.pcnm <- as.data.frame(scores(pcnm(dist(mite.xy))))
dim(mite.pcnm)

## [1] 70 43
```

The patterns represented by the variables in the resulting `pcnm` object are fairly straightforward when samples are evenly spaced in one direction (see Borcard and Legendre, 2002, *Ecological Modelling* 153:51–68). Two-dimensional sampling, especially when samples are not evenly spaced, results in variables with patterns that are not as easy to interpret, but the first axes tend to represent large scale patterns while the later axes tend to represent smaller scale patterns. In figure [1.16](#) we use colour to plot the loadings on the first six PCNM axes to visualise the spatial patterns that they represent.

```
# set up a multipanel graphics window
par(mfrow=c(2, 3))
```

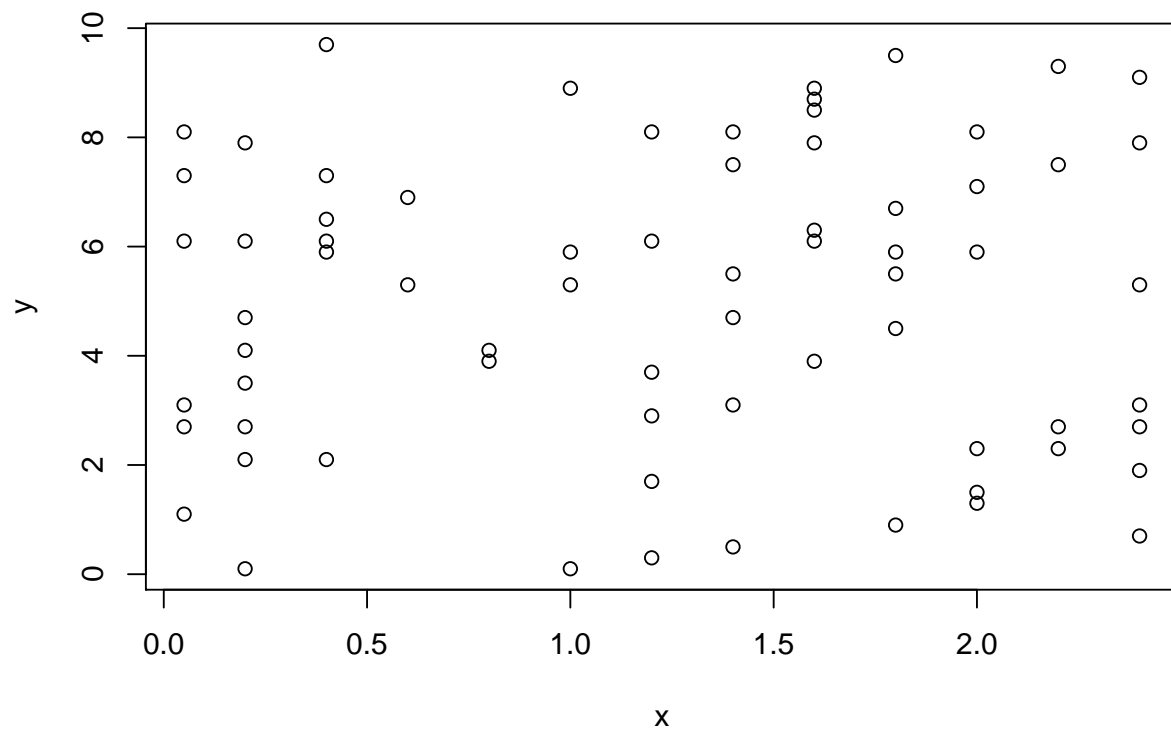


Figure 1.15: Spatial locations of mite samples, along two dimensions (x, y).

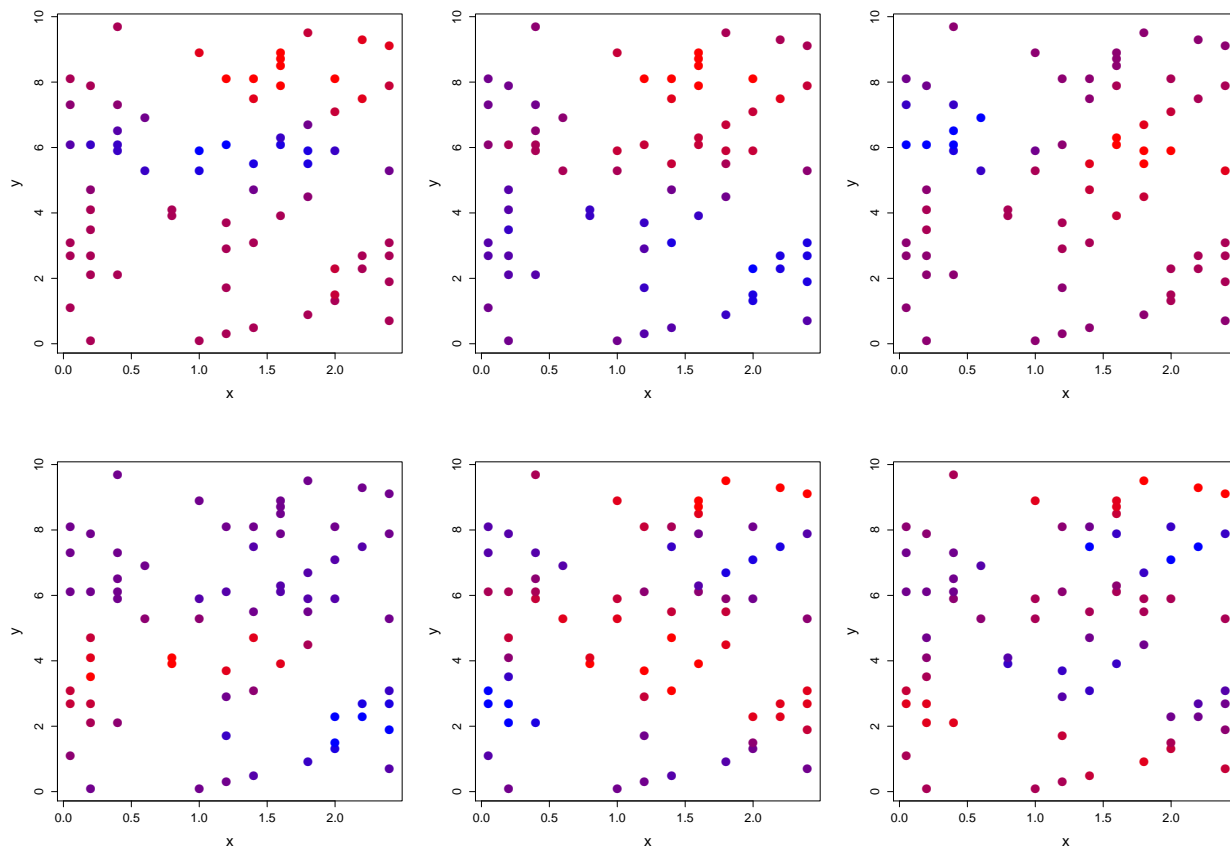


Figure 1.16: Loadings associated with six PCNM axes, plotted against the geographic position of where each sample was collected. Loadings are scaled from positive (red) to negative (blue).

```
# set colour palette with ten levels along a gradient from red to blue
# from Chapter 4
blueredfun <- colorRampPalette(c("blue", "red"))
palette(blueredfun(10))

# for each of the first six PCNM axes, use colour to represent loadings
plot(mite.xy, pch=16, col=cut(mite.pcnm[[1]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
plot(mite.xy, pch=16, col=cut(mite.pcnm[[2]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
plot(mite.xy, pch=16, col=cut(mite.pcnm[[3]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
plot(mite.xy, pch=16, col=cut(mite.pcnm[[4]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
plot(mite.xy, pch=16, col=cut(mite.pcnm[[5]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
plot(mite.xy, pch=16, col=cut(mite.pcnm[[6]], breaks=10), cex.lab=1.5, cex.axis=1.3, cex=2)
```

Now that we have two matrices, one representing variation in measured environmental variables (`mite.env`) and the other representing spatial distributions of samples (`mite.pcnm`), we can estimate the amount of variation in community composition that each explains (shown in Figure 1.17).

```
# do predictor matrices explain community composition, and how much?

