

Reading

Nature Biotechnology **26**, 303 - 304 (2008)

[doi:10.1038/nbt0308-303](https://doi.org/10.1038/nbt0308-303)

What is principal component analysis?

Markus Ringnér¹

Principal component analysis is often incorporated into genome-wide expression studies, but what is it and how can it be used to explore high-dimensional data?

PCA as a tool to Quantify and Visualise

Multivariate Analysis

Multiple Regression / two way-
ANOVA / mixed models / machine
learning algorithms

Ordination methods

What is the difference between these two pairwise correlation matrices?

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.80	0.90	0.78	0.87
X_2	0.80	1.00	0.76	0.87	0.78
X_3	0.90	0.76	1.00	0.78	0.89
X_4	0.78	0.87	0.78	1.00	0.95
X_5	0.87	0.78	0.89	0.95	1.00

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.87	0.96	0.04	0.05
X_2	0.87	1.00	0.95	0.03	0.07
X_3	0.96	0.95	1.00	0.04	0.05
X_4	0.04	0.03	0.04	1.00	0.84
X_5	0.05	0.07	0.05	0.84	1.00

What is the difference between these two pairwise correlation matrices?

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.80	0.90	0.78	0.87
X_2	0.80	1.00	0.76	0.87	0.78
X_3	0.90	0.76	1.00	0.78	0.89
X_4	0.78	0.87	0.78	1.00	0.95
X_5	0.87	0.78	0.89	0.95	1.00

One
dimension

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.87	0.96	0.04	0.05
X_2	0.87	1.00	0.95	0.03	0.07
X_3	0.96	0.95	1.00	0.04	0.05
X_4	0.04	0.03	0.04	1.00	0.84
X_5	0.05	0.07	0.05	0.84	1.00

Two
dimensions

Ordination analyses

- Uncover, organize and summarize the main patterns of variation in a set of variables measured over multiple observations.
- Patterns of variation are structured in a reduced space with smaller number number of dimensions.
- Reduction is possible because often variables are associated (e.g., correlated). Dimensions represent combinations (e.g., linear combinations of variables).

Ordination analyses

A procedure for adapting a multidimensional swarm of data points in such a way that when it is projected onto a reduced number of dimensions any intrinsic pattern will become apparent.

Adapted from Connie Clark

Ordination analyses – uncover and organize data; a quick example:

