

Outline

Multivariate Analysis in Ecology

I: Unconstrained Ordination

Jari Oksanen

Oulu

January 2009

Multivariate Analysis and Ordination

- Basic ordination methods to simplify multivariate data into low dimensional graphics
- Analysis of multivariate dependence and hypotheses
- Analyses can be performed in **R** statistical software using **vegan** package and allies
- Course homepage
<http://cc.oulu.fi/~jarioksa/opetus/metodi/>
- **Vegan** homepage <http://vegan.r-forge.r-project.org/>

1 Introduction

- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Outline

1 Introduction

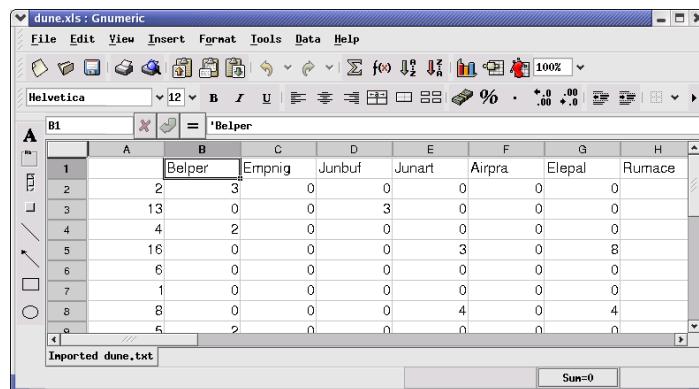
- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

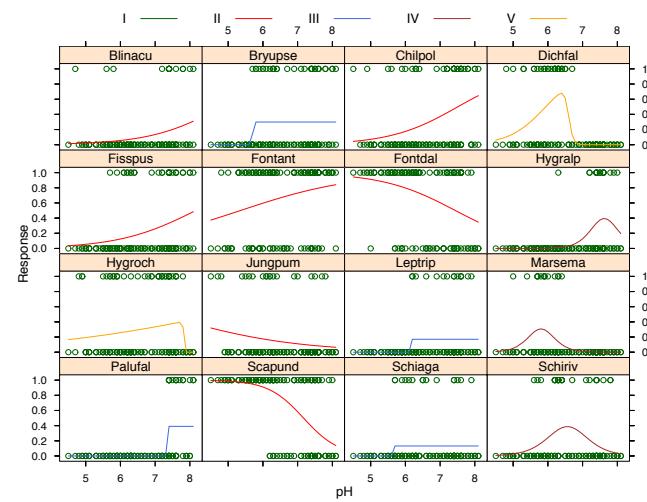
- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Why Ordination?

- **Nobody** should want to make an ordination, but they are desperate with multivariate data
- Map multidimensional table into low-dimensional display

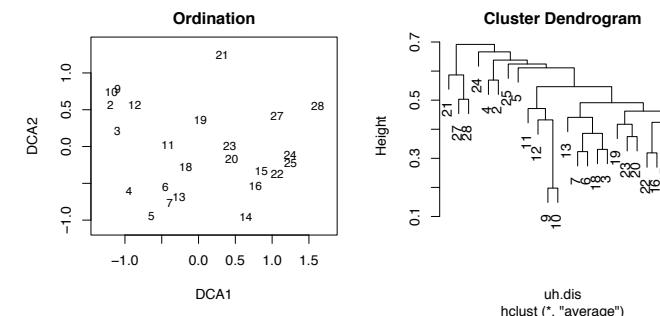


Two Ways of Analysing Data

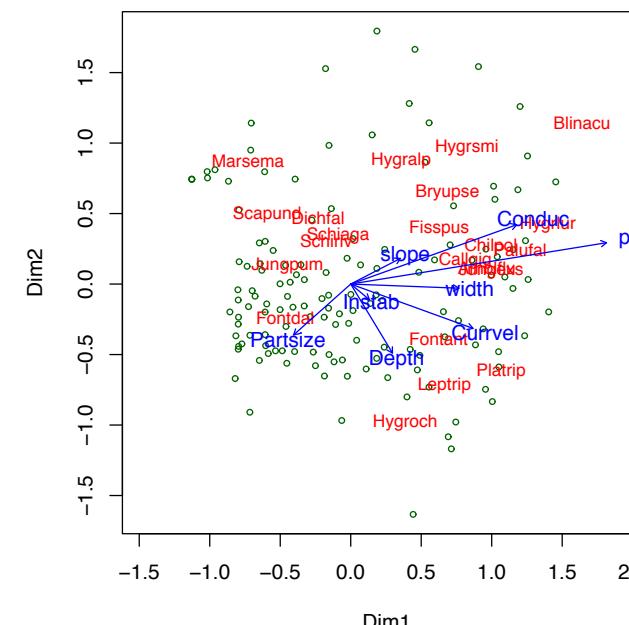


Why Ordination?

- **Nobody** should want to make an ordination, but they are desperate with multivariate data
- Map multidimensional table into low-dimensional display



Two Ways of Analysing Data



Outline

1 Introduction

- What is Ordination?
 - Gradient Analysis

2 Unconstrained Ordination

- NMDS
 - Eigenvector Methods
 - PCA
 - CA
 - Graphics
 - Environmental Variables
 - Gradient Model and Ordination

The Gradient Model

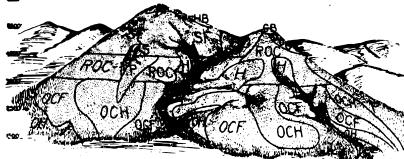


FIG. 21. Topographic disposition of vegetation types. View of idealized mountain and valley, looking east, with 6500-ft peak bearing subalpine forest on left, lower 5500-ft peak covered up to summit bald with deciduous forest on right. Vegetation types:

BG—Beech Gap	OH—Oak-Hickory Forest
CF—Cove Forest	P—Pine Forest and Pine
F—Fraser Fir Forest	Heath
GB—Grassy Bald	ROC—Red Oak-Chestnut
H—Hemlock Forest	Forest
HB—Heath Bald	S—Spruce Forest
OCF—Chestnut Oak-	SF—Spruce-Fir Forest
Chestnut Forest	WOC—White Oak-Chestnut
OCH—Chestnut Oak- Chionanthus Heath	Forest

R. H. Whittaker (1956) Vegetation of The Great Smoky Mountains. *Ecological Monographs* **26**, 1-80.

- Gradient Analysis developed in 1950s in USA, with R. H. Whittaker as the main founding father
 - Only two or three environmental variables, or *Gradients* needed to explain complicated community patterns
 - Against classification: Species responses smooth along gradients
 - Against organism analogies: Species responses individualistic
 - The basis of modern theory and praxis: Ordination and Gradient modelling of communities

The Gradient Model

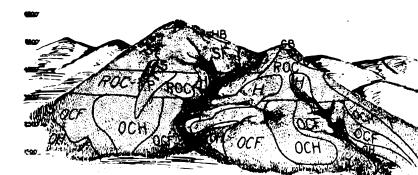
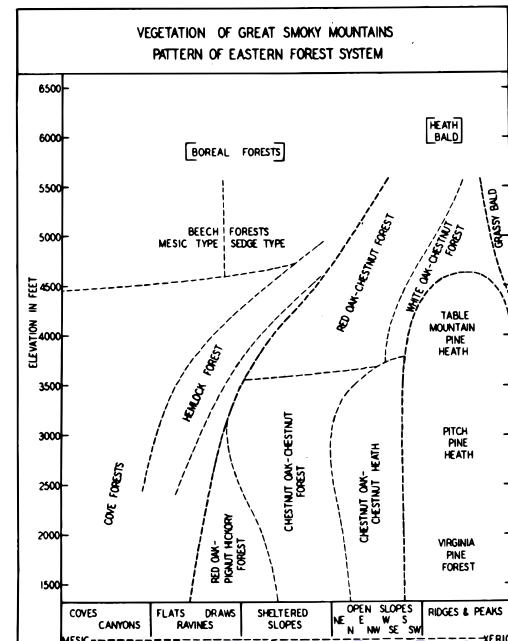


FIG. 21. Topographic disposition of vegetation types. View of idealized mountain and valley, looking east, with 6500-ft peak bearing subalpine forest on left, lower 5500-ft peak covered up to summit bald with deciduous forest on right. Vegetation types:

BG—Beech Gap	OH—Oak-Hickory Forest
CF—Cove Forest	P—Pine Forest and Pine
F—Fraser Fir Forest	Heath
GB—Grassy Bald	ROC—Red Oak-Chestnut
H—Hemlock Forest	Forest
HB—Heath Bald	S—Spruce Forest
OCF—Chestnut Oak	SF—Spruce-Fir Forest
Chestnut Forest	WOC—White Oak-Chestnut
OCH—Chestnut Oak	Forest
Chestnut Heath	



The Gradient Model

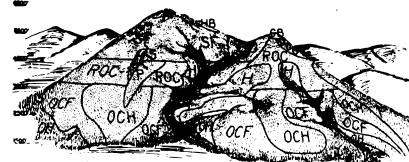


FIG. 21. Topographic disposition of vegetation types. View of idealized mountain and valley, looking east, with 6500-ft peak bearing subalpine forest on left, lower 5500-ft peak covered up to summit bald with deciduous forest on right. Vegetation types:

BG—Beech Gap	OH—Oak-Hickory Forest
CF—Cove Forest	P—Pine Forest and Pine
F—Fraser Fir Forest	Heath
GB—Grassy Bald	ROC—Red Oak-Chestnut Forest
H—Hemlock Forest	S—Spruce Forest
HB—Heath Bald	SF—Spruce-Fir Forest
OCF—Chestnut Oak-Chestnut Forest	WOC—White Oak-Chestnut Forest
OCH—Chestnut Oak-Chestnut Heath	

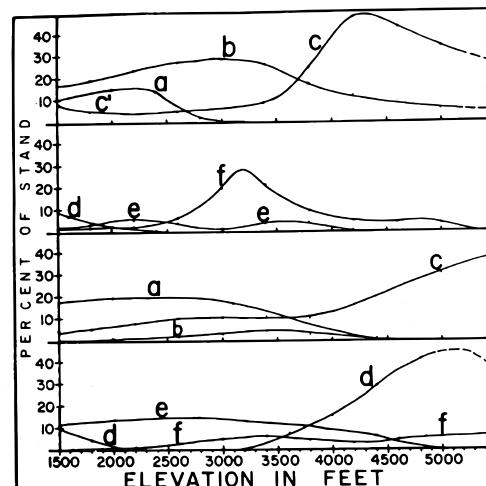


FIG. 9. Elevation transects in submesic and subxeric sites, smoothed curves for tree species. Above—submesic sites: a, *Cornus florida*; b, *Acer rubrum*; c and c', *Quercus borealis* and var. *maxima*; d, *Carya tomentosa*; e, *Carya glabra*; f, *Hamamelis virginiana*. Below—subxeric sites: a, *Quercus prinus*; b, *Sassafras albidum*; c, *Castanea dentata*; d, *Quercus alba*; e, *Oxydendrum arboreum*; f, *Robinia pseudoacacia*.