# partition variation among three predictor tables:
# 1) substrate ('Substrate', 'SubsDens', and 'WatrCont')
# 2) landscape, i.e., shrub density and microtopography ('Shrub' and 'Topo')
```

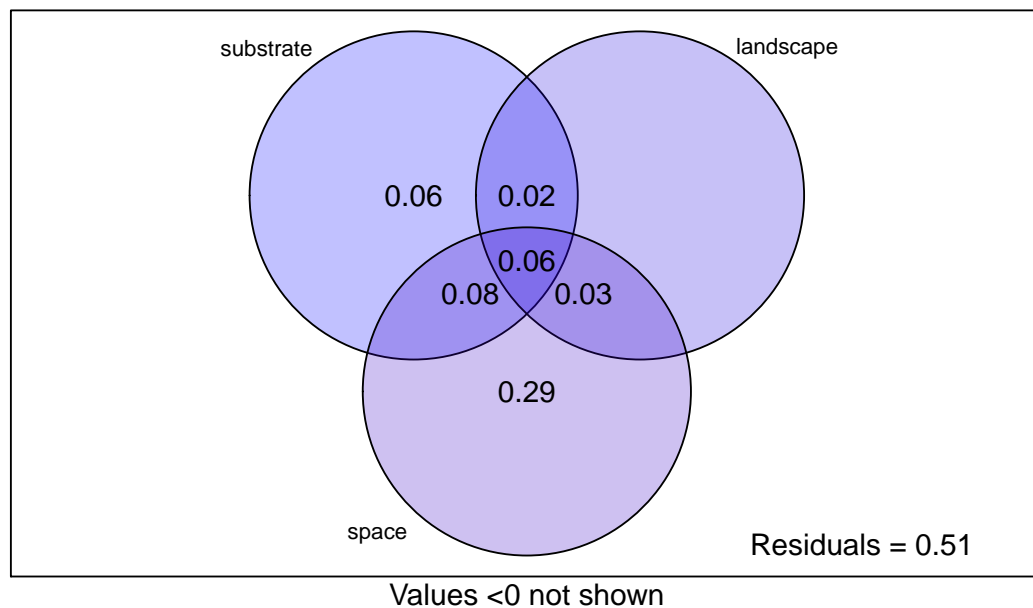


Figure 1.17: Venn diagram showing variation partitioned to variables associated individually with soil chemistry, with soil exposure, with PCNM axes, or across multiple partitions.

```
# 3) space ('mite.pcnm')
mite.var <- varpart(mite,
                    ~ Substrate + SubsDens + WatrCont,
                    ~ Shrub + Topo,
                    mite.pcnm, data=mite.env)
plot(mite.var, bg=1:3, Xnames=c('substrate', 'landscape', 'space'), id.size=0.75)
```

The individual fraction associated with 'landscape' is missing because this number is negative. These numbers represent R^2 values after adjusting for the number of explanatory variables in each partition ('adjusted R^2 ') and will be negative when the raw R^2 is very small

Try this yourself In the above example, use `ordistep()` to select significant PCNM axes and repeat the variation partitioning. Then, use `rda()` to evaluate the significance of variation explained by each of the individual partitions.

1.3.4 Using categorical variables in canonical analyses

In the variation partitioning example above, we used the species-sample matrix as our response matrix (essentially redundancy analysis), not a distance matrix (principal coordinates analysis) even though our responses are counts. This is because there is a debate in the literature about whether it is ap-

propriate to use a distance-based approach for variation partitioning, with the answer tending towards 'no' (see Laliberté, 2008, Ecology 89:3232–3237 and papers cited within). We can perform principal coordinates analysis using the `capscale` function when visualising patterns in the response and relationships with the explanatory variables and evaluating the significance of particular explanatory variables.

```
# use decostand to perform 'Hellinger' transformation of mite data
mite.hel <- decostand(mite, method='hellinger')

# set up full and null models for 'ordistep'
mite.cap1 <- capscale(mite.hel ~ ., data=mite.env, dist='bray')
mite.cap0 <- capscale(mite.hel ~ 1, data=mite.env, dist='bray')

# perform forward and backward selection of explanatory variables
# output not shown
step.env <- ordistep(mite.cap0, scope=formula(mite.cap1))

# look at the significant variables
step.env$anova

##              Df      AIC      F Pr(>F)
## + WatrCont    1 151.62 26.2082 0.005 **
## + SubsDens    1 147.61  6.0122 0.005 **
## + Topo        1 142.31  7.2483 0.005 **
## + Shrub       2 141.48  2.2881 0.005 **
## + Substrate   6 142.65  1.6169 0.005 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# view ordination
plot(step.env)
```

From the ANOVA table, we can see that all of the variables are significant, but we cannot see which levels are important for the factor variables. We can kind of see this from the ordination plot (the levels that are furthest from the centre are most likely to be significant), but it would be useful to do a formal statistical test for each factor level. We can use `ordistep` to do this after transforming the factor levels each into an individual 'dummy' variable. For this example, we use the `dudi.mix` function in the `ade4` library to do this.

```
# load the 'ade4' library
# note the warning due to multiple packages having a function called 'cca'
library('ade4')

# transform the factor variables into individual dummy variables
# modified variables stored in the 'tab' element of the resulting list object
# use 'scannf' and 'nf' arguments to suppress dialogue
mite.env.mod <- dudi.mix(mite.env, scannf=F, nf=2)$tab

# set up full and null models for 'ordistep'
mite.cap1 <- capscale(mite.hel ~ ., data=mite.env.mod, dist='bray')
mite.cap0 <- capscale(mite.hel ~ 1, data=mite.env.mod, dist='bray')

# perform forward and backward selection of explanatory variables
# (output not shown)
step.env <- ordistep(mite.cap0, scope=formula(mite.cap1))
```

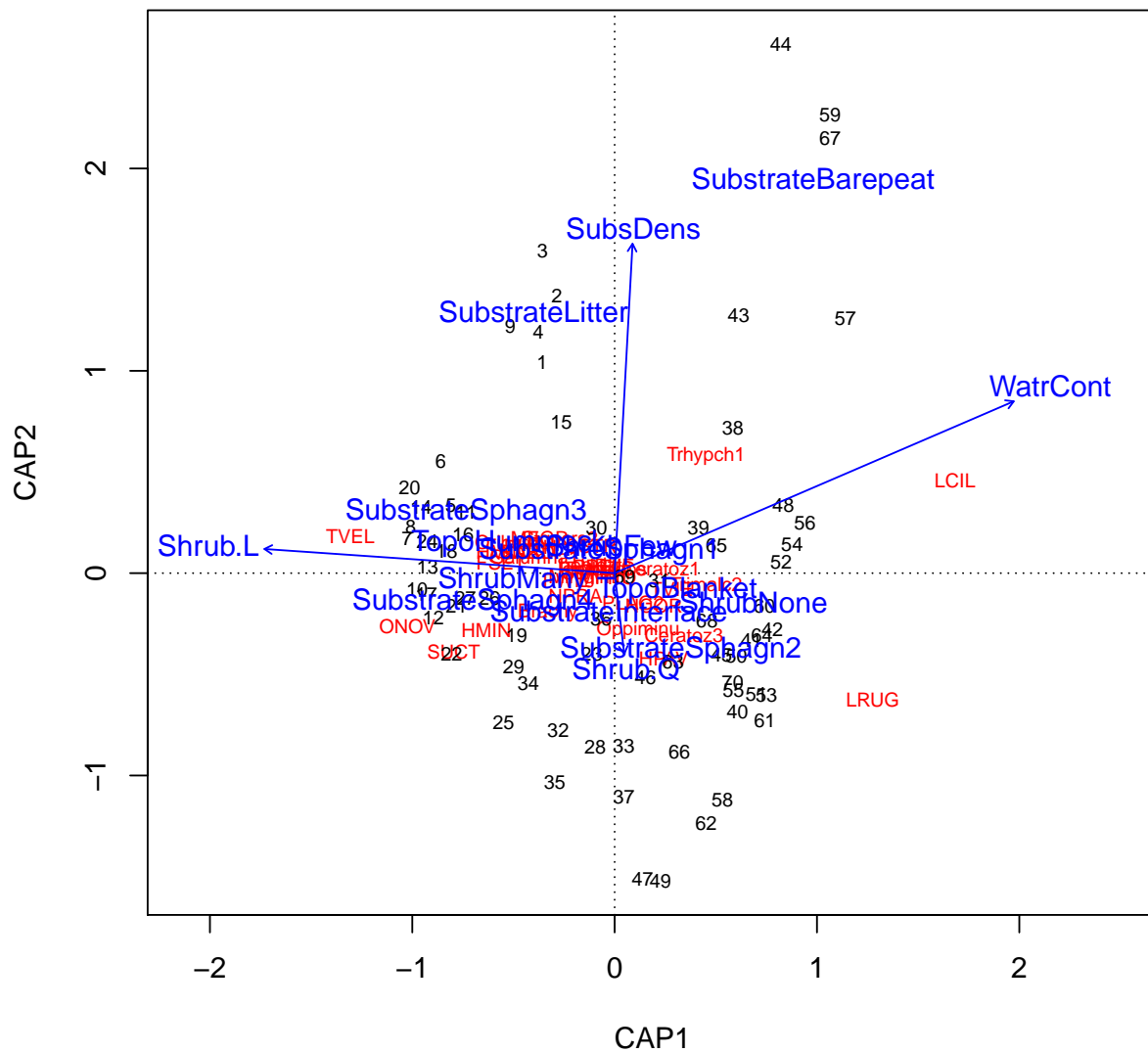


Figure 1.18: Ordination of mite data based on CAP analysis including continuous and categorical variables.