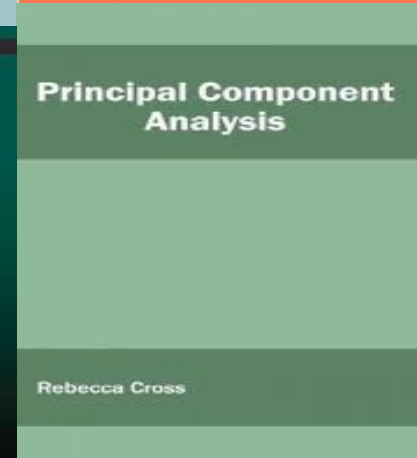
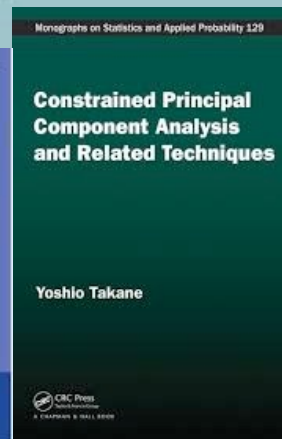
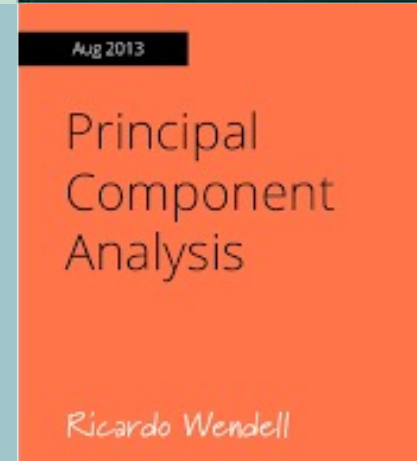
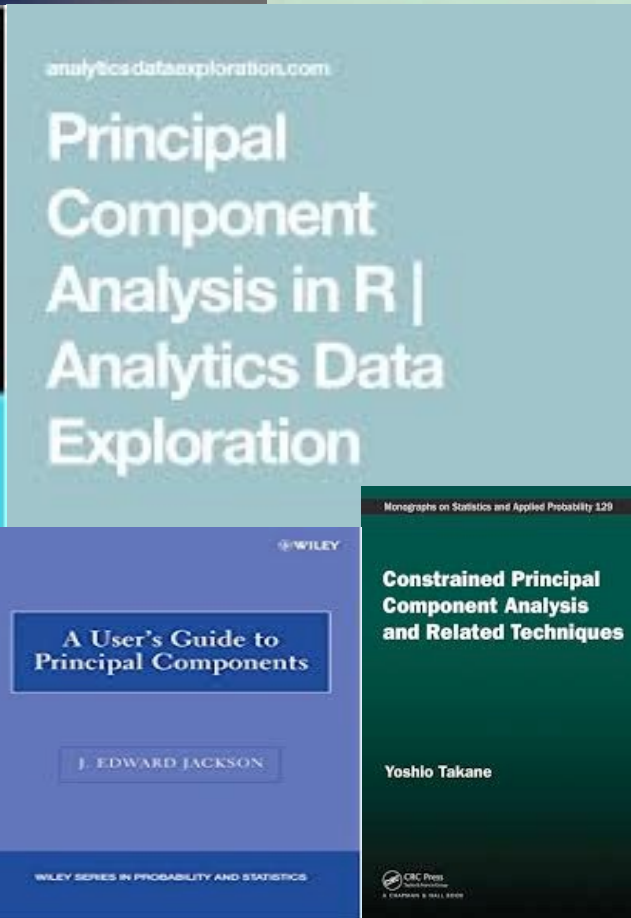
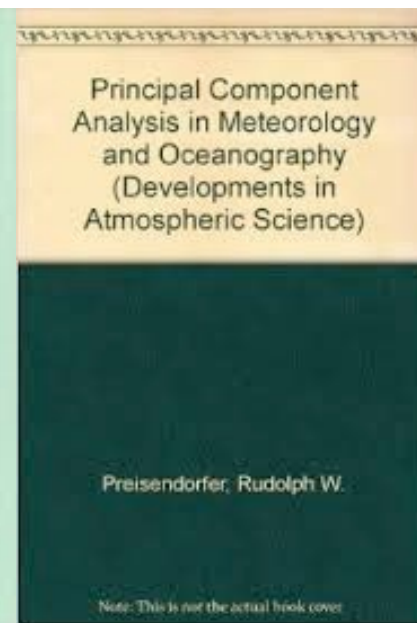
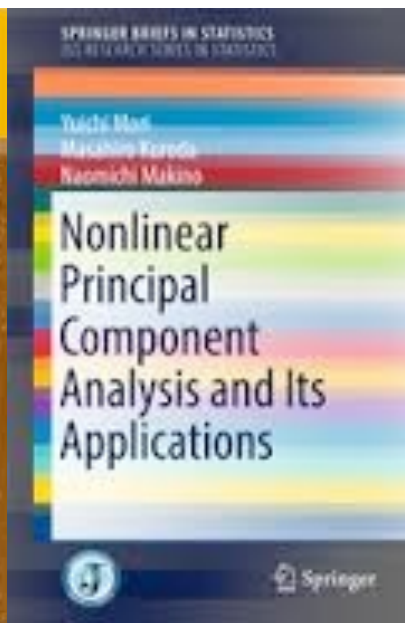
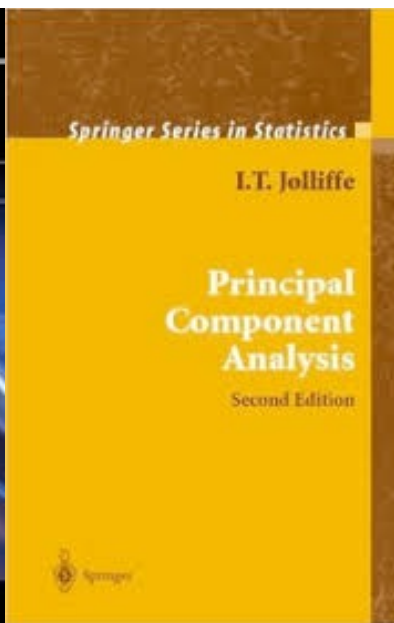
Species

Site	B	I	D	A	H	E	G	C
4	1	0	1	0	0	0	0	1
1	0	0	0	1	0	0	0	0
7	0	0	0	0	1	1	1	0
8	0	1	0	0	1	0	1	0
6	0	0	1	0	0	1	1	0
5	0	0	1	0	0	1	0	1
10	0	1	0	0	0	0	0	0
2	1	0	0	1	0	0	0	0
9	0	1	0	0	1	0	0	0
3	1	0	0	1	0	0	0	1

Ordination methods

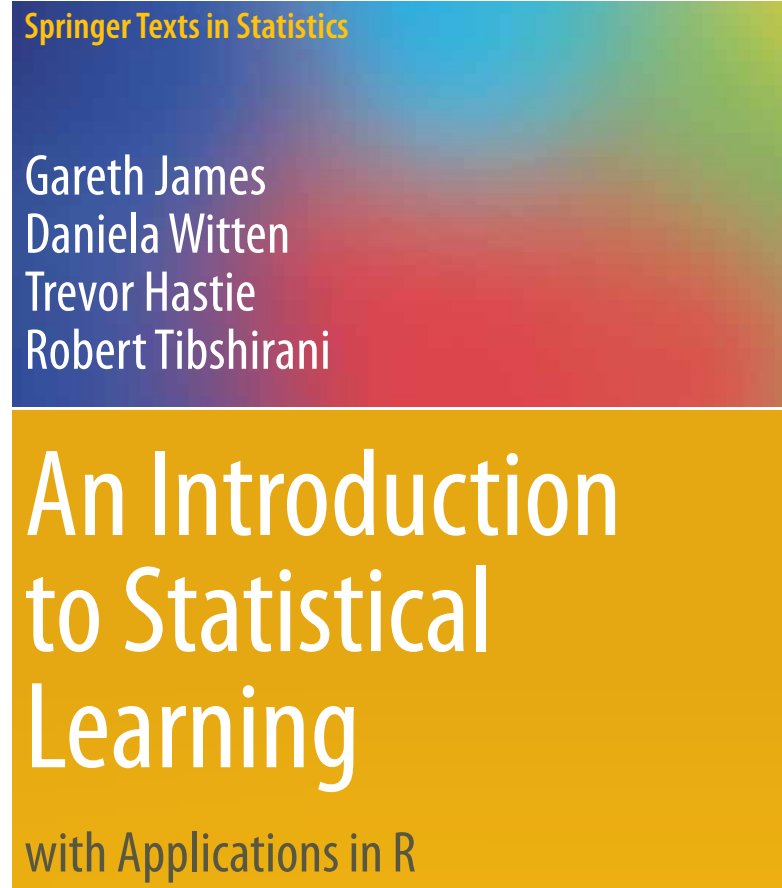
- **Principal Component Analysis (PCA)**
- Correspondence Analysis (CA)
- Principal Coordinate Analysis (PCoA)
- Discriminant Function Analysis (DFA)
- Principal Curve Analysis
- Etc, etc, etc...