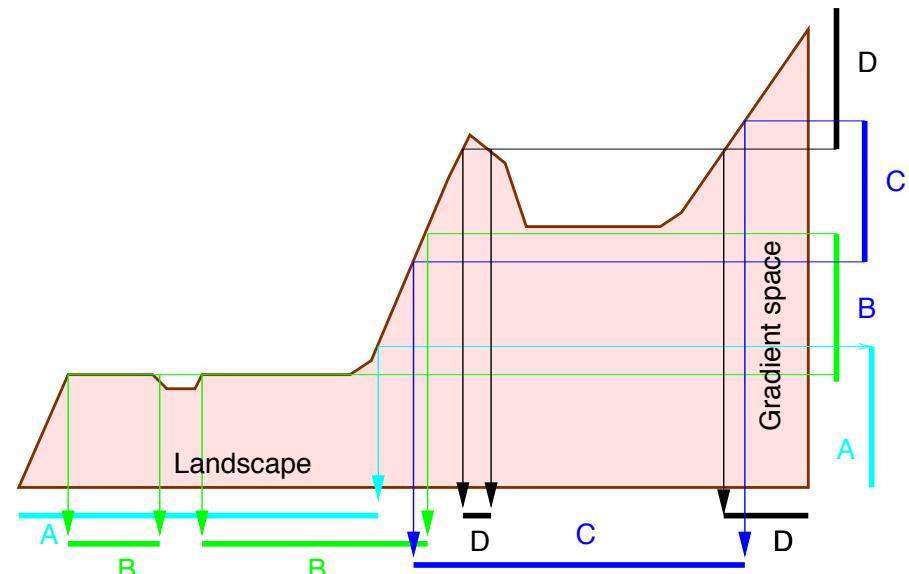
Types of Gradients

- ① **Direct gradients:** Influence organisms but are not consumed.
 - Correspond to conditions.
- ② **Resource gradients:** Consumed
 - Correspond to resources.
- ③ **Complex gradients.** Covarying direct and/or resource gradients:
Impossible to separate effects of single gradients.
 - Most observed gradients.

Types of Gradients

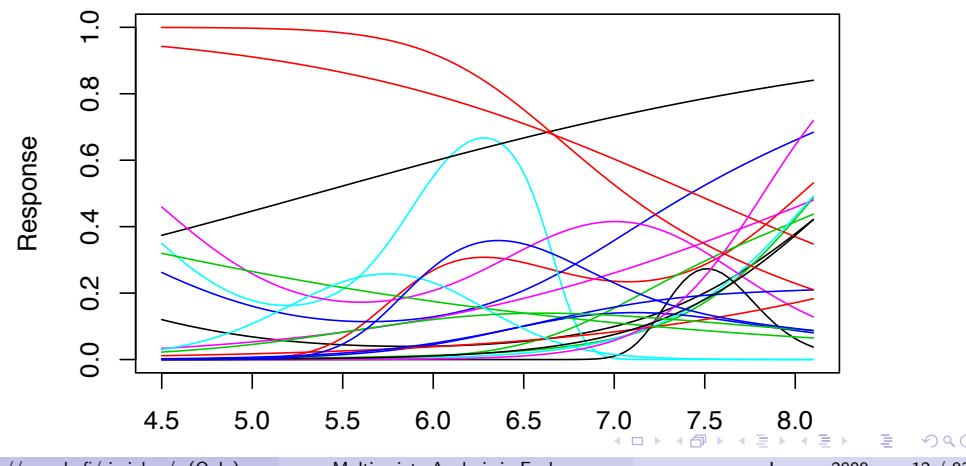
- ① **Direct gradients:** Influence organisms but are not consumed.
 - Correspond to conditions.
- ② **Resource gradients:** Consumed
 - Correspond to resources.

Landscapes and Gradients



Species responses

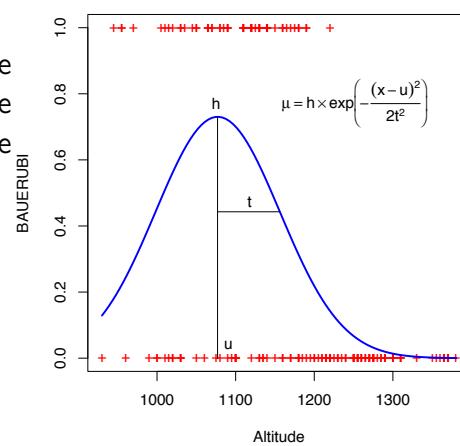
- Species have non-linear responses along gradients.



Gaussian Response Function

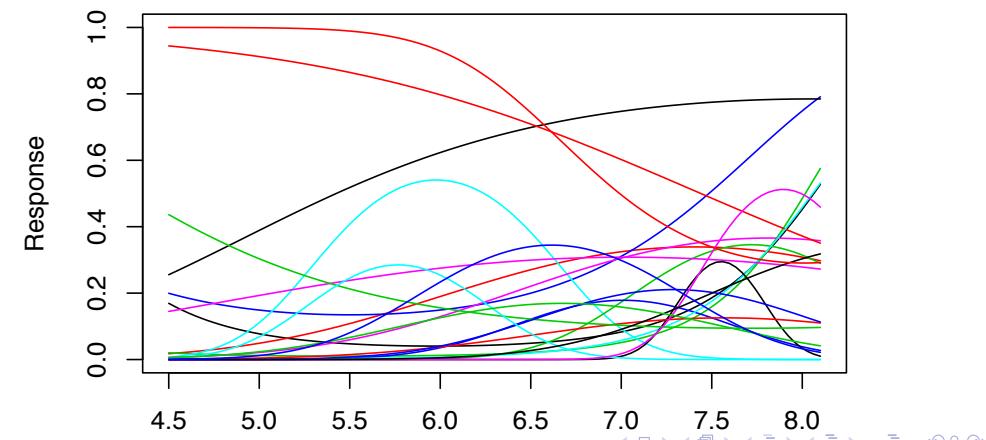
Gaussian Response Function has three interpretable parameters that define the expected response μ along the gradient x

- Location of the optimum u
- Width of the response t
- Height of the response h



Species responses

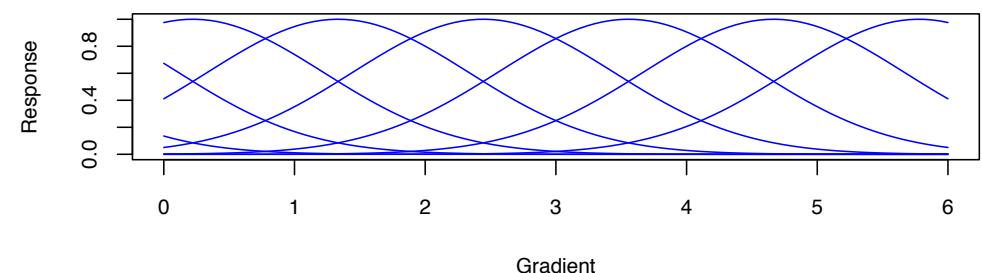
- Species have non-linear responses along gradients.
- Often assumed to be Gaussian...



Dream of species packing

Species have Gaussian responses and divide the gradient optimally:

- Equal heights h .
- Equal widths t .
- Evenly distributed optima u .



Evidence for Gaussian Responses

- Whittaker reported a large number of different response types
- Only a small proportion were symmetric, bell shaped responses
- Still became the standard of our times

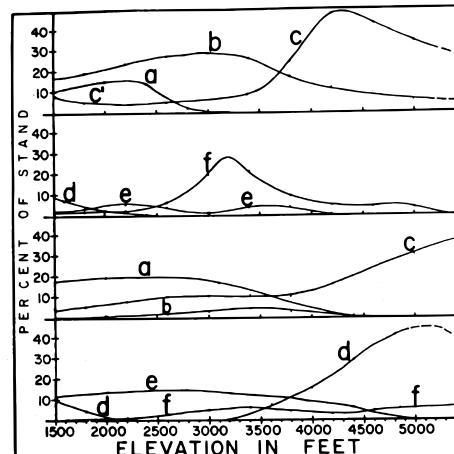


FIG. 9. Elevation transects in submesic and subxeric sites, smoothed curves for tree species. Above—submesic sites: a, *Cornus florida*; b, *Acer rubrum*; c and c', *Quercus borealis* and var. *maxima*; d, *Carya tomentosa*; e, *Carya glabra*; f, *Hamamelis virginiana*. Below—subxeric sites: a, *Quercus prinus*; b, *Sassafras albidum*; c, *Castanea dentata*; d, *Quercus alba*; e, *Oxydendrum arboreum*; f, *Robinia pseudoacacia*.

Outline

1 Introduction

- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

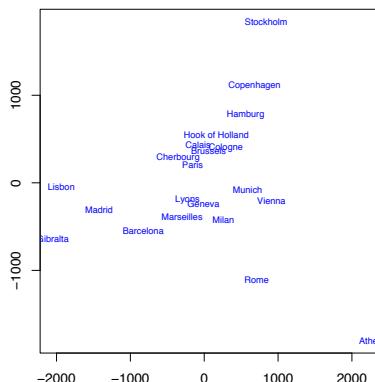
- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Ordination

- Ordination maps multivariate data onto low dimensional displays: "Most data sets have 2.5 dimensions"
- Basic ordination uses only community composition: *Indirect Gradient Analysis*
- Constrained ordination studies only the variation that can be explained by the available environmental variables: Often called *Direct Gradient Analysis*
- Distinct flavours of tools:
 - Nonmetric MDS the most robust method
 - PCA duly despised
 - Flavours of Correspondence Analysis popular
 - Canonical method: Constrained Correspondence Analysis

MDS is a map

- MDS tries to draw a map using distance data.
- MDS tries to find an underlying configuration from dissimilarities.
- Only the configuration counts:
 - No origin, but only the constellations.
 - No axes or natural directions, but only a framework for points.