```

# look at the significant variables
step.env$anova

##           Df      AIC      F Pr(>F)
## + WatrCont      1 151.62 26.2082 0.005 **
## + SubsDens      1 147.61  6.0122 0.005 **
## + Topo.Blanket   1 142.31  7.2483 0.005 **
## + Subst.Barepeat 1 141.47  2.6938 0.005 **
## + Subst.Sphagn1  1 141.03  2.2680 0.005 **
## + Shrub.Q        1 140.37  2.4396 0.005 **
## + Shrub.L        1 139.78  2.3439 0.005 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# view ordination
plot(step.env)

```

The nonsignificant explanatory variables are not plotted, making it easier to see relationships among significant variables, samples, and mite species.

1.3.4.1 Indicator species analysis

We can visualise relationships between species and environmental variables by looking at ordinations, but it may be useful to identify species that are significantly associated with particular environmental variables. This can be done through Indicator Species analysis using functions in the `labdsv` library or the `indicspecies` library. Below we show examples using the `indval` function from the `labdsv` library.

```

# load library
library(labdsv)

## example using a factor variable

# calculate indicator values for each species
ind.topo <- indval(mite, mite.env$Topo)

# calculate adjusted P-values for each species and show significant species
topo.pval <- p.adjust(ind.topo$pval, method='bonferroni')
topo.pval[topo.pval < 0.05]

##      PTHH      RARD      TVEL      ONOV      SUCT Oribatl1 Galumna1      FSET
##      0.035      0.035      0.035      0.035      0.035      0.035      0.035      0.035
##      LRUG
##      0.035

# show indicator values for significant indicator species in each group
topo.ind <- ind.topo[['indval']]
topo.ind[rownames(topo.ind) %in% names(topo.pval[topo.pval < 0.05]), ]

##           Blanket      Hummock
## PTHH      0.03721177 0.48943709
## RARD      0.01369430 0.54122471
## TVEL      0.07548822 0.71867805
## ONOV      0.22046222 0.73782871
## SUCT      0.29871725 0.67942539
## Oribatl1 0.06621842 0.53720770
## Galumna1 0.03379908 0.54522598

```

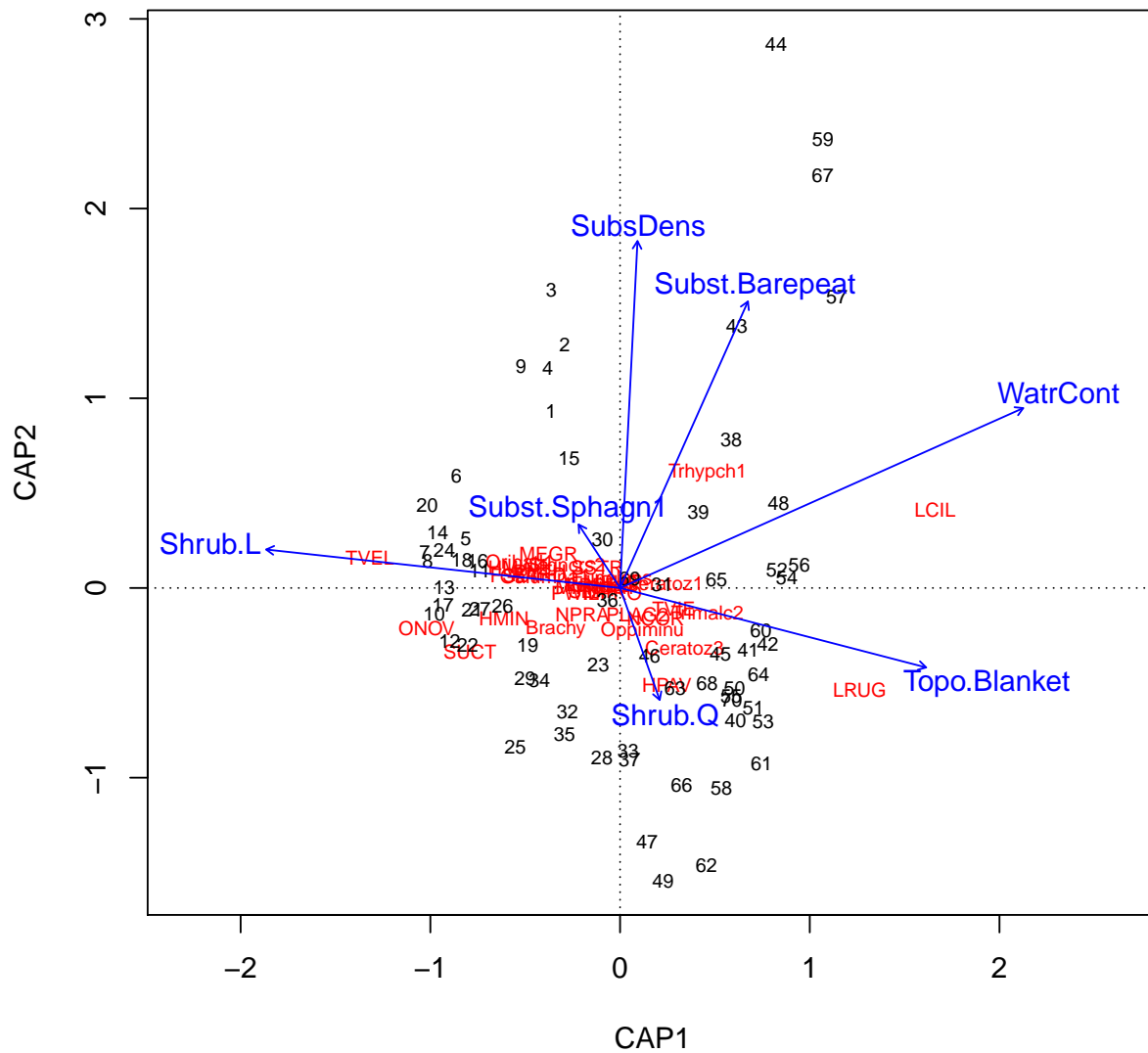


Figure 1.19: Ordination of mite data based on CAP analysis including continuous and categorical variables after removal of nonsignificant levels in categorical variables.

```
## FSET      0.04719869 0.57249563
## LRUG      0.79702367 0.04267176

# show frequencies in each group for significant indicator species
topo.frq <- ind.topo[['relfrq']]
topo.frq[rownames(topo.frq) %in% names(topo.pval[topo.pval < 0.05]), ]

##          Blanket   Hummock
## PTHH      0.1818182 0.6153846
## RARD      0.1136364 0.6153846
## TVEL      0.3409091 0.9230769
## ONOV      0.8409091 1.0000000
## SUCT      0.9318182 1.0000000
## Oribat11  0.2500000 0.7307692
## Galumna1  0.1590909 0.6923077
## FSET      0.2727273 0.6923077
## LRUG      0.9090909 0.3461538

## example using a continuous variable

# create grouping categories from the continuous variable
groups <- cut(mite.env$WatrCont, breaks=3)

# calculate indicator values for each species
ind.dens <- indval(mite, groups)

# calculate adjusted P-values for each species and show significant species
dens.pval <- p.adjust(ind.dens$pval, method='bonferroni')
dens.pval[dens.pval < 0.05]

## TVEL  ONOV  SUCT  NCOR
## 0.035 0.035 0.035 0.035

# can use code from the first example to interpret group membership for indicator species
```

1.3.5 'Experimental' frameworks: more working with factors

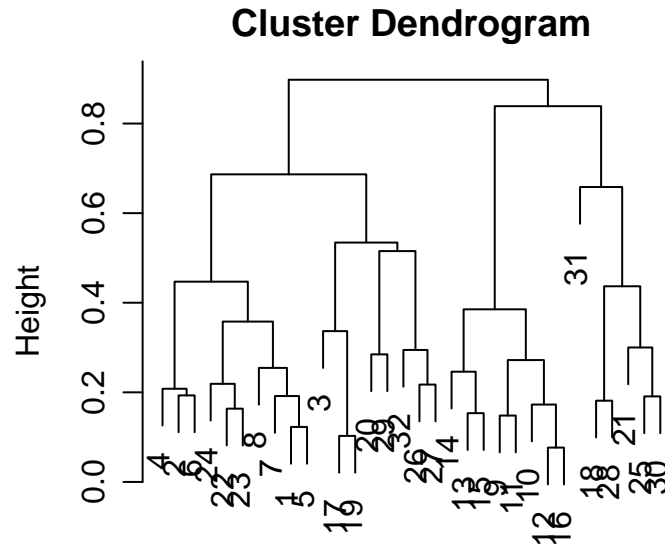
1.3.5.1 Cluster analysis

We often collect multivariate data in the context of experimental studies or in the context of observed categorical explanatory variables that are of interest (as seen in section 1.3.4). We could perform ANOVA for each variable individually, but this amplifies the potential for Type I error and does not account for potential collinearity among response variables. In addition, we may not be interested in particular response variables, but in how multiple variables are responding as a whole. In this example, we use the Tibetan Plateau plant community data from Section ???. We will use the complete-linkage algorithm, which is the default and aims to maximise the distances among clusters, but other algorithms are available and may be more appropriate.