Principal components analysis (PCA) is perhaps the most common technique used to summarize patterns among variables in multivariate datasets.



Some treat Principal Component Analysis (PCA) as an unsupervised learning method (an exploratory technique such as k-means)

10 Unsupervised Learning	373
10.1 The Challenge of Unsupervised Learning	373
10.2 Principal Components Analysis	374
10.2.1 What Are Principal Components?	375
10.2.2 Another Interpretation of Principal Components	379
10.2.3 More on PCA	380
10.2.4 Other Uses for Principal Components	385
10.3 Clustering Methods	385
10.3.1 <i>K</i> -Means Clustering	386
10.3.2 Hierarchical Clustering	390
10.3.3 Practical Issues in Clustering	399



Supervised *versus* unsupervised learning techniques

- Techniques for unsupervised learning are fast growing in a number of fields, particularly biology.
- A cancer researcher might assay gene expression levels in 100 patients with breast cancer. They might then look for subgroups among the breast cancer samples, or among the genes, in order to obtain a better understanding of the disease.
- A search engine might choose what search results to display to a particular individual based on the click histories of other individuals with similar search patterns. These statistical learning tasks, and many more, can be performed via unsupervised learning techniques.

Supervised *versus* unsupervised learning techniques

In contrast, unsupervised learning is often much more challenging. The exercise tends to be more subjective, and there is no simple goal for the analysis, such as prediction of a response.

Unsupervised learning is often performed as part of an exploratory data analysis.

Hard to assess the results obtained given that there is no universally accepted mechanism for performing cross-validation or validating results on an independent data set; there is no way to check how the models does because we don't know the true answer—the problem is unsupervised.

Adapted from James et al. 2013

Examples of Principal Component Analysis



Principal components analysis (PCA) - example 1

A subtype of childhood acute lymphoblastic leukaemia with poor treatment outcome: a genome-wide classification study



*Monique L Den Boer**, *Marjon van Slegtenhorst**, *Renée X De Menezes*, *Meyling H Cheok*, *Jessica G C A M Buijs-Gladdines*, *Susan T C J M Peters*, *Laura J C M Van Zutven*, *H Berna Beverloo*, *Peter J Van der Spek*, *Gaby Escherich†*, *Martin A Horstmann†*, *Gritta E Janka-Schaub†*, *Willem A Kamps‡*, *William E Evans*, *Rob Pieters‡*

Summary

Background Genetic subtypes of acute lymphoblastic leukaemia (ALL) are used to determine risk and treatment in children. 25% of precursor B-ALL cases are genetically unclassified and have intermediate prognosis. We aimed to use a genome-wide study to improve prognostic classification of ALL in children.

Lancet Oncol 2009; 10: 125-34

Published [Online](#)

January 9, 2009

DOI:10.1016/S1470-

Quantification and Visualisation

Principal components analysis (PCA) - example 1

A subtype of childhood acute lymphoblastic leukaemia with poor treatment outcome: a genome-wide classification study



Monique L Den Boer, Marjon van Slegtenhorst*, Renée X De Menezes, Meyling H Cheok, Jessica G C A M Buijs-Gladdines, Susan T C J M Peters, Laura J C M Van Zutven, H Berna Beverloo, Peter J Van der Spek, Gaby Escherich†, Martin A Horstmann†, Gritta E Janka-Schaub†, Willem A Kamps‡, William E Evans, Rob Pieters‡*

Summary

Background Genetic subtypes of acute lymphoblastic leukaemia (ALL) are used to determine risk and treatment in children. 25% of precursor B-ALL cases are genetically unclassified and have intermediate prognosis. We aimed to use a genome-wide study to improve prognostic classification of ALL in children.