Map of Europe from road distances.

metaMDS I

```
> vare.mds <- metaMDS(varespec)

Square root transformation
Wisconsin double standardization
Run 0 stress 18.4
Run 1 stress 20.7
Run 2 stress 20.6
Run 3 stress 20.5
Run 4 stress 18.3
... New best solution
... procrustes: rmse 0.0451 max resid 0.169
Run 5 stress 21.1
Run 6 stress 20.5
Run 7 stress 18.4
Run 8 stress 19.7
Run 9 stress 23.0
Run 10 stress 19.8
Run 11 stress 20.9
Run 12 stress 19.5
```

Recommended procedure

NMDS may be good, but its use needs special care: Not every NMDS automatically is good

- ① Use adequate dissimilarity indices: An adequate index gives a good rank-order relation between community dissimilarity and gradient distance.
- ② No convergence guaranteed: Start with several random starts and inspect those with lowest stress.
- ③ Satisfied only if minimum stress configurations are similar.

metaMDS II

```
Run 13 stress 21.4
Run 14 stress 21.3
Run 15 stress 25.2
Run 16 stress 21.7
Run 17 stress 20.9
Run 18 stress 20.3
Run 19 stress 18.3
... New best solution
... procrustes: rmse 5.51e-05 max resid 0.00015
*** Solution reached
```

```
> vare.mds
```

metaMDS III

```
Call:  
metaMDS(comm = varespec)
```

Nonmetric Multidimensional Scaling using isoMDS (MASS package)

Data: wisconsin(sqrt(varespec))

Distance: bray

Dimensions: 2

Stress: 18.3

Two convergent solutions found after 19 tries

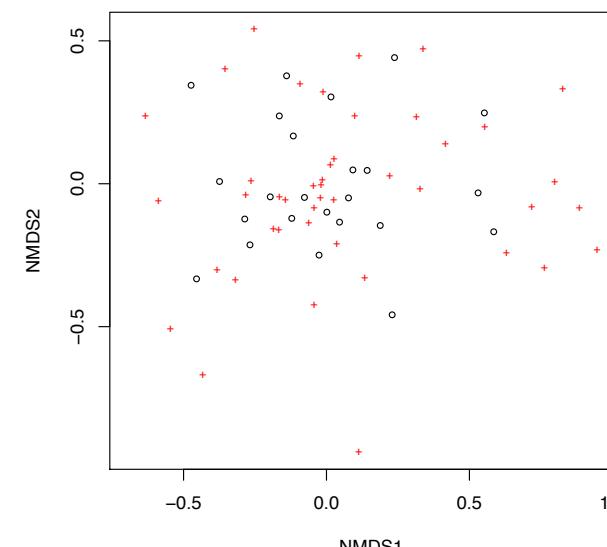
Scaling: centring, PC rotation, halfchange scaling

Species: expanded scores based on 'wisconsin(sqrt(varespec))'

Numbers

- Goodness of fit measure **stress** is based on the residuals from the non-linear regression
- Orientation, rotation, scale and origin of the coordinates (scores) are indeterminate: only the constellation matters
- Vegan arbitrarily fixes some of these:
 - Axes are centred, but the origin has no special meaning
 - Axes are rotated so that the first is the longest (technically: rotated to principal components)
 - Axes are scaled so that one unit corresponds to halving of similarity from the "replicate similarity"
 - The sign (direction) of the axes still undefined

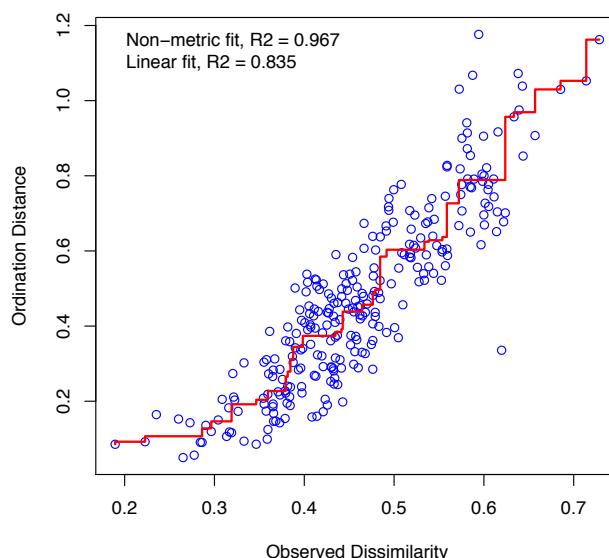
Plot metaMDS



What happened in metaMDS?

- ➊ Square root transformation and Wisconsin double standardization
- ➋ Bray–Curtis dissimilarities
- ➌ isoMDS with several random starts and stopping after finding two identical minimum stress solutions
- ➍ Solution rotated to PCs
- ➎ Solution scaled to half-change units
- ➏ Species scores as weighted averages

Shepard Diagram

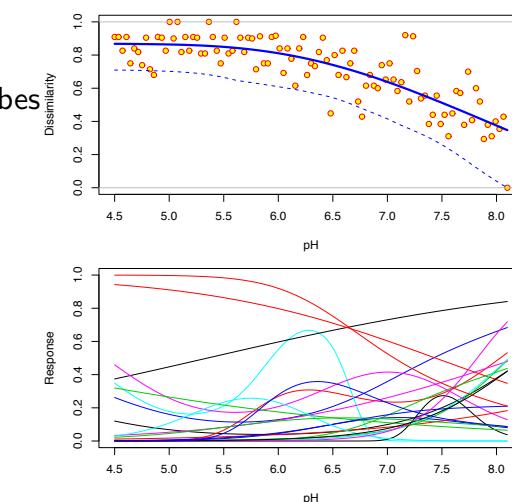


Procrustes rotation

- Procrustes rotation to maximal similarity between two configurations:
 - Translate the origin.
 - Rotate the axes.
 - Deflate or inflate the axis scale.
- Single points can move a lot, although the stress is fairly constant:
Especially in large data sets.

Dissimilarity measures

- Use a dissimilarity that describes correctly gradient separation
- Bray–Curtis (Steinhaus), Jaccard, Kulczyński
- Wisconsin double standardization often helpful
- Euclidean and Chi-square distances are poor



Procrustes Rotation

```
> tmp <- wisconsin(sqrt(varespec))
> dis <- vegdist(tmp)
> vare.mds0 <- isoMDS(dis, trace = 0)
> pro <- procrustes(vare.mds, vare.mds0)
> pro
```

Call:

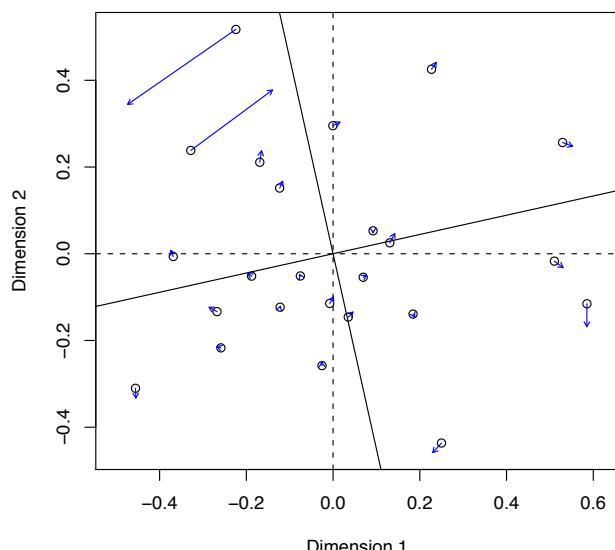
```
procrustes(X = vare.mds, Y = vare.mds0)
```

Procrustes sum of squares:

0.157

Plot Procrustes Rotations

Procrustes errors



Outline

1 Introduction

- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Number of dimensions

- In NMDS, 2D solution is not a plane in 3D space
- Solution must be found separately for each dimensionality
- Some people very disturbed: how do they *know* the correct number
- Answer is easy: there is no correct number, although some numbers may be worse than others
- "Most data sets have 2.5 dimensions"
- Typically you try with 2 and 3
- Do you need more dimensions to explain species patterns and environmental data?
- Is convergence very slow? Try another number of dimensions

Simplified mapping: Eigen analysis

- NMDS uses non-linear mapping for any dissimilarity measure: This is very difficult
- Things are much simpler if we accept only certain dissimilarity indices and map them linearly onto ordination
- It is only a rotation, and can be solved using eigenvector techniques
- Euclidean distances + rotation: Principal Components Analysis (PCA)
- Chi-square distances + weighted rotation: Correspondence Analysis (CA)
- All ordination methods are distance based

Why Not PCA?

- We admit that PCA is just a rotation, but it is a linear method

Why Not PCA?

- We admit that PCA is just a rotation, but it is a linear method
- PCA works with species space, but we boldly go to gradient space
- CA is an optimal scaling method

Why Not PCA?

<http://cc.oulu.fi/~jarioksa/> (Oulu) Multivariate Analysis in Ecology January 2009 35 / 93
Unconstrained Ordination Eigenvector Methods

- We admit that PCA is just a rotation, but it is a linear method
- PCA works with species space, but we boldly go to gradient space
- CA is an optimal scaling method
 - Sites with similar species composition packed close to each other
 - Species that occur together simultaneously packed close to each other
- CA can handle unimodal species responses, even approximate one dimensional species packing model

<http://cc.oulu.fi/~jarioksa/> (Oulu) Multivariate Analysis in Ecology January 2009 35 / 93
Unconstrained Ordination PCA

Outline

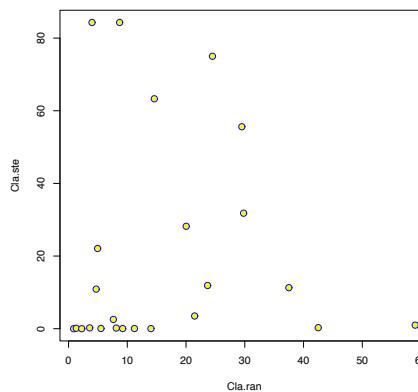
- 1 Introduction
 - What is Ordination?
 - Gradient Analysis

- 2 Unconstrained Ordination

- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

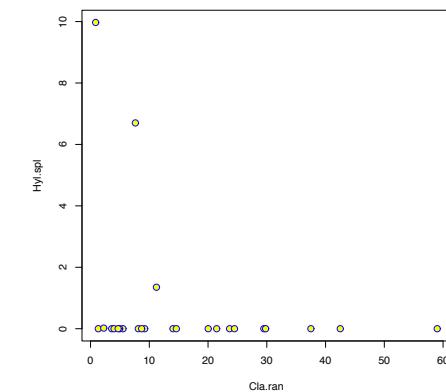
Species space

- Graphical presentations of data matrix: Species are axes and span the space where sites are points



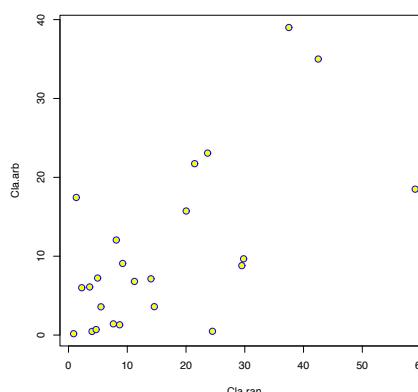
Species space

- Graphical presentations of data matrix: Species are axes and span the space where sites are points
- Some species show more of the configuration than others



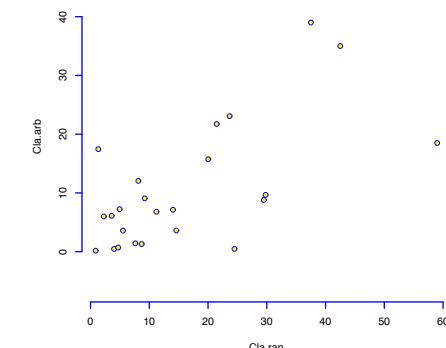
Species space

- Graphical presentations of data matrix: Species are axes and span the space where sites are points
- Some species show more of the configuration than others
- What is the ideal viewing angle to the species space?