```
tibplat <- read.csv('tibplat.csv')

# Species data are in columns 3:10
tib.clust <- hclust(vegdist(tibplat[,3:10], 'bray'))
plot(tib.clust)
```

The dendrogram in Fig. 1.20 indicates that the plant communities appear to be hierarchically clustered



```
vegdist(tibplat[, 3:10], "bray")
hclust (*, "complete")
```

Figure 1.20: Hierarchical clustering of Tibetan Plateau plant communities based on Bray-Curtis dissimilarities.

into groups. The label on the tips represent the rownames in the dataframe but are not very informative on their own. Here it is useful to provide a more informative vector with which to label the tips. We will use the information on the management practices associated with each plant community (Fig. 1.21).

```
tib.clust <- hclust(vegdist(tibplat[,3:10], 'bray'))
plot(tib.clust, labels=tibplat[['fertilization']], main='fertilization')
plot(tib.clust, labels=tibplat[['enclosure']], main='enclosure')
```

Groups do appear to be linked to fertilization and grazing, but there are a few examples in which a plant community of one management type is quite divergent from others of that type. We need a statistical tool to test the hypothesis that plant communities within each management type are more similar to each other than to those under other forms of management.

Analysis of similarities (ANOSIM) is one approach to test this hypothesis, using the `anosim()` function in the `vegan` package. This is a nonparametric test; distances are transformed into ranks and the mean of the ranks for distances classified as within the groups is compared to that for the between-group distance ranks. However, there are concerns about how ANOSIM results should be interpreted (also mentioned on the help page for the function) and PerMANOVA (Permutational ANOVA) is usually preferred.

PerMANOVA is potentially more powerful than ANOSIM because the distances are preserved and not rank-transformed; the approach partitions variance among and within groups through calculations of sums of squared distances and calculates pseudo *F*-statistics. In addition, it is possible to estimate the effects associated with multiple factors and their interactions. The function to perform PerMANOVA is `adonis` within `vegan` and the result is interpreted in the same way as an ANOVA table.

```
adonis(vegdist(tibplat[,3:10]) ~ fertilization*enclosure, data=tibplat)
##
```

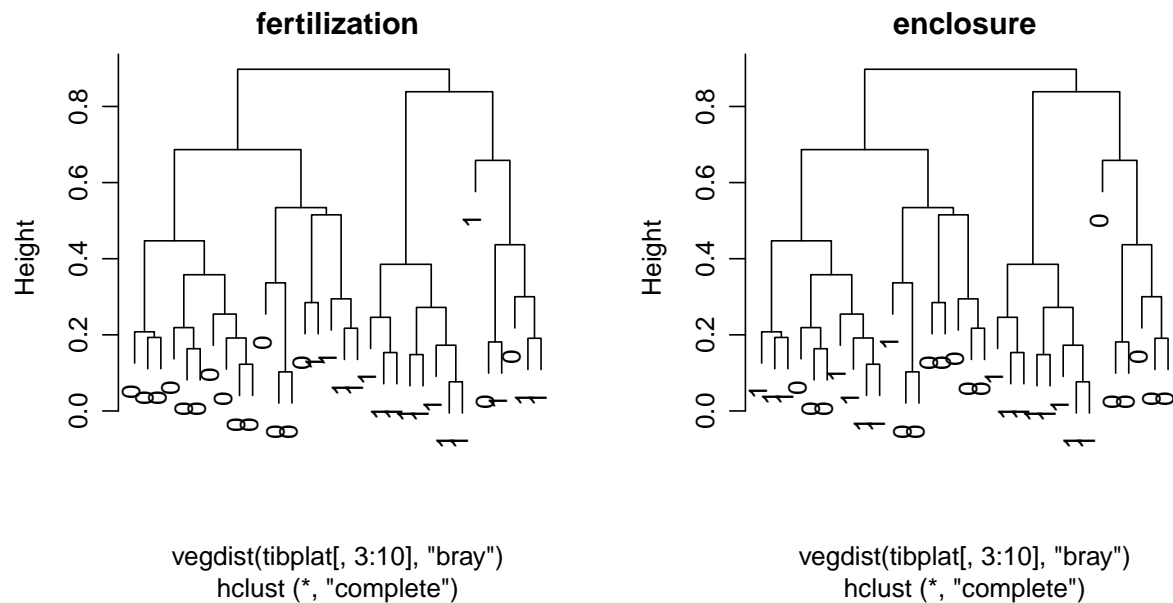


Figure 1.21: Hierarchical clustering of Tibetan Plateau plant communities based on Bray-Curtis dissimilarities, with tips labelled by management system (fertilization, enclosed).

```
## Call:
## adonis(formula = vegdist(tibplat[, 3:10]) ~ fertilization * enclosure,      data = tibplat)
##
## Permutation: free
## Number of permutations: 999
##
## Terms added sequentially (first to last)
##
##              Df SumsOfSqs MeanSqs F.Model    R2 Pr(>F)
## fertilization    1    1.2320 1.23202   20.757 0.25944  0.001 ***
## enclosure        1    1.1358 1.13583   19.137 0.23918  0.001 ***
## fertilization:enclosure 1    0.7190 0.71901   12.114 0.15141  0.001 ***
## Residuals       28    1.6619 0.05935     0.34997
## Total           31    4.7488              1.00000
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

1.3.5.2 More manipulating graphics from *vegan* objects

Let's produce a figure that visualises the relationship between the experimental treatments and the results of the PCoA used in the PERMANOVA in the previous section.

```
# perform PCoA and store results in object
tib.pco <- wcmdscale(vegdist(tibplat[,3:10]), eig=T)

# return the proportions of variance explained
tib.pco$eig[tib.pco$eig >= 0] / sum(tib.pco$eig[tib.pco$eig >= 0])

## [1] 0.4569200850 0.1601330893 0.1144687053 0.0828594613 0.0520696801
## [6] 0.0447359716 0.0256728499 0.0198370466 0.0136060804 0.0091392099
```

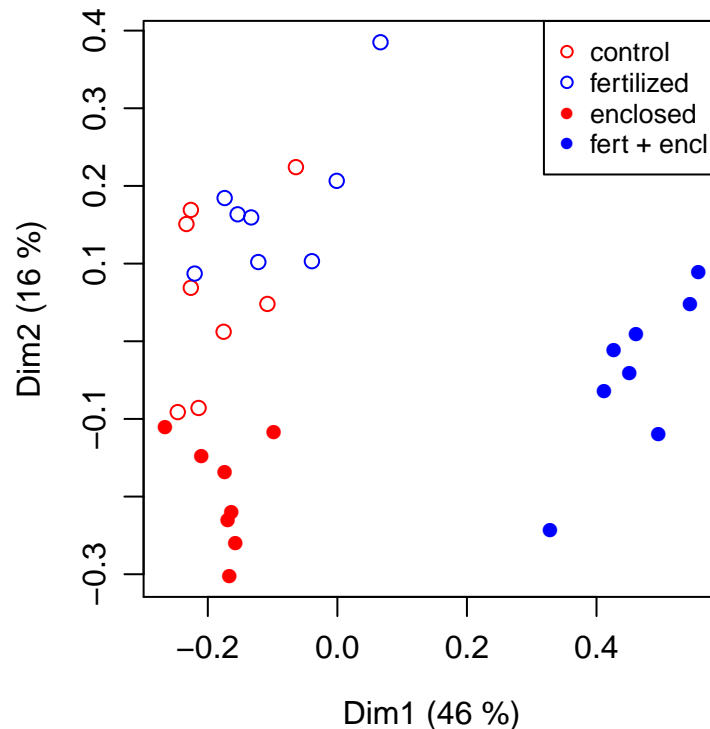


Figure 1.22: MDS ordination of sites based on variation in plant communities

```
## [11] 0.0067712217 0.0057549299 0.0038943360 0.0027786236 0.0011320982
## [16] 0.0002266113

palette(c('red','blue')) # set up the colour palette
with(tibplat, plot(scores(tib.pco, display='sites'), # identifies coordinates
                  pch=c(1,16)[enclosure], col=fertilization, # assign symbols, colours
                  xlab='Dim1 (46 %)', ylab='Dim2 (16 %)')) # change axis labels
legend('topright', # position the legend
      legend=c('control', 'fertilized', 'enclosed', 'fert + encl'), # set legend text
      pch=c(1,1,16,16), col=c('red','blue','red','blue'), # assign symbols and colours
      cex=0.8) # slightly reduce the size of the text and points
```

Try this yourself In the above example, adjust the `par` settings to modify the plot to your liking.