Lancet Oncol 2009; 10: 125-34
Published Online
January 9, 2009
DOI:10.1016/S1470-

Data matrix: 190 observations by 22283 columns

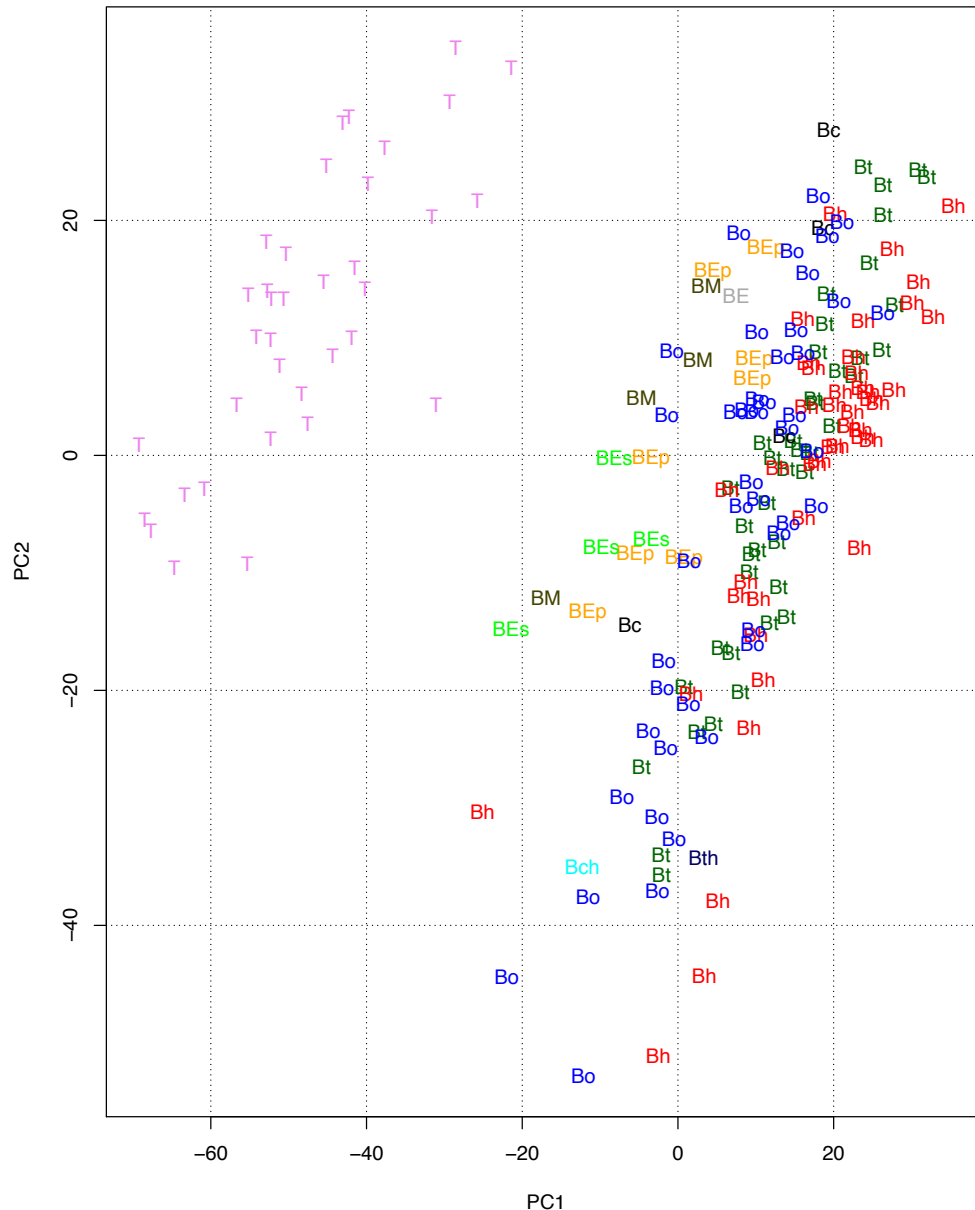
Gene expression (22283 genes)

Gene expression
(190 patients)



Principal components analysis (PCA) - example 1

PCA; Den Boer (2009); 190 samples * 22283 genes



Each letter is a patient.
Labels stand for
different lymphoblastic
leukaemia (ALL) types.

Data matrix: 190
observations by 22283
columns.

Principal components analysis (PCA) - example 2



PCA – A Powerful Method for Analyze Ecological Niches

Franc Janžekovič and Tone Novak
*University of Maribor, Faculty of Natural Sciences and Mathematics,
Department of Biology, Maribor
Slovenia*

Principal components analysis (PCA) - example 2

2.1 Environmental niche of three hymenopteran and two spider species

Between 1977 and 2004, 63 caves and artificial tunnels were ecologically investigated in Slovenia; the three most abundant Hymenoptera species found in these studies have been ecologically evaluated (details in Novak et al. 2010a). In the caves, many environmental data were collected, as follows. The following abbreviations of the environmental variables are used: Dist-E = distance from entrance; Dist-S = distance from surface; Illum = illumination; PCS = passage cross-section; Tair = air temperature; RH = relative air humidity; Tgr = ground temperature; HY = substrate moisture. The hymenopteran spatial niche breadth was originally represented by nine variables.

Data matrix: 63 observations (caves) by 9 columns

Environmental variables (9)

63 caves



**PCA – A Powerful Method
for Analyze Ecological Niches**

Franc Janžekovič and Tone Novak
University of Maribor, Faculty of Natural Sciences and Mathematics,
Department of Biology, Maribor
Slovenia

Principal components analysis (PCA) - example 2

(pairwise correlation among environmental variables)