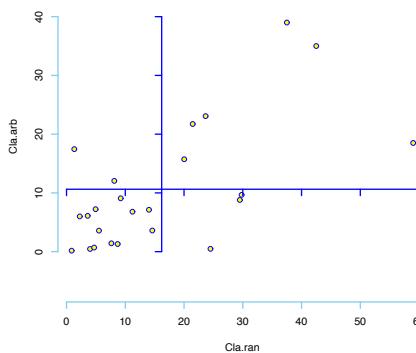
Rotation in species space

- Put sites into species space



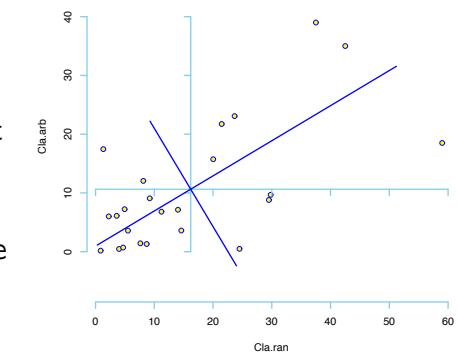
Rotation in species space

- ➊ Put sites into species space
- ➋ Move the origin to the centroid



Rotation in species space

- ➊ Put sites into species space
- ➋ Move the origin to the centroid
- ➌ Rotate the axes so that the first axis (1) is as close to all points as possible, and (2) explains as much of the variance as possible



Goodness of Fit

- ➊ The total variation is the sum of squared distances from the origin $\Lambda = \sum_j^N x_{jk}^2$: the Euclidean distance
- ➋ This can be sum of squares (SS) or variance (SS/n or $SS/(n - 1)$)
- ➌ The sum of squared distances u projected to an axis is the eigenvalue of the axis $\lambda_i = \sum_j^N u_{ij}^2$
- ➍ The eigenvalues are ordered $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$, non-negative $\lambda_i \geq 0$ and sum up to total variance $\Lambda = \sum_i^p \lambda_i$
- ➎ The proportion λ_i/Λ gives the proportion an axis explains of the total variance, and λ_1 explains the largest proportion
- ➏ PCA is often used to reduce data into a few linearly independent components that explain the most of the original variables

Running PCA I

```
> (ord <- rda(dune))
Call: rda(X = dune)

Inertia Rank
Total      84.1
Unconstrained 84.1   19
Inertia is variance

Eigenvalues for unconstrained axes:
PC1   PC2   PC3   PC4   PC5   PC6   PC7   PC8
24.80 18.15  7.63  7.15  5.70  4.33  3.20  2.78
(Showed only 8 of all 19 unconstrained eigenvalues)

> head(summary(ord), 3, 1)
```

Running PCA II

Call:
rda(X = dune)

Partitioning of variance:

	Inertia	Proportion
Total	84.1	1
Unconstrained	84.1	1

Eigenvalues, and their contribution to the variance

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
Eig.value	24.795	18.15	7.629	7.153	5.695	4.333	3.199	2.782
Accounted	0.295	0.51	0.601	0.686	0.754	0.805	0.843	0.876
	PC9	PC10	PC11	PC12	PC13	PC14	PC15	PC16
Eig.value	2.482	1.854	1.747	1.314	0.991	0.638	0.551	0.351
Accounted	0.906	0.928	0.949	0.964	0.976	0.984	0.990	0.994
	PC17	PC18	PC19					
Eig.value	0.200	0.149	0.116					
Accounted	0.997	0.999	1.000					

Running PCA IV

Site scores (weighted sums of species scores)

	PC1	PC2	PC3	PC4	PC5	PC6
2	-1.6448	-1.230	0.887	-0.986	-2.0346	-1.811
13	0.6994	-2.184	-2.213	-0.423	-1.0650	-0.656
4	0.0479	-2.046	1.274	-0.974	0.6421	0.721
....						
7	-1.7926	0.322	-0.220	1.471	-0.0125	0.426

Running PCA III

Accumulated constrained eigenvalues
numeric(0)

Scaling 2 for species and site scores

- * Species are scaled proportional to eigenvalues
- * Sites are unscaled: weighted dispersion equal on all dimensions
- * General scaling constant of scores: 6.32

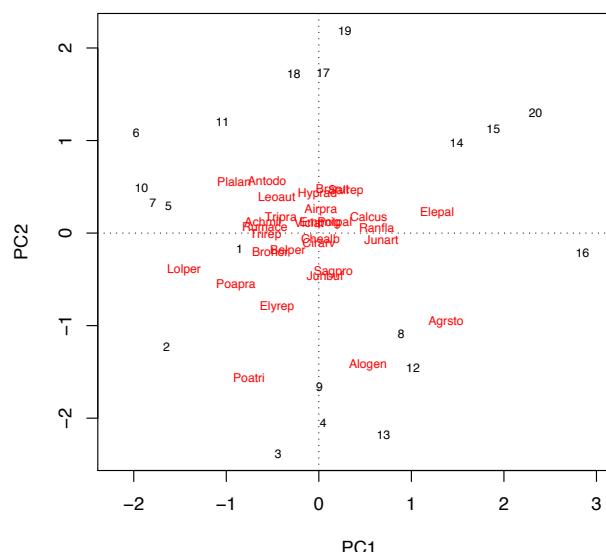
Species scores

	PC1	PC2	PC3	PC4	PC5	PC6
Belper	-0.3336	-0.189	0.1406	-0.08418	-0.1254	-0.13477
Empnig	0.0141	0.110	-0.0994	-0.16179	0.0229	0.00120
Junbuf	0.0656	-0.460	-0.5489	-0.01890	0.1057	0.08717
....						
Brohor	-0.5235	-0.197	0.1642	0.00567	-0.3861	-0.25763

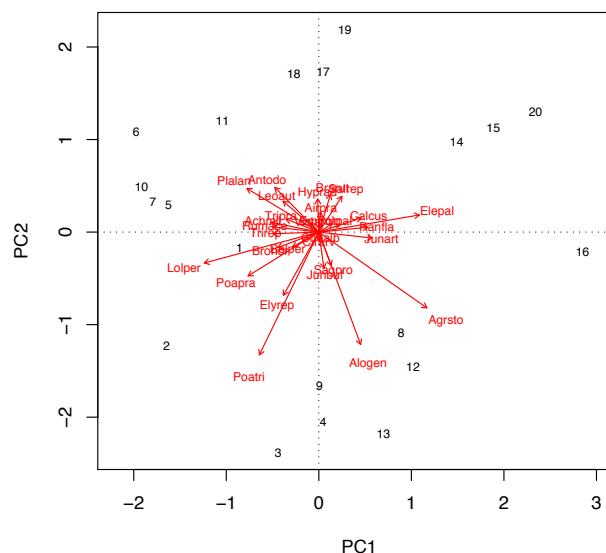
Row and Column scores

- The scores are either normalized $\sum_j^N u_{ij}^2 = 1$ or proportional to eigenvalues $\sum_j^N u_{jk}^2 \propto \lambda_j$
- Normalized scores give the regression coefficients between the axis and the variables: often used for species
- Scores proportional to the eigenvalue give the true configuration of points in the space defined by normalized scores: often used for sites (hence in species space)
- Together these scores give a linear least square approximation of the data
- Graphical presentation called **biplot**
- However, there are many alternative scaling systems

Default Plot



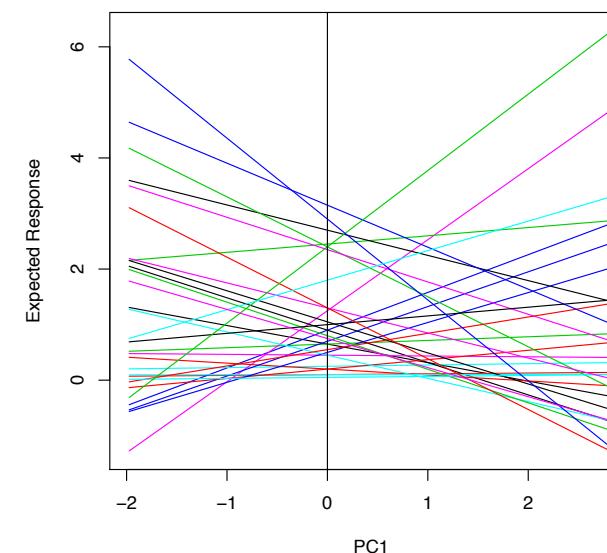
Arrow Biplot



Reading the Plot

- Origin: all species (variables) at their average values
- The *distance* from the origin for a row (site) implies how much the point differs from the average
- The *distance* from the origin for a column (species, variable) implies how much the point increases to that *direction*
- The change is measured in absolute scale: big changes, long distances from the origin
- Implies a linear model of species response against axes
- The *angle* between two points implies correlations
- 90° means zero correlation, $< 90^\circ$ positive correlation, $> 90^\circ$ negative correlation, 0° implies $r = 1$
- Arrow biplots often used instead of point biplot

Linear Model



Variances and Correlations

- Analysis of raw data explains variances: variables with high variance are most important
- If the variables are standardized to unit variance before analysis $z = (x - \bar{x})/s_x$ all variables are equally important and the analysis explains correlations among variables
- Standardization can be used when we want all variables have equal weights
- Standardization must be used when variables are measured in different scales, such as for environmental measurements

The Number of Components

- PCA is a rotation in species (character) space and retains the original configuration
- The number of PC's is $\min(n, m)$, and all together give the original data
- First axes are most important and we may ignore the minor axes
- We can either use the axes as variables in other models, or use them to identify major (almost) independent variables
- Often we want to retain a certain proportion of the variance, say 50 %
- Sometimes we would like to retain "significant" axes
- There really is no way of doing this, but some people suggest comparing eigenvalues against *broken stick* distribution