Patterns among treatments may be difficult to see in the ordination in some cases. The `vegan` package includes functions to include shapes that delimit boundaries around groups of points (the `ordihull` function) or a certain distance from the centroid (the `ordiellipse` function). Another `vegan` function (`ordispider`) connects each point to the group centroid with lines.

```
# perform PCoA and store results in object
tib.pco <- wcmdscale(vegdist(tibplat[,3:10]), eig=T)

# return the proportions of variance explained
tib.pco$eig[tib.pco$eig >= 0] / sum(tib.pco$eig[tib.pco$eig >= 0])

## [1] 0.4569200850 0.1601330893 0.1144687053 0.0828594613 0.0520696801
## [6] 0.0447359716 0.0256728499 0.0198370466 0.0136060804 0.0091392099
```

```
## [11] 0.0067712217 0.0057549299 0.0038943360 0.0027786236 0.0011320982
## [16] 0.0002266113

# set plot options
par(mfrow=c(3, 1)) # set plot layout
palette(c('red','blue')) # set up the colour palette

# create a vector containing all treatment combinations
trts <- with(tibplat, interaction(fertilization, enclosure, sep='-'))
levels(trts)

## [1] "0-0" "1-0" "0-1" "1-1"

# draw a polygon around all points in a treatment
with(tibplat, plot(scores(tib.pco, display='sites'), # identifies coordinates
                  pch=c(1,16)[enclosure], col=fertilization, # assign symbols, colours
                  xlab='Dim1 (46 %)', ylab='Dim2 (16 %)',
                  main='ordihull')) # change axis labels
legend('topright', # position the legend
      legend=c('control', 'fertilized', 'enclosed', 'fert + encl'), # set legend text
      pch=c(1,1,16,16), col=c('red','blue','red','blue'), # assign symbols and colours
      cex=0.8) # slightly reduce the size of the text and points
ordihull(tib.pco, trts, draw='polygon', # draw polygon
        col=c('red','blue','red','blue'), lty=c('dashed', 'dashed', 'solid', 'solid'))

# draw an ellipse representing the 95% confidence interval from the centroid
with(tibplat, plot(scores(tib.pco, display='sites'), # identifies coordinates
                  pch=c(1,16)[enclosure], col=fertilization, # assign symbols, colours
                  xlab='Dim1 (46 %)', ylab='Dim2 (16 %)',
                  main='ordiellipse')) # change axis labels
legend('topright', # position the legend
      legend=c('control', 'fertilized', 'enclosed', 'fert + encl'), # set legend text
      pch=c(1,1,16,16), col=c('red','blue','red','blue'), # assign symbols and colours
      cex=0.8) # slightly reduce the size of the text and points
ordiellipse(tib.pco, trts, kind='se', conf=0.95, # draw ellipse
            col=c('red','blue','red','blue'), lty=c('dashed', 'dashed', 'solid', 'solid'))

# draw a spider graph linking each point to the group centroid
with(tibplat, plot(scores(tib.pco, display='sites'), # identifies coordinates
                  pch=c(1,16)[enclosure], col=fertilization, # assign symbols, colours
                  xlab='Dim1 (46 %)', ylab='Dim2 (16 %)',
                  main='ordispider')) # change axis labels
legend('topright', # position the legend
      legend=c('control', 'fertilized', 'enclosed', 'fert + encl'), # set legend text
      pch=c(1,1,16,16), col=c('red','blue','red','blue'), # assign symbols and colours
      cex=0.8) # slightly reduce the size of the text and points
ordispider(tib.pco, trts, spiders='centroid', # draw spiders
           col=c('red','blue','red','blue'), lty=c('dashed', 'dashed', 'solid', 'solid'))
```

1.3.5.3 Model-based approaches

Up to this point, we have used distance-based approaches to detect patterns and to make inferences about the drivers of these patterns. 'Distance-based' essentially means that the data are transformed prior to analysis so that the similarities and dissimilarities among each of the observations can be rep-

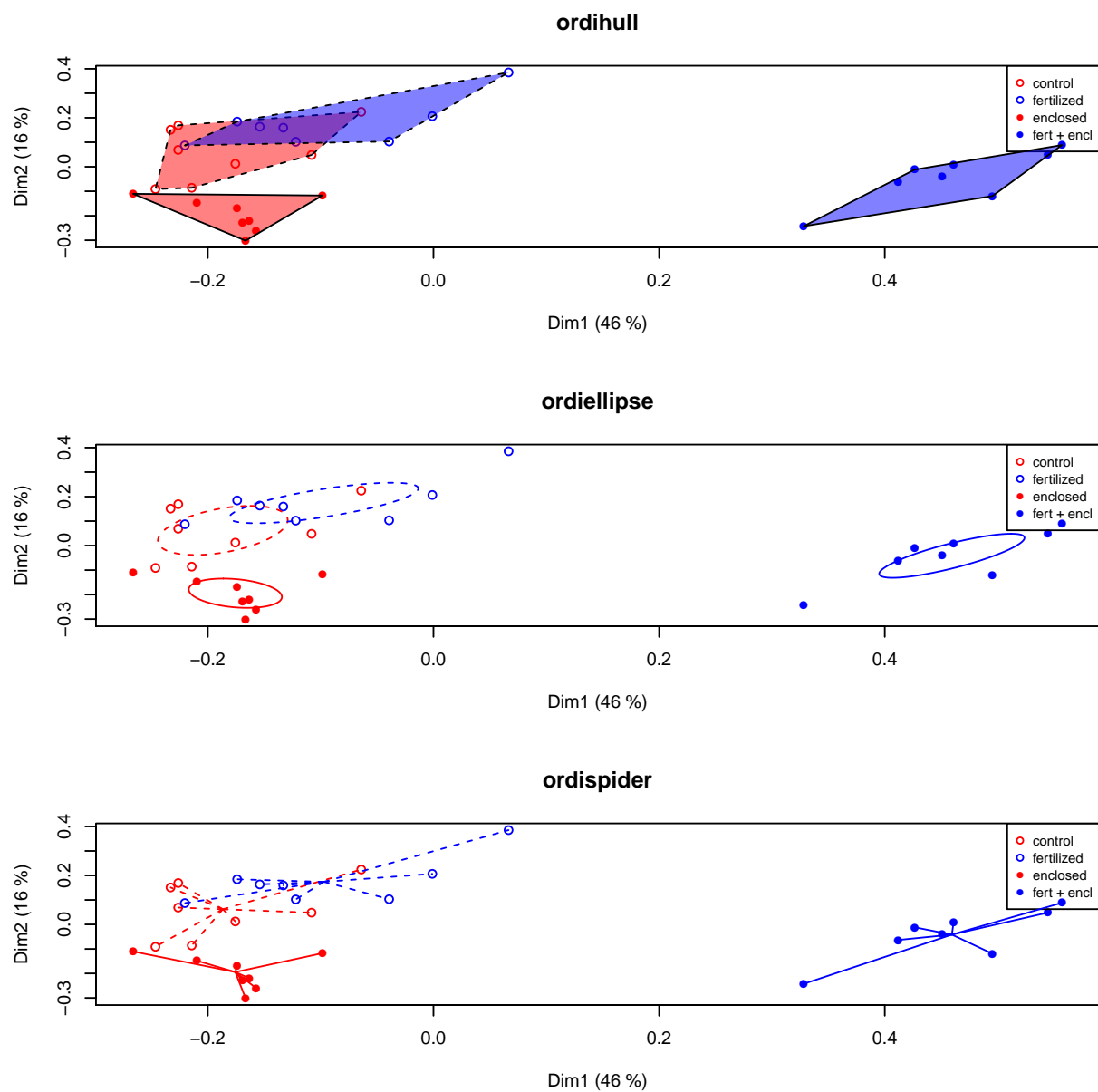


Figure 1.23: MDS ordination of sites based on variation in plant communities, further highlighting treatment groupings.

resented by their distances in multivariate space. These approaches were developed to account for the limited computational power available at the time, which would have made model-based approaches infeasible.