	1	2	3	4	5	6	7	8	9
1 Air temperature	1.00 ---								
2 <i>arc-sin</i> relative air humidity	0.15 0.133	1.00 ---							
3 Ground temperature	0.94 <0.001	0.18 0.079	1.00 ---						
4 <i>arc-sin</i> substrate moisture	0.388 <0.001	0.59 <0.001	0.37 <0.001	1.00 ---					
5 Airflow	-0.48 <0.001	-0.36 <0.001	-0.43 <0.001	-0.55 <0.001	1.00 ---				
6 Distance from entrance	-0.34 <0.001	0.14 0.153	-0.41 <0.001	0.10 0.312	0.04 0.712	1.00 ---			
7 Distance from surface	-0.02 0.837	0.24 0.017	-0.04 0.683	0.46 <0.001	-0.11 0.275	0.67 <0.001	1.00 ---		
8 Passage cross-section	0.35 <0.001	0.17 0.089	0.23 0.025	0.39 <0.001	-0.40 <0.001	-0.11 0.274	0.05 0.656	1.00 ---	
9 <i>log</i> illumination	0.45 <0.001	-0.18 0.077	0.46 <0.001	-0.04 0.690	-0.07 0.494	-0.821 <0.001	-0.679 <0.001	0.37 <0.001	1.00 ---

Table 1. Pearson correlations coefficient among nine environmental variables. Significant correlations in bold. (Upper row r, lower row p).

PCA – A Powerful Method for Analyze Ecological Niches

Franc Janžekovič and Tone Novak
University of Maribor, Faculty of Natural Sciences and Mathematics,
Department of Biology, Maribor
Slovenia

Principal components analysis (PCA) - example 2

(niche differences – dots represent different caves ellipsoids are confidence intervals for where species is found)

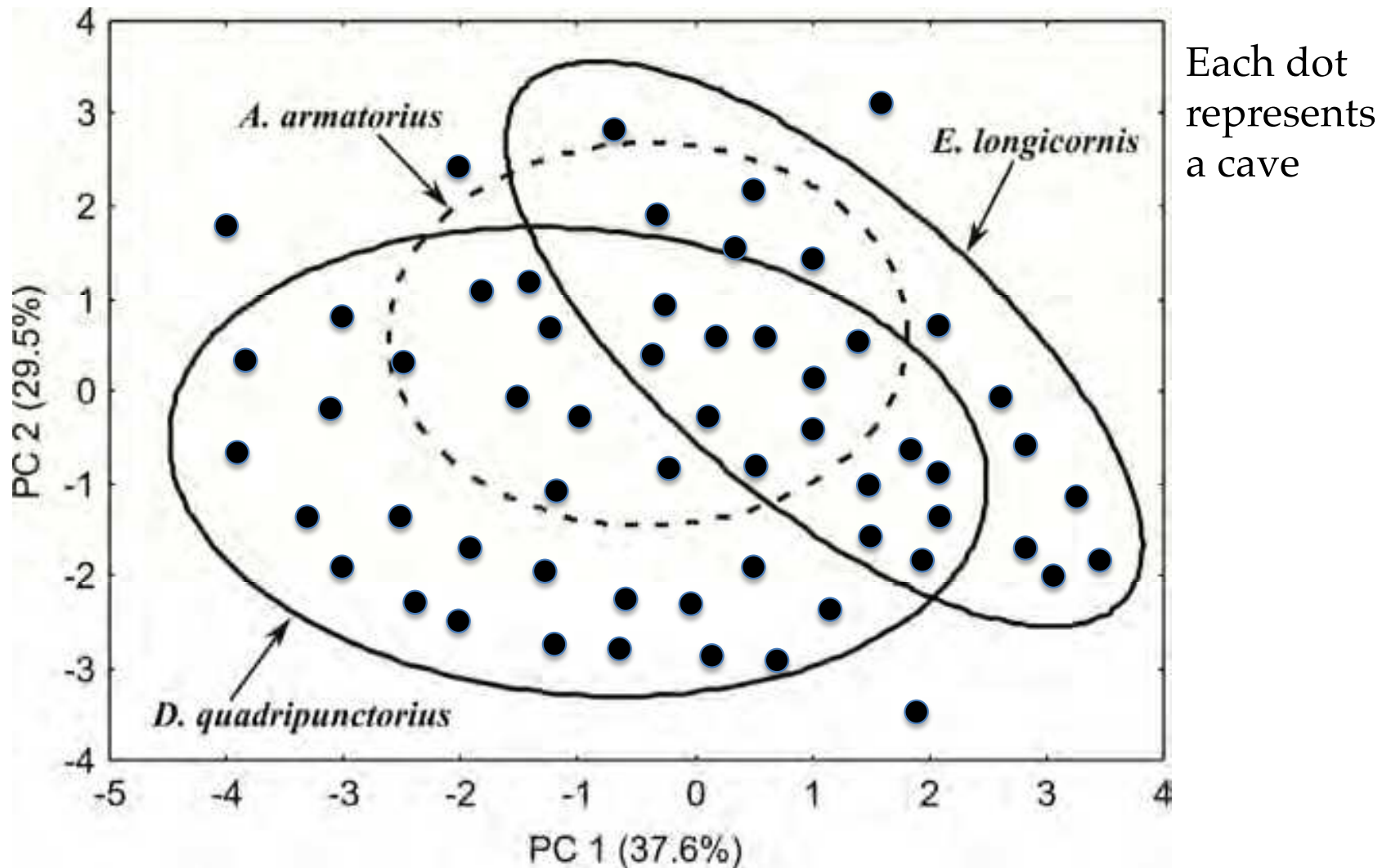
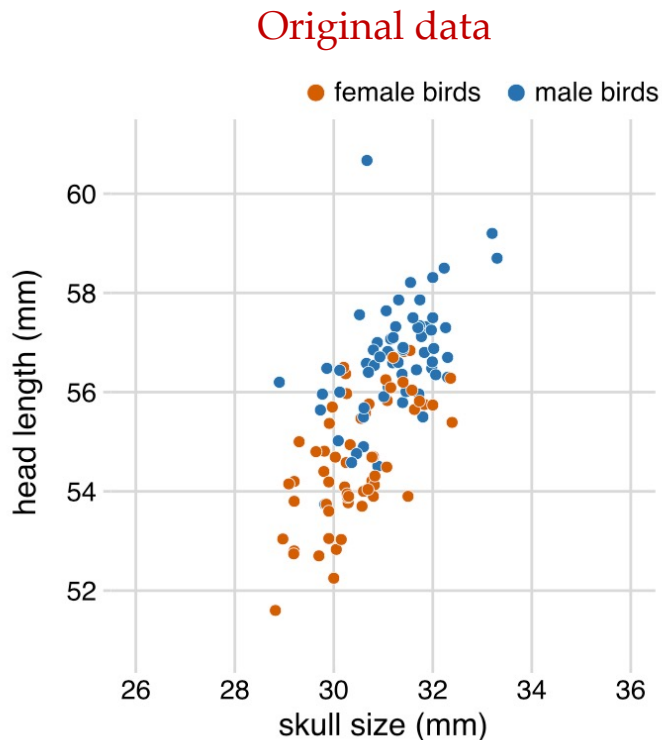