Reducing the Number of Correlated Environmental Variables I

```
> (pc <- rda(varechem, scale = TRUE))
```

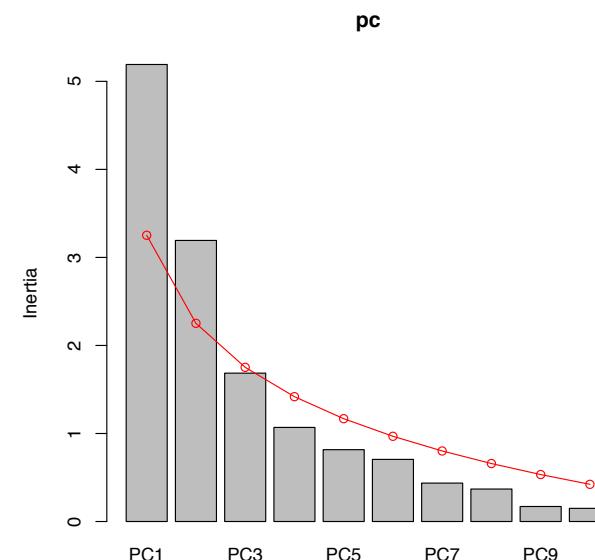
```
Call: rda(X = varechem, scale = TRUE)
```

	Inertia	Rank
Total	14	
Unconstrained	14	14
Inertia is correlations		

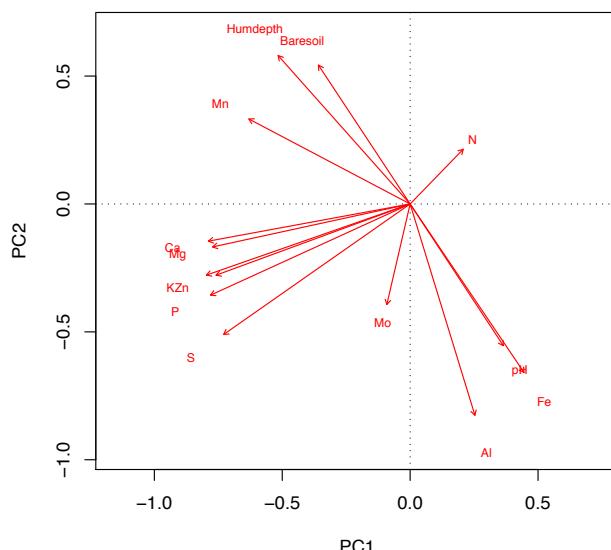
Eigenvalues for unconstrained axes:

PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9
5.1916	3.1928	1.6855	1.0690	0.8160	0.7058	0.4364	0.3688	0.1707
PC10	PC11	PC12	PC13	PC14				
0.1495	0.0853	0.0699	0.0351	0.0236				

Broken Stick and Eigenvalues



Two Dimensions, but which?



Outline

1 Introduction

- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Methods Related to PCA

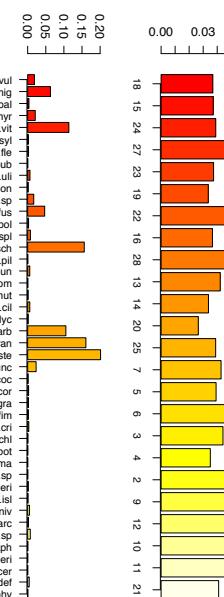
• Metric Scaling a.k.a. Principal Coordinates Analysis

- Used dissimilarities instead of raw data
- With Euclidean distances equal to PCA, but can use other dissimilarities

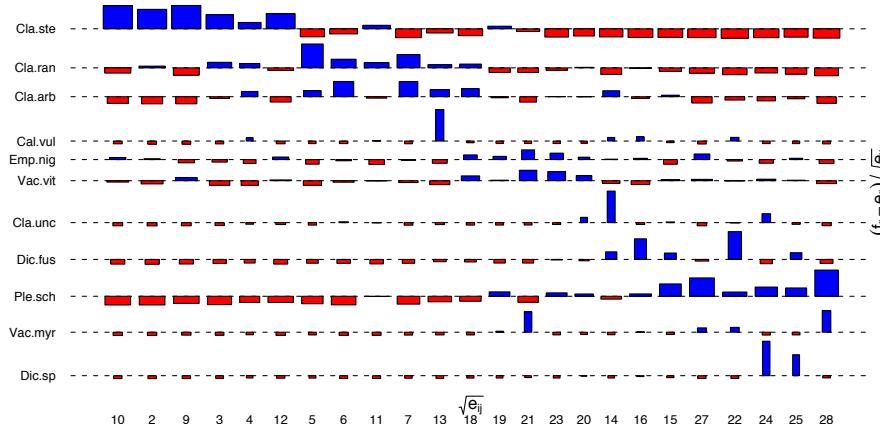
• Factor Analysis

- A *statistical* method that makes a difference between systematic components and random error
- In PCA we just ignore latter components, but here we really identify the real components

Correspondence Analysis

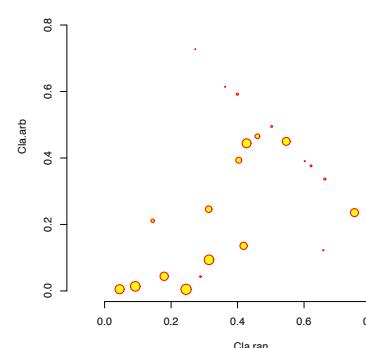


Chi-squared metric



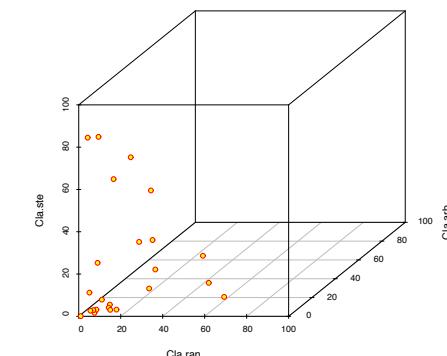
CA Rotation

- ① Sites in a species space
- ② Relative proportions are axes and points have weights



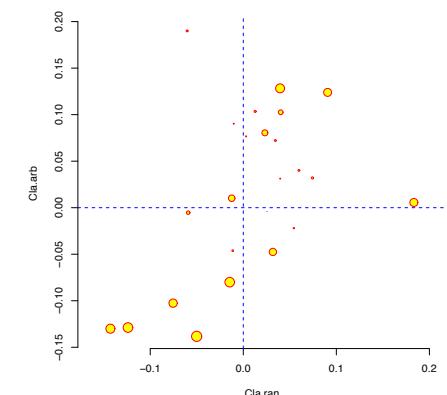
CA Rotation

- ① Sites in a species space



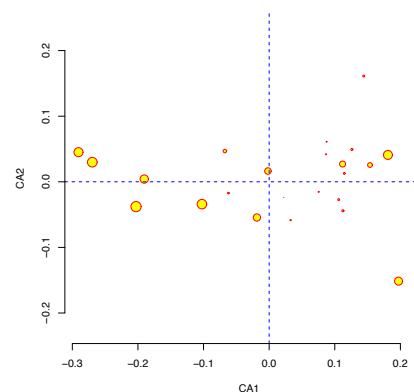
CA Rotation

- ① Sites in a species space
- ② Relative proportions are axes and points have weights
- ③ Chi-square transformation



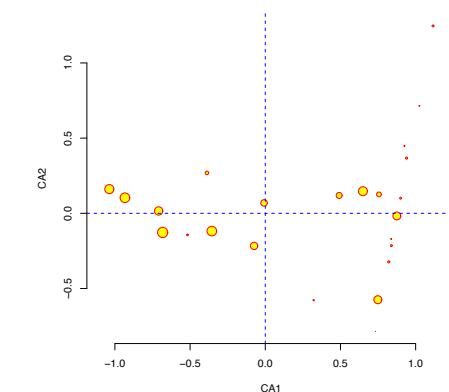
CA Rotation

- ① Sites in a species space
- ② Relative proportions are axes and points have weights
- ③ Chi-square transformation
- ④ Weighted rotation



CA Rotation

- ① Sites in a species space
- ② Relative proportions are axes and points have weights
- ③ Chi-square transformation
- ④ Weighted rotation
- ⑤ De-weighting



Running CA I

```
> (ord <- cca(dune))

Call: cca(X = dune)

Inertia Rank
Total      2.12
Unconstrained 2.12 19
Inertia is mean squared contingency coefficient

Eigenvalues for unconstrained axes:
  CA1   CA2   CA3   CA4   CA5   CA6   CA7   CA8
0.5360 0.4001 0.2598 0.1760 0.1448 0.1079 0.0925 0.0809
(Showed only 8 of all 19 unconstrained eigenvalues)
```

```
> head(summary(ord), 2)
```

Running CA II

Call:
cca(X = dune)

Partitioning of mean squared contingency coefficient:
Inertia Proportion

Total	2.12	1
Unconstrained	2.12	1

Eigenvalues, and their contribution to the mean squared contingency coefficient:

	CA1	CA2	CA3	CA4	CA5	CA6	CA7	CA8
Eig.value	0.536	0.400	0.260	0.176	0.145	0.108	0.0925	0.0809
Accounted	0.253	0.443	0.565	0.649	0.717	0.768	0.8118	0.8500
	CA9	CA10	CA11	CA12	CA13	CA14	CA15	
Eig.value	0.0733	0.0563	0.0483	0.0412	0.0352	0.0205	0.0149	
Accounted	0.8847	0.9113	0.9341	0.9536	0.9702	0.9800	0.9870	
	CA16	CA17	CA18	CA19				
Eig.value	0.00907	0.00794	0.007	0.00348				
Accounted	0.99129	0.99505	0.998	1.00000				

Running CA III

Accumulated constrained eigenvalues
numeric(0)

Scaling 2 for species and site scores

* Species are scaled proportional to eigenvalues

* Sites are unscaled: weighted dispersion equal on all dimensions

Species scores

	CA1	CA2	CA3	CA4	CA5	CA6
Belper	-0.50	-0.355	-0.152	-0.704	-0.0585	-0.0731
Empnig	-0.69	3.264	1.957	-0.177	-0.0735	0.1608
...						