Model-based approaches are what we used in Chapter ??, which involved fitting models the observed data and estimating effect sizes assuming that the data correspond to particular distributions. While computationally difficult for multivariate data, model-based approaches are potentially more powerful and less susceptible to error, at least when focussing on the effects of particular environmental drivers. An important question to ask when deciding between these general types of approaches is whether you are interested in observing/predicting patterns in data ('distance-based') or changes in variables ('model-based').

For responses exhibiting normal error distributions, it is appropriate to use MANOVA (multivariate analysis of variance). We will use the I x F data from Section ?? to demonstrate MANOVA. First, we need to convert the response variables in the dataframe into a matrix object.

```
hfeIxF <- read.csv("HFEIFplotmeans.csv")
str(hfeIxF)

## 'data.frame': 320 obs. of 5 variables:
## $ plotnr : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Date : Factor w/ 20 levels "1/1/2008","1/1/2009",...: 3 13 7 5 15 2 6 17 1 19 ...
## $ diameter: num NA 3.96 7.38 NA 4.36 ...
## $ height : num NA 4.05 6.18 NA 4.22 ...
## $ treat : Factor w/ 4 levels "C","F","I","IL": 4 4 4 4 4 4 4 4 4 4 ...

Y <- as.matrix(hfeIxF[,3:4])
m1 <- manova(Y ~ treat*Date, data=hfeIxF)
summary(m1)

## Df Pillai approx F num Df den Df Pr(>F)
## treat 3 0.87681 53.083 6 408 < 2.2e-16 ***
## Date 16 1.83960 146.227 32 408 < 2.2e-16 ***
## treat:Date 48 0.97563 4.048 96 408 < 2.2e-16 ***
## Residuals 204
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

summary.aov(m1)

## Response diameter :
## Df Sum Sq Mean Sq F value Pr(>F)
## treat 3 168.9 56.314 88.9642 < 2.2e-16 ***
## Date 16 4727.4 295.465 466.7754 < 2.2e-16 ***
## treat:Date 48 115.9 2.414 3.8136 1.417e-11 ***
## Residuals 204 129.1 0.633
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Response height :
## Df Sum Sq Mean Sq F value Pr(>F)
## treat 3 150.1 50.042 134.6903 < 2.2e-16 ***
## Date 16 3461.7 216.358 582.3402 < 2.2e-16 ***
## treat:Date 48 157.1 3.273 8.8098 < 2.2e-16 ***
## Residuals 204 75.8 0.372
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## 48 observations deleted due to missingness
```

The analysis provides two sets of significance tests. Using `summary` gives test statistics and their significance for the multivariate response to the predictor variables. Using `summary.aov` gives univariate ANOVA tables for each response variable. We can see that both variables are responding to the treatments, which makes sense given the variables are strongly correlated.

Quite often, our data are not normally distributed (e.g., abundance data) and we rely on GLMs to fit models to appropriate error families. Recent advances in model-based approaches for multivariate data are included in the `mvabund` package. Let's reanalyse the Tibetan Plateau plant community data.

```
library(mvabund)
tibplat <- read.csv("tibplat.csv")
str(tibplat)

## 'data.frame': 32 obs. of 10 variables:
## $ fertilization      : int  0 0 0 0 0 0 0 0 1 1 ...
## $ enclosure          : int  1 1 1 1 1 1 1 1 1 1 ...
## $ elymus_nutans       : num  3.53 4.85 2.2 2.87 3.53 ...
## $ poa_crymophila      : num  0.89 0.6 0 0.3 1.79 ...
## $ kobresia_setchwanensis : num  15.1 14.3 12.9 10.8 17.9 ...
## $ anemone_rivularis    : num  5.12 13.65 0 18.76 5.12 ...
## $ potentilla_fragarioides : num  3.62 3.62 2.36 3.15 2.05 2.52 2.05 1.57 0.86 0.22 ...
## $ stipa_aliena        : num  3.8 3.8 0 1.52 6.08 5.32 4.56 7.6 0 0 ...
## $ astragalus_polycladus : num  2.08 2.08 2.08 0 1.04 8.3 0 0 0 0 ...
## $ anemone_obtusiloba   : num  3.8 3.8 2.66 2.66 2.85 2.47 3.23 0.38 1.89 0 ...

X <- tibplat[,1:2] # create matrix of predictor variables
Y <- mvabund(tibplat[,3:10]) # convert to 'mvabund' object, for use in downstream functions
plot(Y, ylab="")

## Kicking off BoxPlot sequence

## Overlapping points were shifted along the y-axis to make them visible.

##
##
## ABOUT TO PLOT THE FUNCTION
```

We can see in Fig. 1.24 that the data roughly fits a log-normal distribution based on the similar spread of the abundances for most species when plotted on a log scale. However, the data also contain several zero values. We will use a negative binomial error distribution, which will result in a warning since the variables represent biomass and not true abundances. A gamma distribution may be more appropriate here but is not available in the current release. Model diagnostics, as in Figure 1.25 can be viewed using the `plot` function on the model output.

```
# fit model
m1 <- manyglm(Y ~ fertilization*enclosure, data=X, family='negative.binomial')

## Warning in manyglm(Y ~ fertilization * enclosure, data = X, family = "negative.binomial"):
## Non-integer data are fitted to the negative.binomial model.

# view model diagnostic plot
plot(m1)

# look at model summary and ANOVA table
summary(m1)

##
## Test statistics:
```

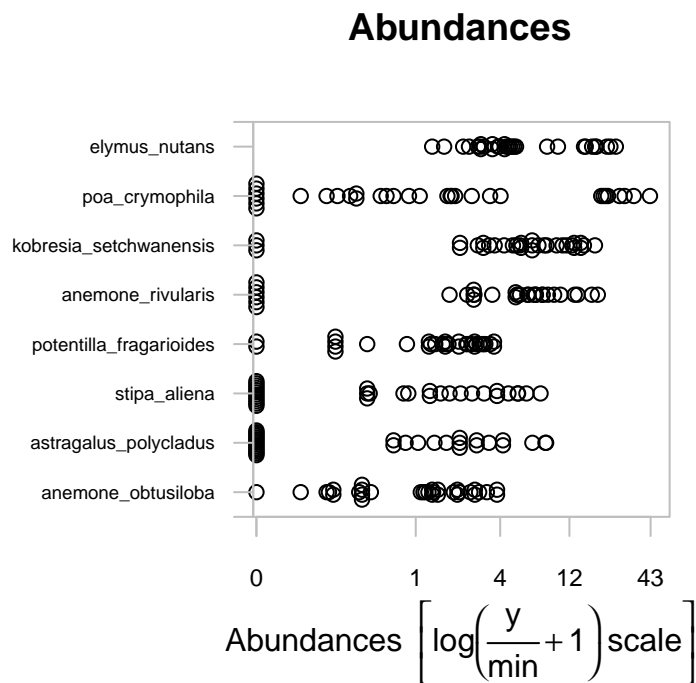


Figure 1.24: Species biomasses per plot in the Tibetan Plateau.