Fig. 5. Ordination of the nine environmental variables in 1st and 2nd PC axes. Ellipses (95% confidence) represent spatial niches in the three hymenopteran species.

Principal Component Analysis (PCA): A geometric interpretation

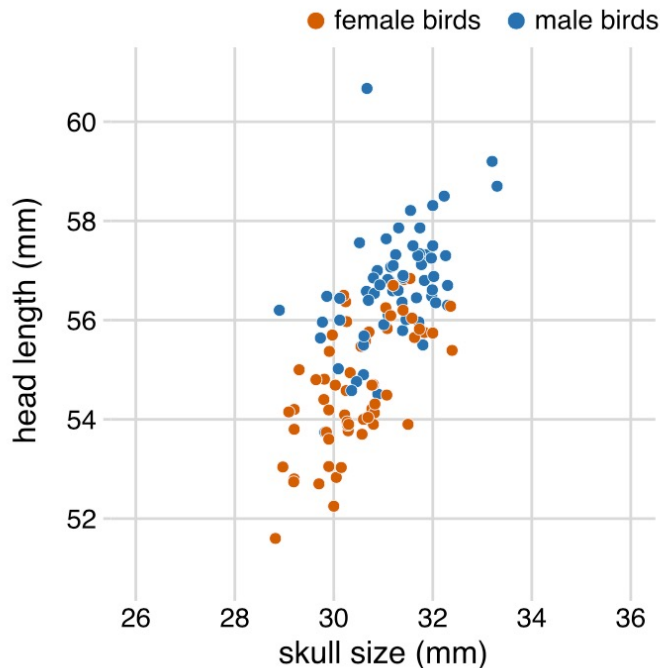
PCA finds the coordinate system (called principal components) that best represents the internal variability in the data, essentially re-projecting the data on these coordinate system. As such, PCA represents associations among variables (gene, environmental variables) and data points are re-projected so that the correlations among variables is maximized.



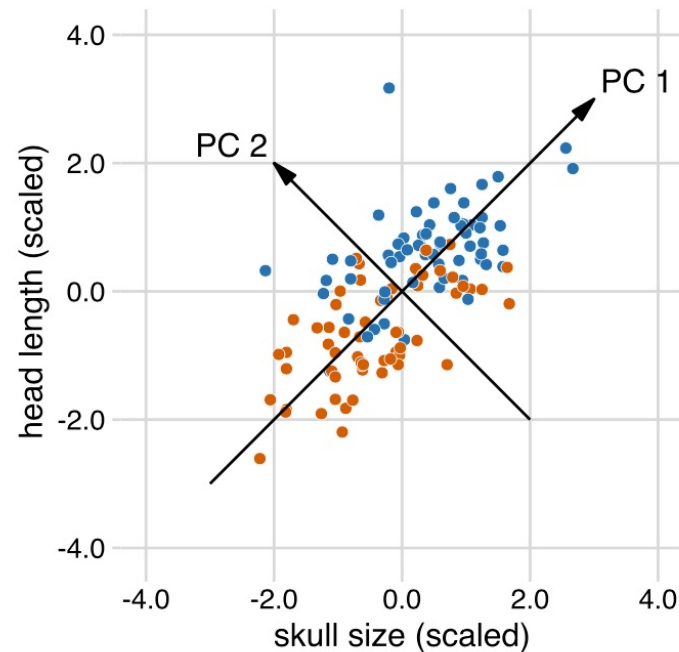
Principal Component Analysis (PCA): A geometric interpretation

PCA finds the coordinate system (called principal components) that best represents the internal variability in the data, essentially re-projecting the data on these coordinate system. As such, PCA represents associations among variables (gene, environmental variables) and data points are re-projected so that the correlations among variables is maximized.