Site scores (weighted averages of species scores)

Goodness of Fit of Scores

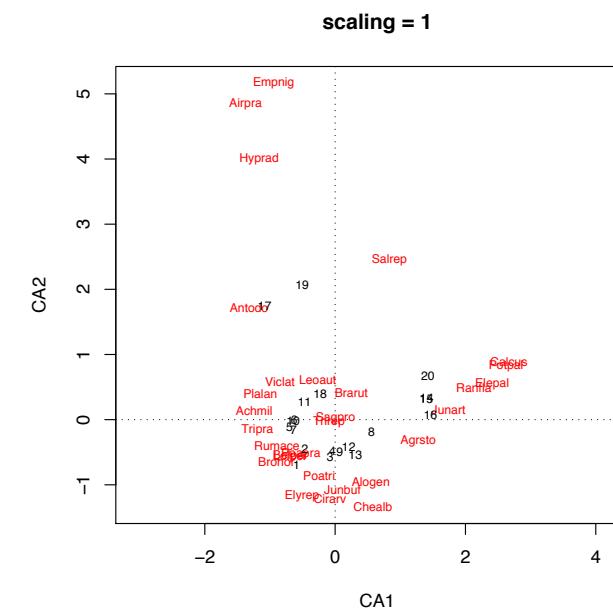
- Inertia is “mean square contingency coefficient”: Chi-squared of a matrix standardized to unit sum, or Chi-square of $x/\sum x$
- Eigenvalues are non-negative and ordered like in PCA, but they are bound to maximum 1
- The origin gives the expected abundances for all species and all sites
- The deviant species and deviant species are further away from the origin
- CA is weighted analysis, and the sum of weighted squared scores is the eigenvalue
- The species and site scores are (scaled) weighted averages of each other: proximity matters
- Rare species have low weights: they are further away from the origin

Running CA IV

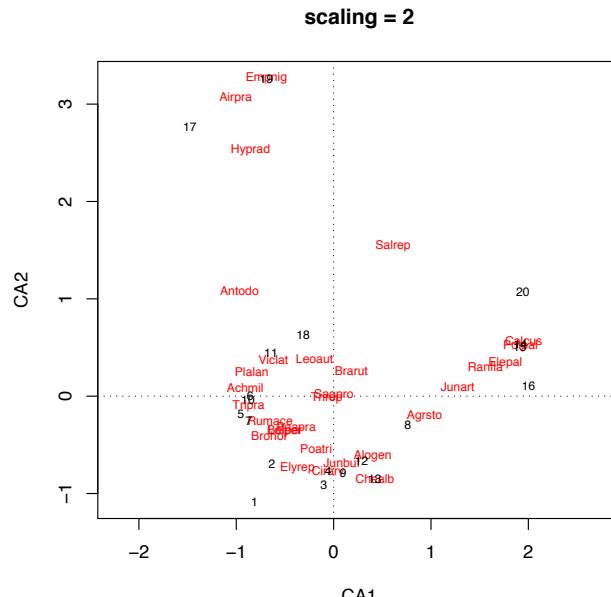
	CA1	CA2	CA3	CA4	CA5	CA6
2	-0.633	-0.696	-0.097	-1.19	-0.977	-0.0658
13	0.424	-0.844	1.590	1.25	-0.207	-0.8757
...						

....

Default Plot and Effect of Scaling

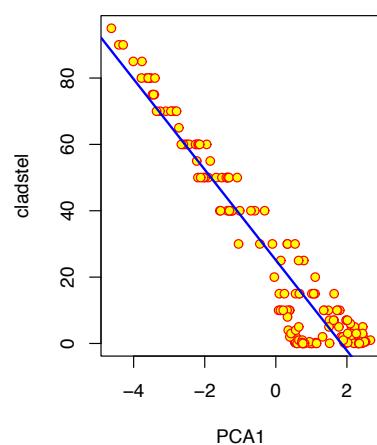


Default Plot and Effect of Scaling



Linear and Unimodal Models

- PCA implies linear relations between axes and species abundances

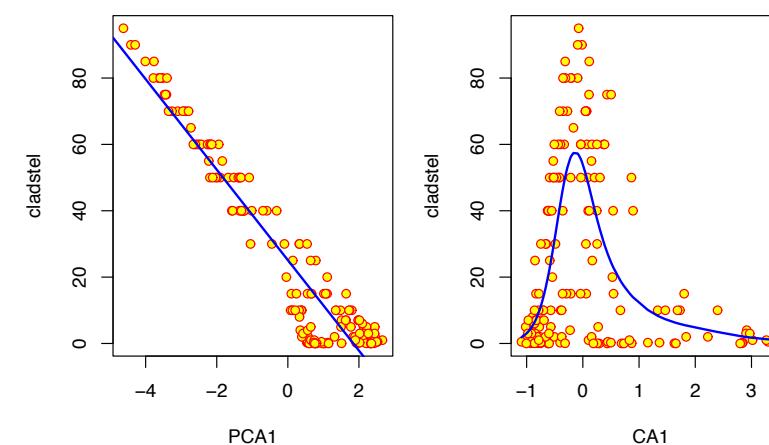


Weighted Averages

- Species scores are [proportional to] weighted averages of site scores, and simultaneously
- Site scores are [proportional to] weighted averages of species scores
- Either one (but not both) of these can be a direct weighted average of other
- If sites scores are weighted averages of species scores, site point is in the middle of points of species that occurs in the site
- The *location* of the point is meaningful whereas in PCA the main things were *distance* and *direction* from the origin (but these, too, matter)
- Can approximate unimodal response model and therefore CA is **much better** for community ordination than PCA

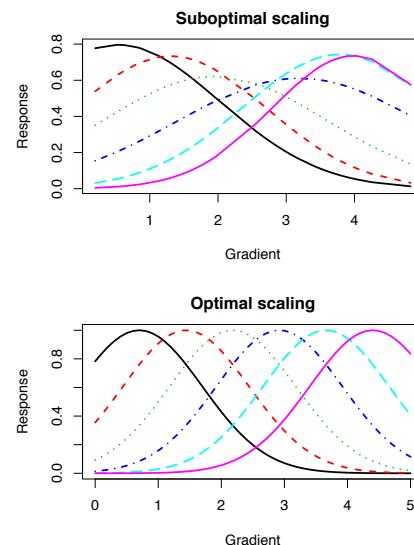
Linear and Unimodal Models

- PCA implies linear relations between axes and species abundances
- CA packs species and approximates a unimodal model



Optimal Scaling

- The locations of species optima (tops) should be widespread: measured as SS_B
- The species responses should be narrow: measured as SS_w
- The total variance is their sum $SS_T = SS_B + SS_w$
- Scaling is optimal if most of variance is between species
- The criterion variance is the eigenvalue maximized in CA: $\lambda = SS_B / SS_T$



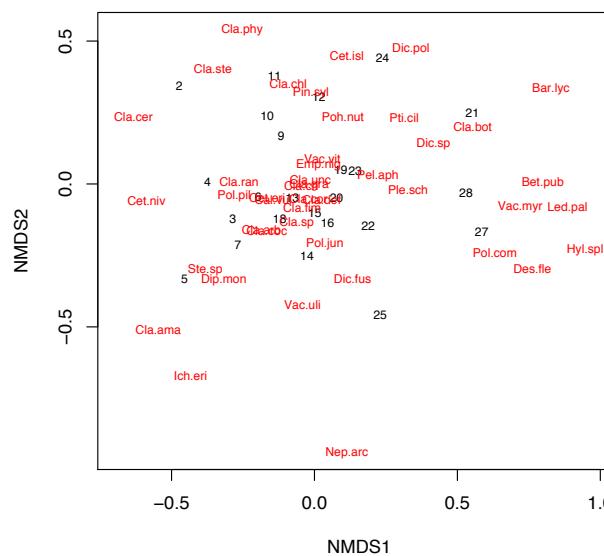
Outline

- 1 Introduction
 - What is Ordination?
 - Gradient Analysis

- 2 Unconstrained Ordination

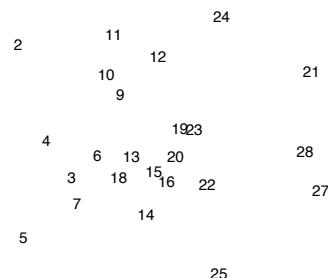
- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Anatomy of a Plot