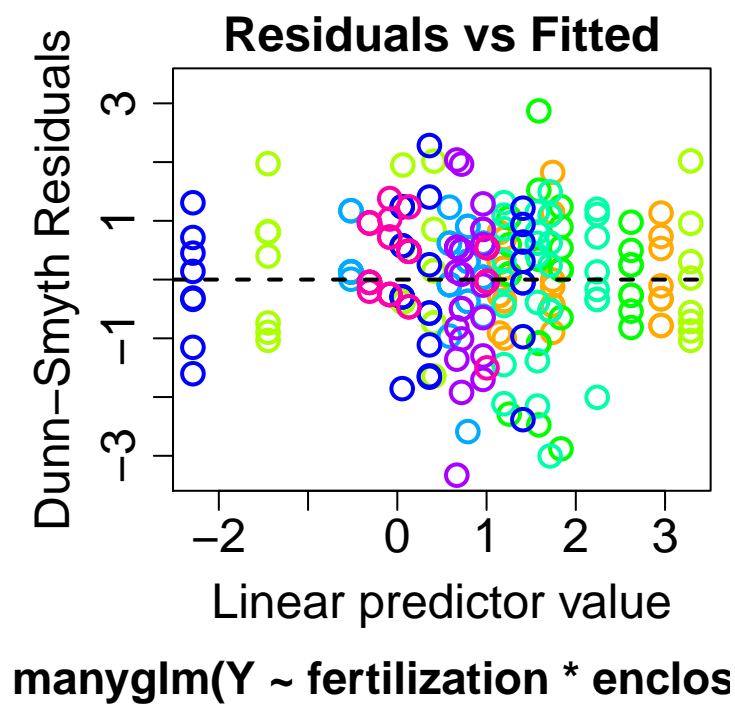


Figure 1.25: Model diagnostic plot for multivariate GLM fit to Tibetan plateau plant data.

```

##                                wald value Pr(>wald)
## (Intercept)                    13.753  0.000999 ***
## fertilization                   4.076  0.013986 *
## enclosure                       6.434  0.000999 ***
## fertilization:enclosure         5.372  0.000999 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Test statistic:    20.2, p-value: 0.000999
## Arguments:
## Test statistics calculated assuming response assumed to be uncorrelated
## P-value calculated using 1000 resampling iterations via pit.trap resampling (to account for correlat

anova(m1, nBoot=100, p.uni='adjusted') # we use only 100 permutations here, to save time

## Time elapsed: 0 hr 0 min 0 sec
## Analysis of Deviance Table
##
## Model: manyglm(formula = Y ~ fertilization * enclosure, family = "negative.binomial",
## Model:      data = X)
##
## Multivariate test:
##
##              Res.Df Df.diff    Dev Pr(>Dev)
## (Intercept)          31
## fertilization        30      1 75.03    0.01 **
## enclosure            29      1 80.96    0.01 **
## fertilization:enclosure 28      1 55.17    0.01 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Univariate Tests:
##
##              elymus_nutans              poa_crymophila
##              Dev Pr(>Dev)              Dev Pr(>Dev)
## (Intercept)
## fertilization        29.644    0.010          14.06    0.010
## enclosure            23.598    0.010          37.826    0.010
## fertilization:enclosure 12.118    0.040          1.498    0.515
##
##              kobresia_setchwanensis              anemone_rivularis
##              Dev Pr(>Dev)              Dev
## (Intercept)
## fertilization        13.008    0.010              2.129
## enclosure            11.308    0.010              3.473
## fertilization:enclosure 2.37    0.465              0.061
##
##              potentilla_fragarioides
##              Pr(>Dev)              Dev Pr(>Dev)
## (Intercept)
## fertilization        0.238              3.81    0.158
## enclosure            0.287              0.557    0.733
## fertilization:enclosure 0.812              11.274    0.040
##
##              stipa_aliena              astragalus_polycladus
##              Dev Pr(>Dev)              Dev
## (Intercept)
## fertilization        4.823    0.109              0.485
## enclosure            0.016    0.960              2.688

```

```
## fertilization:enclosure      15.039    0.020          8.987
##                               anemone_obtusiloba
##                               Pr(>Dev)          Dev Pr(>Dev)
## (Intercept)
## fertilization                0.455          7.075    0.050
## enclosure                    0.287          1.496    0.446
## fertilization:enclosure      0.069          3.82     0.287
## Arguments:
## Test statistics calculated assuming uncorrelated response (for faster computation)
## P-value calculated using 100 resampling iterations via PIT-trap resampling (to account for correlation)
# (not all output shown)
```

The first part of the `anova()` result gives the significance of the predictor variables, while the second part gives their significance in univariate tests, to determine which species are particularly responsive.

Large datasets with lots of species can lead to calls to `summary()` and `anova()` taking a very long time. An alternative approach to assessing the significance of treatments and/or their interactions is via sequential model selection and assessment using the `drop1` function. The output includes AIC scores for the global model ("`<none>`") and each model for which a single term can be dropped. If the AIC is lower for any of the latter models than it is for the global model, this is evidence in support of the simpler model and justifies attempting to simplify the model further. Let's demonstrate using all plant species in the Tibetan Plateau dataset.

```
# set up predictor and response tables
X <- tibplat[,1:2] # create matrix of predictor variables
Y <- mvabund(tibplat[,3:ncol(tibplat)]) # convert to 'mvabund' object, for use in downstream functions

# fit model (ignore warning)
m2 <- manyglm(Y ~ fertilization*enclosure, data=X, family='negative.binomial')

## Warning in manyglm(Y ~ fertilization * enclosure, data = X, family = "negative.binomial"):
## Non-integer data are fitted to the negative.binomial model.

# evaluate effect of dropping interaction from model on overall fit (ignore warning)
drop1(m2)

## Warning in manyglm(formula = Y ~ fertilization + enclosure, family = "negative.binomial",
## : Non-integer data are fitted to the negative.binomial model.

## Single term deletions
##
## Model:
## Y ~ fertilization * enclosure
##           Df    AIC
## <none>      1008.1
## fertilization:enclosure  8 1048.7
```

The AIC score is much higher when the `fertilization:enclosure` term is removed. Therefore, we conclude that the global model containing both fixed effects and the interaction performs best.

1.4 Functions used in this chapter

For functions not listed here, please refer to the index at the end of this book.

All functions in this table are from the `vegan` package unless otherwise noted.

Function	What it does	Example use
<code>summary</code> (base R)	Simple summary table for dataframes and <code>rda</code> and <code>cca</code> objects	<code>summary(pupae)</code>
<code>rda</code>	Principal components analysis (as used in this chapter), and redundancy analysis	Section 1.2
<code>cca</code>	Correspondence analysis. Similar to <code>rda</code> , but using χ^2 distances to estimate divergence between samples (and thus appropriate for count data)	Section 1.2
<code>decorana</code>	Estimate the 'gradient length', to decide between PCA or CA	Section 1.2.2.1
<code>wcmdscale</code>	Principle coordinates analysis, used in conjunction with <code>vegdist</code>	Section 1.2.3
<code>vegdist</code>	Calculate a distance matrix, used primarily as input to <code>wcmdscale</code> , <code>mantel</code> and <code>hclust</code>	Section 1.2.3
<code>rankindex</code>	Compare dissimilarity indices for gradient detection	Section 1.2.3
<code>metaMDS</code>	Non-metric multi-dimensional scaling. Suitable for a wide variety of data.	Section 1.2.4
<code>scores</code>	Extract species or site loadings from a fitted NMDS or other ordination method.	Section 1.2.4
<code>mantel</code>	Test for correlation between two dissimilarity matrices	Section 1.3.1
<code>envfit</code>	Identify variables from an environmental matrix that are correlated with the site loadings	Section 1.3.2.1
<code>ordistep</code>	Use forward and/or backward selection to select environmental variables associated with variation in the response table	Section 1.3.2.1
<code>varpart</code>	Partition variation in a response table to variables in two or more tables of explanatory variables	Section 1.3.3
<code>pcnm</code>	Calculate principal coordinates of neighbour matrices for partitioning variation to spatial patterns	Section 1.3.3.1
<code>hclust</code>	Hierarchical clustering on a dissimilarity matrix	Section 1.3.5.1
<code>adonis</code>	Permutational multivariate analysis of variance with distance matrices (PerMANOVA)	Section 1.3.5.1
<code>manova</code>	(base R) Multivariate analysis of variance	Section 1.3.5.3
<code>manyglm</code>	(<code>mvabund</code> package) Generalized linear models for multivariate abundance data	Section 1.3.5.3
<code>mvabund</code>	(<code>mvabund</code> package) Construct a multivariate abundance object, for use in e.g. <code>manyglm</code>	Section 1.3.5.3

1.5 Exercises

In these exercises, we use the following colour codes:

- **Easy:** make sure you complete some of these before moving on. These exercises will follow examples in the text very closely.
- ◆ **Intermediate:** a bit harder. You will often have to combine functions to solve the exercise in two steps.
- ▲ **Hard:** difficult exercises! These exercises will require multiple steps, and significant departure from examples in the text.

We suggest you complete these exercises in an **R** markdown file. This will allow you to combine code chunks, graphical output, and written answers in a single, easy-to-read file.

1.5.1 Ordination

1.5.1.1 Allometry data

File: 'Allometry.csv'

This dataset contains measurements of tree dimensions and biomass. Data kindly provided by John Marshall, University of Idaho.

Variables:

- `species` - The tree species (PSME = Douglas fir, PIMO = Western white pine, PIPO = Ponderosa pine).
 - `diameter` - Tree diameter at 1.3m above ground (*cm*).
 - `height` - Tree height (*m*).
 - `leafarea` - Total leaf area (*m*²).
 - `branchmass` - Total (oven-dry) mass of branches (*kg*).
1. ■ Perform a PCA of the data in the allometry dataset ('Allometry.csv') without changing the `scale` argument and plot the ordination result. Note that the first column contains species identities so you need to exclude this column from the analysis. See Section 1.2.1 for help, if necessary.
 2. ■ In the previous ordination, the length of each vector varies extensively. This is because the variance in 'branchmass' and 'leafarea' is much larger than that of 'diameter' and 'height'. Use the `var` function to calculate the variance for each variable in `allom` to confirm this. Then use the `scale` argument to repeat the analysis after standardising the variables. (See Section 1.2.1 for help, if necessary.) What effect did this have on the vector lengths?
 3. ◆ The ordination plot shows very strong overlap for most of the sites and a few very divergent sites. In addition, the plot is partly off-centre. These properties indicate that there is lack of normality in the response variables. Log-transform the variables to see the effects on the ordination, trying this two ways: (a) using `log()` to generate new vectors containing transformed data and (b) using `decostand()` with an appropriate 'method' argument). See Section 1.2.2 for help, if necessary.