Original data



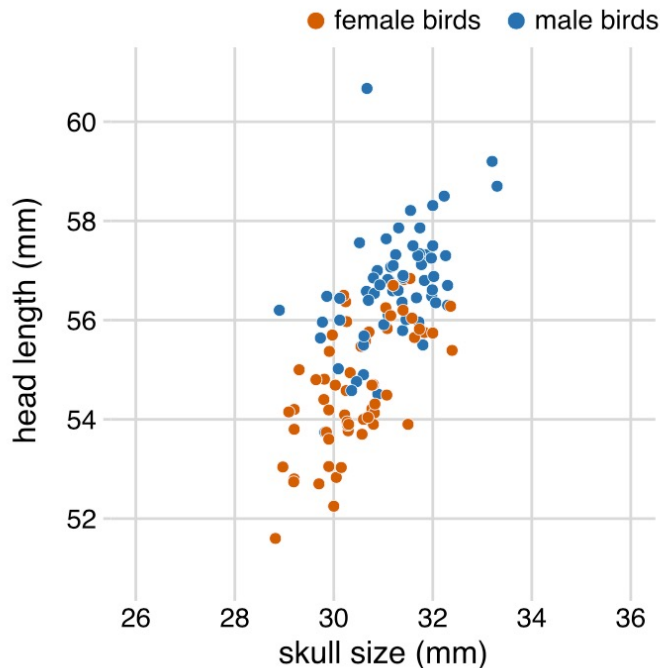
Standardization and PCA fitting



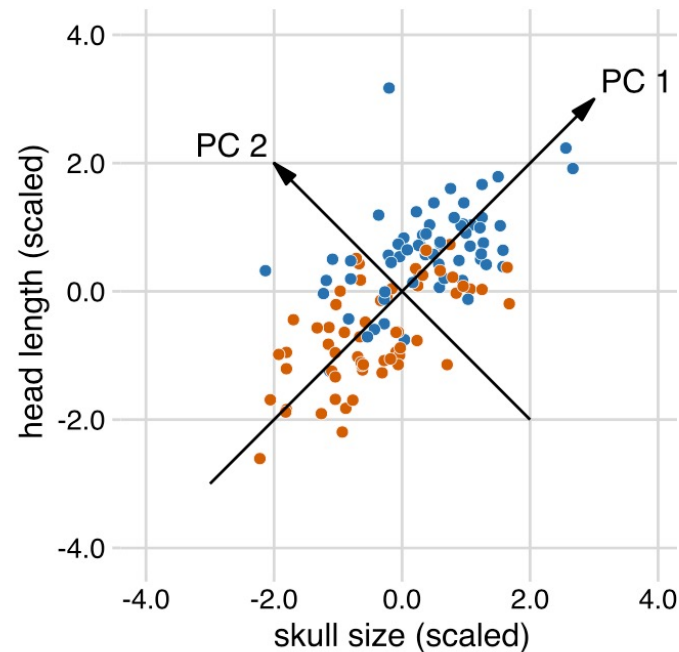
Principal Component Analysis (PCA): A geometric interpretation

PCA finds the coordinate system (called principal components) that best represents the internal variability in the data, essentially re-projecting the data on these coordinate system. As such, PCA represents associations among variables (gene, environmental variables) and data points are re-projected so that the correlations among variables is maximized.

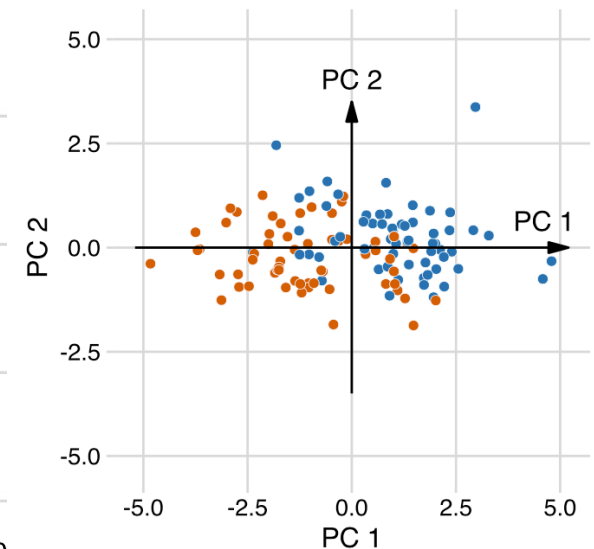
Original data



Standardization and PCA fitting



Rotation



PCA aligns their axes with directions of maximum variation in the data

Principal Component Analysis (PCA): A geometric interpretation

- PCA constructs a new coordinate system (new variables, PCs) which are linear combinations of the original data and which are defined to align the samples along their major axes of variation (assuming linearity).
- Thus, PCA determines the coordinate system that best represents the internal variability in the data, essentially re-projecting the data.

The association among variables need to be measured by either (in most cases):

Correlation Matrix (for variables that have different units or scales, e.g., pH, temperature).

Covariance Matrix (variables have the same units, e.g., body length & body width in cm).

Raw data when variables are in the same units (more difficult to interpret) and calculations differ (very rare to find applications in the literature); rarely used.