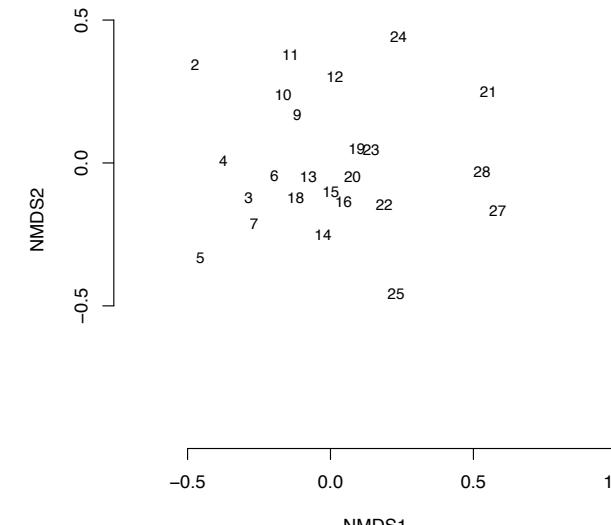


Anatomy of a Plot

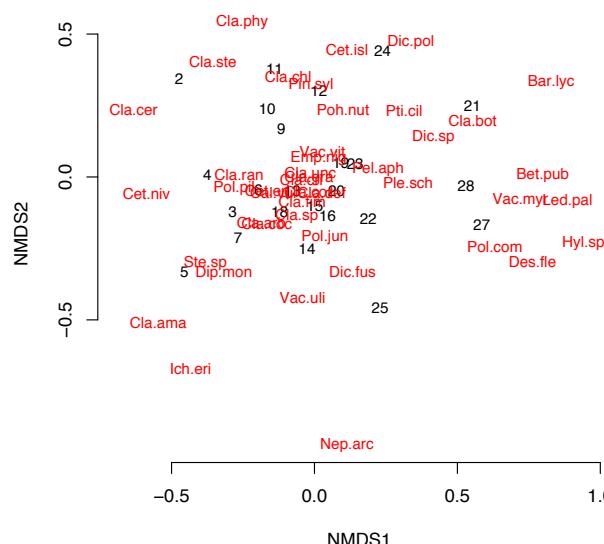
Anatomy of a Plot



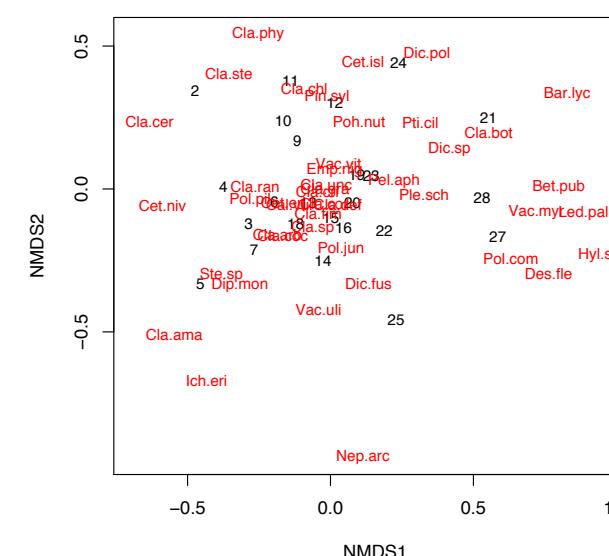
Anatomy of a Plot



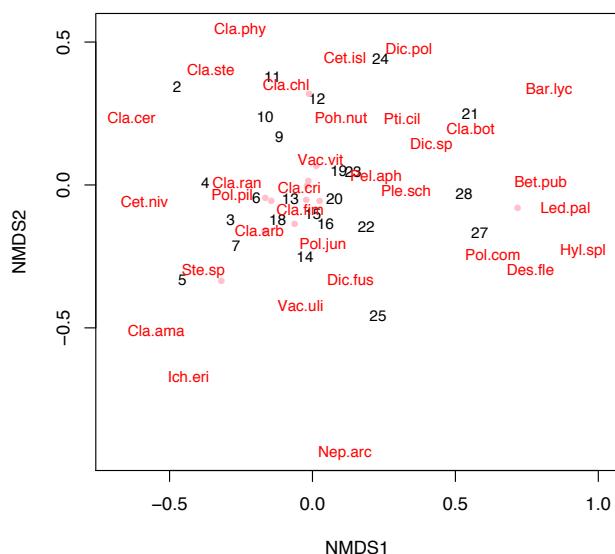
Anatomy of a Plot



Anatomy of a Plot



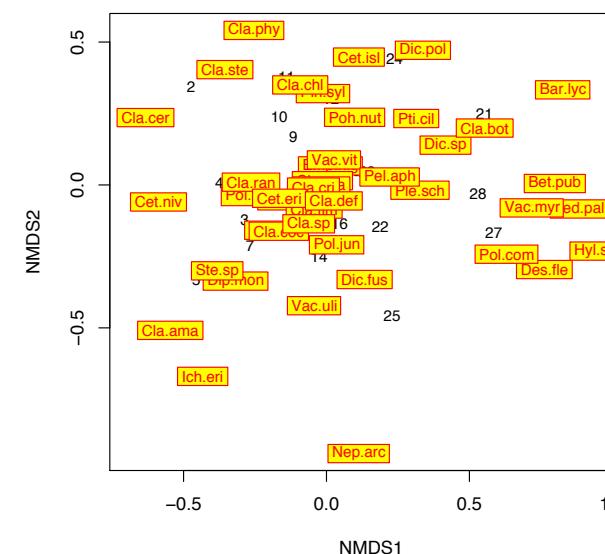
Anatomy of a Plot



Plotting functions

- All vegan ordination functions have a plot function, and ordiplot can be used for other functions as well
- For full control, use first plot(x, type="n") and then add configurable points or text
- Congested plots can be displayed with orditorp or edited with orditkplot
- Lattice graphics can be made with ordixyplot, ordicloud or ordisplom
- Dynamic, spinnable 3D plots can be made with ordirgl
- Items can be added to the plots with ordiarrows, ordihull, ordispider, ordihull, ordiellipse, ordisegments, or ordigrid

Anatomy of a Plot



Outline

- 1 Introduction
 - What is Ordination?
 - Gradient Analysis
- 2 Unconstrained Ordination
 - NMDS
 - Eigenvector Methods
 - PCA
 - CA
 - Graphics
 - Environmental Variables
 - Gradient Model and Ordination

Ordination and Environment

We take granted that vegetation is controlled by environment, so

- ① Two sites close to each other in ordination have similar vegetation

Ordination and Environment

We take granted that vegetation is controlled by environment, so

- ① Two sites close to each other in ordination have similar vegetation
- ② If two sites have similar vegetation, they have similar environment
- ③ Two sites far away from each other in ordination have dissimilar vegetation,

Ordination and Environment

We take granted that vegetation is controlled by environment, so

- ① Two sites close to each other in ordination have similar vegetation
- ② If two sites have similar vegetation, they have similar environment

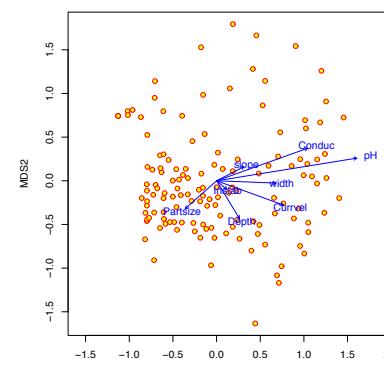
Ordination and Environment

We take granted that vegetation is controlled by environment, so

- ① Two sites close to each other in ordination have similar vegetation
- ② If two sites have similar vegetation, they have similar environment
- ③ Two sites far away from each other in ordination have dissimilar vegetation, and perhaps
- ④ If two sites have different vegetation, they have different environment

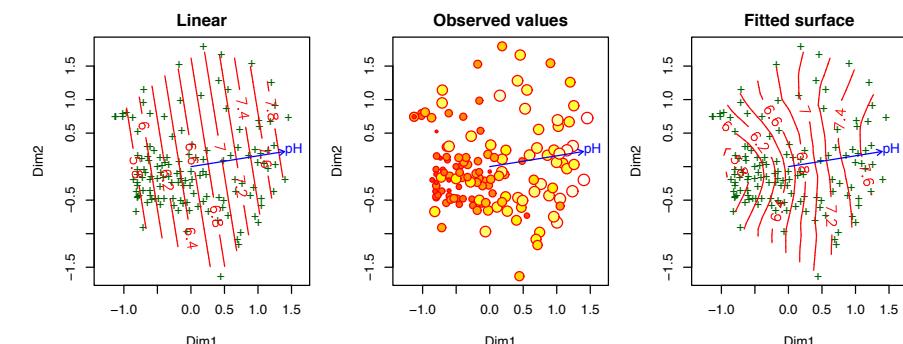
Fitted Vectors

- Direction of fitted vector shows the gradient, length shows its importance.
- For every arrow, there is an equally long arrow into opposite direction: Decreasing direction of the gradient.
- Implies a linear model: Project sample plots onto the vector for expected value.



Alternatives to Vectors

- Fitted vectors natural in constrained ordination, since these have linear constraints.
- Distant sites are different, but may be different in various ways: Environmental variables may have a non-linear relation to ordination.



Fitting Environmental Vectors I

```
> (ef <- envfit(vare.mds, varechem, permu = 1000))
```

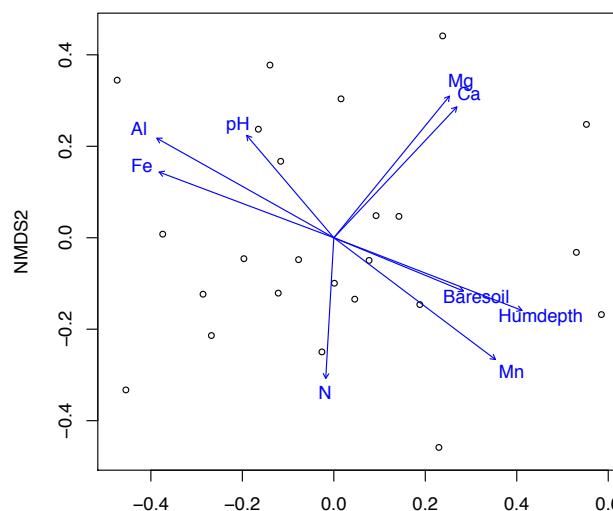
***VECTORS

	NMDS1	NMDS2	r2	Pr(>r)
N	-0.0572	-0.9984	0.25	0.056 .
P	0.6197	0.7849	0.19	0.102
K	0.7664	0.6424	0.18	0.120
Ca	0.6851	0.7284	0.41	0.009 **
Mg	0.6324	0.7746	0.43	0.002 **
S	0.1913	0.9815	0.18	0.131
Al	-0.8716	0.4901	0.53	<0.001 ***
Fe	-0.9361	0.3518	0.45	0.001 ***
Mn	0.7987	-0.6017	0.52	<0.001 ***
Zn	0.6175	0.7865	0.19	0.118
Mo	-0.9031	0.4294	0.06	0.527
Baresoil	0.9250	-0.3801	0.25	0.038 *
Humdepth	0.9329	-0.3602	0.52	<0.001 ***
pH	-0.6481	0.7616	0.23	0.054 .

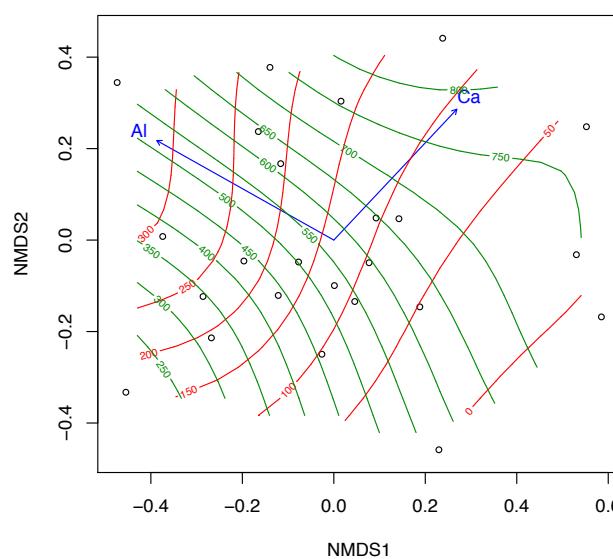
Fitting Environmental Vectors II

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1
P values based on 1000 permutations.

Plotting Environmental Vectors

Limit $p < 0.1$ 

Plotting Environmental Surfaces



Fitting Environmental surfaces

```
> ef <- envfit(vare.mds ~ Al + Ca, varechem)
> plot(vare.mds, display = "sites")
> plot(ef)
> tmp <- with(varechem, ordisurf(vare.mds, Al, add = TRUE))
```

This is mgcv 1.4-1

```
> tmp <- with(varechem, ordisurf(vare.mds, Ca, add = TRUE,
+ col = "green4"))
```

Factor Fitting I

```
> dune.ca <- cca(dune)
> ef <- envfit(dune.ca ~ A1 + Management, data = dune.env,
+ perm = 1000)
> ef
```

***VECTORS

CA1	CA2	r2	Pr(>r)
A1	0.9982	0.0606	0.31 0.055 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
P values based on 1000 permutations.