1.5.1.2 Scotch data

File: 'scotch.csv'

These data are a guide to Scottish malt whiskeys as compiled and analysed by Lapointe and Legendre (1994; <https://doi.org/10.2307/2986124>). The dataframe contains 68 aesthetic attributes measured on 109 malt whiskies along with additional metadata.

Variables:

- NAME - Whisky name
 - COLOR_* - Colour attributes
 - NOSE_* - Smell attributes (binary)
 - BODY_* - Body attributes (binary)
 - PAL_* - Palatte attributes (binary)
 - FIN_* - Finish attributes (binary)
 - AGE - Whisky age, with NA values scored as -9.
 - REGION - Distillery region of origin.
 - DISTRICT - Distillery district of origin.
1. ■ Read in the scotch dataset ('scotch.csv'). The dataframe includes variables that represent assessments of taste, smell and appearance (binary response data) and other attributes (metadata) of malt whiskeys. Generate a new dataframe that includes only the binary response data. Hint: names of data variables include the _ symbol to separate variable category from the variable name while metadata variables do not include this symbol, use `grep('_', names(df))` to index the dataframe.
 2. ■ Perform PCA on the binary response data and inspect a screeplot to evaluate whether it is useful to interpret the result of the PCA.
 3. ♦ Produce a dendrogram from the binary response data, using the 'jaccard' dissimilarity index instead of euclidean distances.
 4. ♦ Perform PCoA on the binary response data, using the 'jaccard' dissimilarity index. Produce a screeplot for the non-negative eigenvalues. Calculate the proportion of variation explained by each of the non-negative eigenvalues.

1.5.1.3 Endophyte data

File: 'endophytes_env.csv' and 'endophytes.csv'

Community fingerprints from fungi colonising living leaves and litter from nine Eucalyptus spp. in the HFE common garden and surrounding area. 'endophytes.csv' contains a 'species-sample' matrix, with 98 samples in rows and 874 operational taxonomic units (OTUs) in columns.

The variables below refer to the data in 'endophytes_env.csv'. The rows are matched across the two tables.

Variables:

- species - Tree species
- type - Whether leaves came from canopy (fresh) or ground (litter)
- percentC - Leaf carbon content, per gram dry mass

- `percentN` - Leaf nitrogen content, per gram dry mass
 - `CNratio` - Ratio of C to N in leaves
1. ■ Read in the data from 'endophytes.csv'; see Section ?? (p. ??) for a description of the data. Use the `decorana` function to calculate gradient lengths and determine whether PCA is appropriate for these data (see Section 1.2.2.1 for help, if necessary).
 2. ♦ Perform CA these data and plot the results. Notice the strong skew in the data along both axes. Try again after standardising the community matrix using the `decostand` function (try the 'hellinger' and 'max' methods). Notice the undesirable parabolic pattern in the ordination and strong skew; this suggests that CA is not an improvement over PCA (common for data matrices that contain many zeros, collected along long environmental gradients).
 3. ♦ Perform PCoA on these data, using the 'hellinger' method for the `decostand` function and 'bray' method for the `vegdist()` function, and plot the results. See Section 1.2.3 for help, if necessary. Repeat as before but use the `binary` argument in the `vegdist` function to convert the matrix to a 'presence/absence' matrix.

1.5.2 Analysis of Structure 1: two-table analysis

1.5.2.1 Scotch data

1. ♦ Read in the scotch dataset ('scotch.csv') and generate a separate dataframe for each of the five types of binary response data (colour, nose, body, palate and finish). (Hint: use the `grep` function to index the dataframe for each category.) Calculate Mantel correlation statistics and their statistical significance for each pair of response types. Which types exhibited statistically significant similarities in their responses, and what pair had the strongest similarity? Which type was most unique in its response?

1.5.2.2 Endophyte data

1. ♦ Look at the help page for the `capscale` function. Use `capscale` to perform distance-based RDA (constrained PCoA) using the continuous variables in 'endophytes_env.csv' (`percentC`, `percentN`, `CNratio`) as predictors, then plot the results. First use the `envfit` function to determine which variables to include in db-RDA (Section 1.3.2).
2. ♦ Repeat the analysis in the previous exercise but use the `ordistep` function to determine which variables to include in db-RDA.

1.5.3 Analysis of Structure 2: variation partitioning

1.5.3.1 Endophyte data

1. ♦ Perform variation partitioning to determine whether leaf species, leaf chemistry, or sample type explains the most variation in fungal community composition.
2. ♦ Use the geographic coordinates of each plot to estimate the contribution of space to variation in fungal community composition. Is this estimate greater than the variation partitioned to the measured leaf variables?
3. ♦ Test the significance of each individual partition.

-
4. ♦ Generate dummy variables (using 'dudi.hillsmith') for each of the levels of 'species' and check whether there are particular leaf species that explain variation in fungal community composition.

1.5.4 Analysis of Structure 3: 'experimental' systems

1.5.4.1 Allometry data

1. ♦ For the allometry data, plot a dendrogram of multivariate distances (euclidean) among individual trees based on the four growth parameters, labelling the tips of the dendrogram with the species level. Use ANOSIM and PERMANOVA to test the hypothesis that clusters can be explained by interspecific variation. See Section 1.3.5.1 for help, if necessary.
2. ♦ Using your knowledge from Chapter ?? and Sections 1.2.4 and 1.3.5.1, plot the ordination results using coloured circles to represent the different tree species and include a legend.
3. ▲ Overlay the plot with the centroid (i.e., average) for each species, using a different symbol than for the individual points. Modify the axes to reflect the percentage of inertia (i.e., variance) explained by each axis. Refer to Chapter ?? for help tabulating mean values, if necessary.

1.5.4.2 Scotch data

1. ♦ For the scotch data, produce a dendrogram from the binary response data using the jaccard dissimilarity index. Label the tips of the dendrogram with the district from which the malt whisky originated.
2. ♦ For the scotch data, use ANOSIM and PERMANOVA to test the hypothesis that clusters can be explained by district See Section 1.3.5.1 for help, if necessary.

1.5.4.3 Endophyte data

1. ♦ Use the `adonis` function to test for main and interactive effects of tree species and tissue type (fresh vs litter) on fungal community composition (see Section 1.3.5.1). The predictor variables can be found in 'endophytes_env.csv'. What terms were significant? Which term explained the most variation?
2. ▲ Plot the PCoA results and use different symbols and colours to reflect the identities of tree species and tissue types. Add a legend to the plot. Use information from Chapter ?? and Sections 1.2.4 and 1.3.5.1 for help, if necessary. Hint: automatic functions for generating vectors of colours, such as `rainbow`, can lead to very similar colours with so many treatment levels. Check out the `randomcolor` package for alternative approaches.
3. ▲ Overlay the plot with ellipses representing 95% confidence intervals for each species and sample type using functions seen in section 1.3.5.1.

1.5.4.4 Mite data

1. ■ Use `manova` to estimate the responses of mite community composition to the environmental variables associated with the `mite` data.
2. ♦ There is not a built in function to view diagnostic plots for `manova` output. Use the `resid` function to obtain the residuals for each response variable and use the `qqnorm` and `qqline` functions to produce quantile-quantile plots for a few of the response variables to determine whether it is appropriate to model the responses using a normal error distribution.

-
3. ♦ Use `manyglm` to estimate the responses of mite community composition to the environmental variables associated with the `mite` data. Which error family, poisson or negative binomial, provides the best fit to the data? Look at the results of the best fitting model.

Index

ade4, 2, 30

adonis, 45

capscale, 30

cca, 45

decorana, 45

drop1, 44

dudi.mix, 30

envfit, 45

grep, 48

grep('_', names(df)), 47

hclust, 8, 45

indicspecies, 32

indval, 32

labdsv, 2, 32

manova, 45

mantel, 45

manyglm, 45

metaMDS, 45

mvabund, 45

mvabund, 41, 45

ordiellipse, 37

ordihull, 37

ordispider, 37

ordistep, 30, 45

ordistep(), 29

pcnm, 45

rainbow, 49

randomcoloR, 49

rankindex, 12, 45

rda, 45

rda(), 29

scores, 45

summary, 6, 45

varpart, 21, 45

vegan, 2, 4, 8, 10, 13, 16, 18–21, 26, 35, 37, 45

vegdist, 45

wcmdscale, 45