***FACTORS:

Centroids:

	CA1	CA2
ManagementBF	-0.73	-0.14
ManagementHF	-0.39	-0.30

Factor Fitting II

```
ManagementNM 0.65 1.44
ManagementSF 0.34 -0.68
```

Goodness of fit:

r^2 Pr(> r)

Management 0.44 0.001 ***

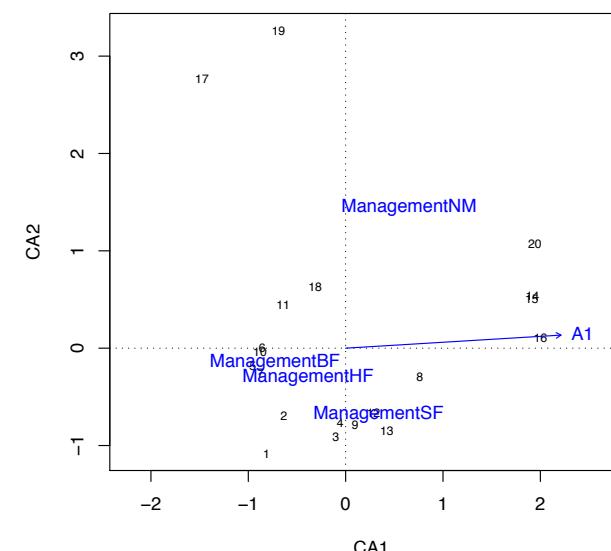
Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

P values based on 1000 permutations.

Environmental Interpretation

- Environmental variables need not be parallel to ordination axes.

Plotting Fitted Factors



Environmental Interpretation

- Environmental variables need not be parallel to ordination axes.
- Axes cannot be taken as gradients, but gradients are oblique to axes:
You cannot tear off an axis from an ordination.

Environmental Interpretation

- Environmental variables need not be parallel to ordination axes.
- Axes cannot be taken as gradients, but gradients are oblique to axes:
You cannot tear off an axis from an ordination.
- **Never** calculate a correlation between an axis and an environmental variable.

Outline

1 Introduction

- What is Ordination?
- Gradient Analysis

2 Unconstrained Ordination

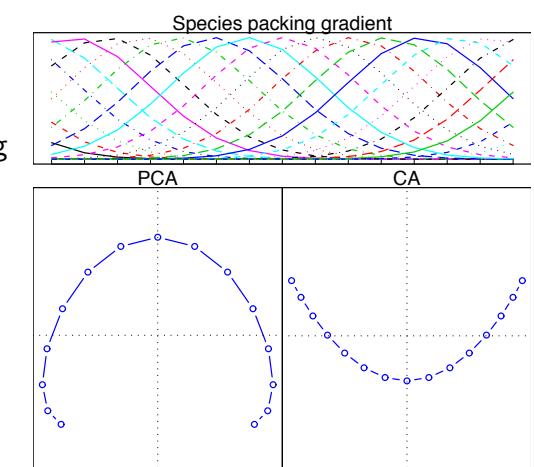
- NMDS
- Eigenvector Methods
- PCA
- CA
- Graphics
- Environmental Variables
- Gradient Model and Ordination

Environmental Interpretation

- Environmental variables need not be parallel to ordination axes.
- Axes cannot be taken as gradients, but gradients are oblique to axes:
You cannot tear off an axis from an ordination.
- **Never** calculate a correlation between an axis and an environmental variable.
- Environmental variables need not be linearly correlated with the ordination, but locations in ordination can be exceptional.

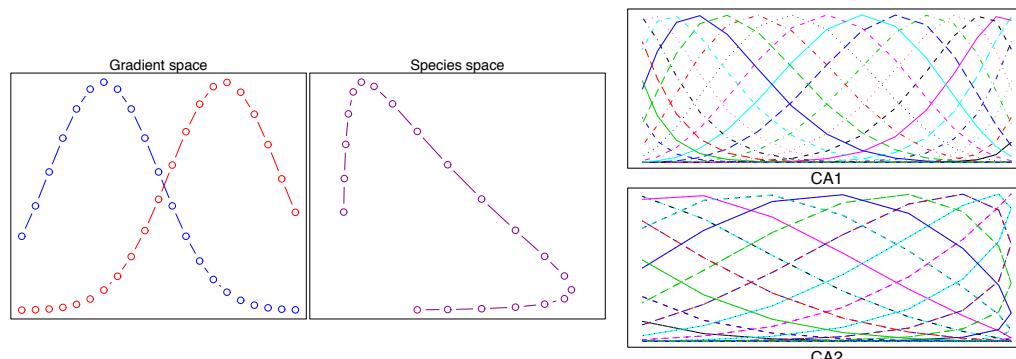
Gradient Model and Ordination

- Single gradients appear as curves in linear ordination methods
- PCA *horseshoe*: curve bends inward and gives wrong ordering of points on axis 1
- CA *arch*: axis 1 retains the correct ordering of sites despite the curve
- Environmental interpretation by vector fitting or surface bound to be biased
- Axes cannot be interpreted as "gradients"



The birth of the curve

- There is a curve in the species space and PCA shows it correctly
- CA deals better with unimodal responses, but the second optimal scaling axis is folded first axis



Solutions to the Curvature

Detrended Correspondence Analysis (DCA)

- CA axis retains the correct ordering: keep that, but instead of orthogonal axes, use detrended axes
- Programme DECORANA additionally rescales axes to sd units approximating t parameter of the Gaussian model
- Distorts space, introduces new artefacts and probably should be avoided

• Nonmetric Multidimensional Scaling (NMDS) should be able to cope with moderately long gradients

- Constrained ordination may linearize the responses

Running Detrended Correspondence Analysis

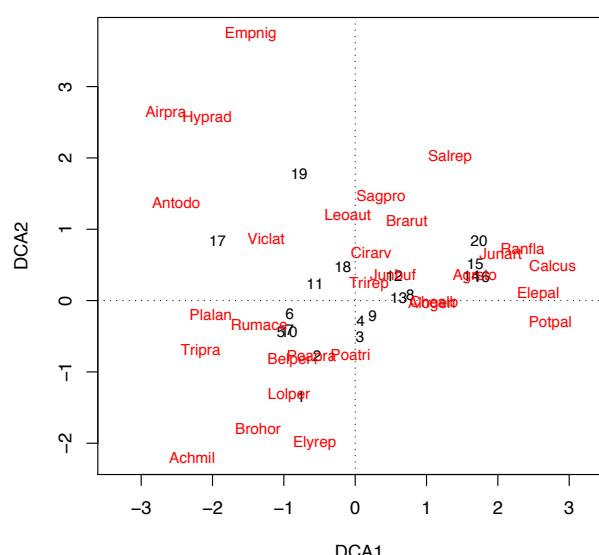
```
> (ord <- decorana(dune))

Call:
decorana(veg = dune)

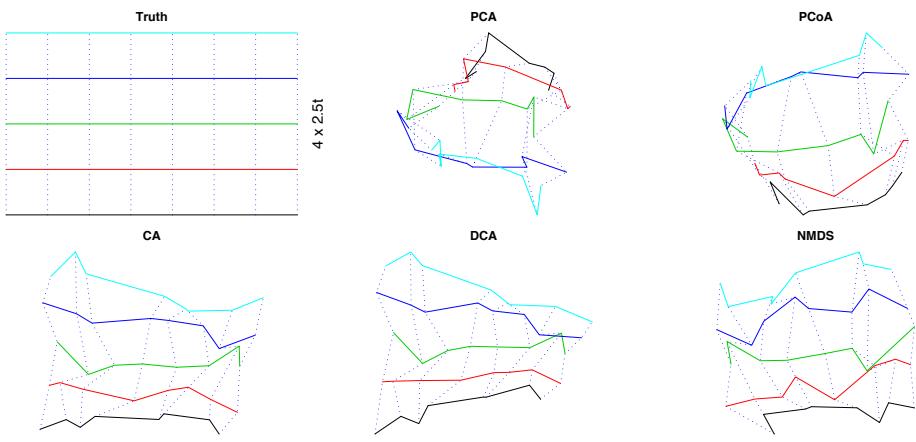
Detrended correspondence analysis with 26 segments.
Rescaling of axes with 4 iterations.
```

	DCA1	DCA2	DCA3	DCA4
Eigenvalues	0.512	0.304	0.1213	0.1427
Decorana values	0.536	0.287	0.0814	0.0481
Axis lengths	3.700	3.117	1.3006	1.4788

Default plot

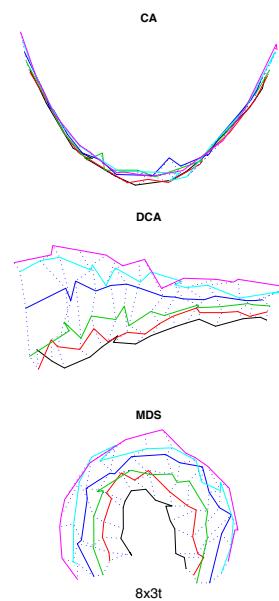


Community Pattern Simulation



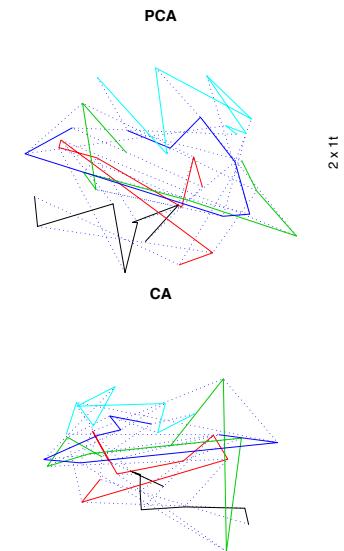
Long Gradients: DCA or NMDS

- Curvature with long gradients: Need either DCA or NMDS.
- NMDS is a test winner: More robust than DCA.
- DCA more popular.
- DCA may produce new artefacts, since it twists the space.



Short Gradients: Is There a Niche for PCA?

- Folklore: PCA with short gradients ($\leq 2t$).
- Not based on research, but simulation finds PCA uniformly worse than CA: At the best case about as good as CA.
- There should be no species optimum within gradient: Shortness alone not sufficient.
- PCA best used for really linear cases (environment) or for reduction of variables into principal components (but see FA).
- Noise dominates over signal in homogeneous data.



Extended Dissimilarities and Step-across

- How different are sites that have nothing in common?
- Use step-across points to estimate their distance
- Flexible shortest path or their approximations, extended dissimilarities
- Extended dissimilarity: use only one-site steps, do not update dissimilarities below a threshold