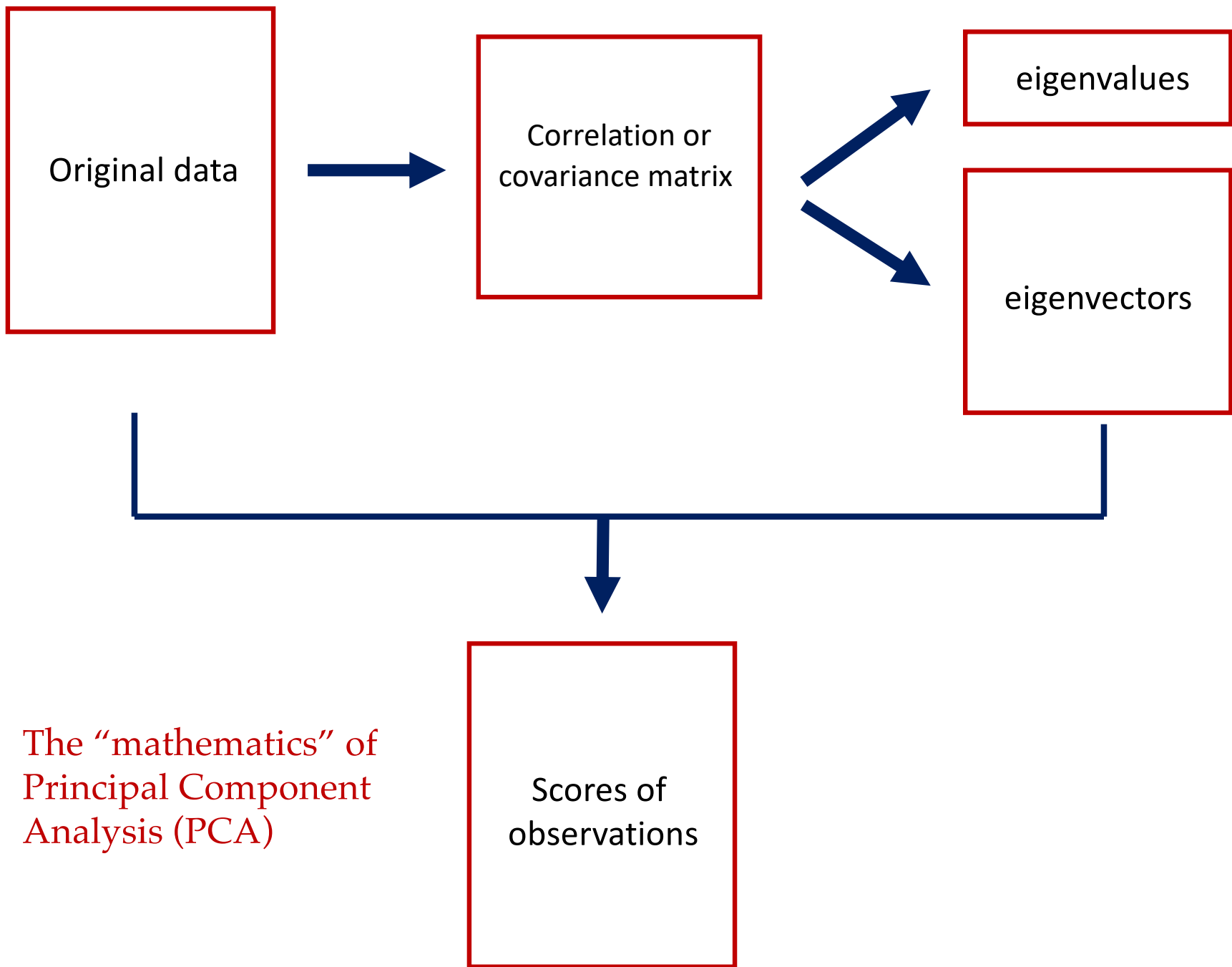
Correlation *versus* covariance

$$COV_{xy} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

$$\bar{X} = 0 \text{ \& } \bar{Y} = 0 \therefore s_x = s_x \text{ \& } s_y = s_y$$

$$COR_{xy} = \frac{COV_{xy}}{s_x s_y}$$

$$\bar{X} = 0 \text{ \& } \bar{Y} = 0 \therefore s_x = 1 \text{ \& } s_y = 1$$



The "mathematics" of
Principal Component
Analysis (PCA)

The mathematics of Principal Component Analysis (PCA):

Eigen-analysis is a mathematical operation on a *square symmetric* matrix (e.g., pairwise correlation matrix, pairwise covariance matrix).

A *square* matrix has the same number of rows as columns.

A *symmetric* matrix is the same if you switch rows and columns.

square and symmetric matrix

(e.g., pairwise correlation matrix)

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.80	0.90	0.78	0.87
X_2	0.80	1.00	0.76	0.87	0.78
X_3	0.90	0.76	1.00	0.78	0.89
X_4	0.78	0.87	0.78	1.00	0.95
X_5	0.87	0.78	0.89	0.95	1.00

The important components of Principal Component Analysis (pun intended)



Principal component analysis presents three important structures:

1 – **Eigenvalues:** represent the amount of variation in the original data summarized by each principal component. The first principal component (PC-1) presents the largest amount, PC-2 presents the second largest amount, and so on.

Eigenvalues

	X ₁	X ₂	X ₃	X ₄	X ₅
X ₁	1.00	0.80	0.90	0.78	0.87
X ₂	0.80	1.00	0.76	0.87	0.78
X ₃	0.90	0.76	1.00	0.78	0.89
X ₄	0.78	0.87	0.78	1.00	0.95
X ₅	0.87	0.78	0.89	0.95	1.00

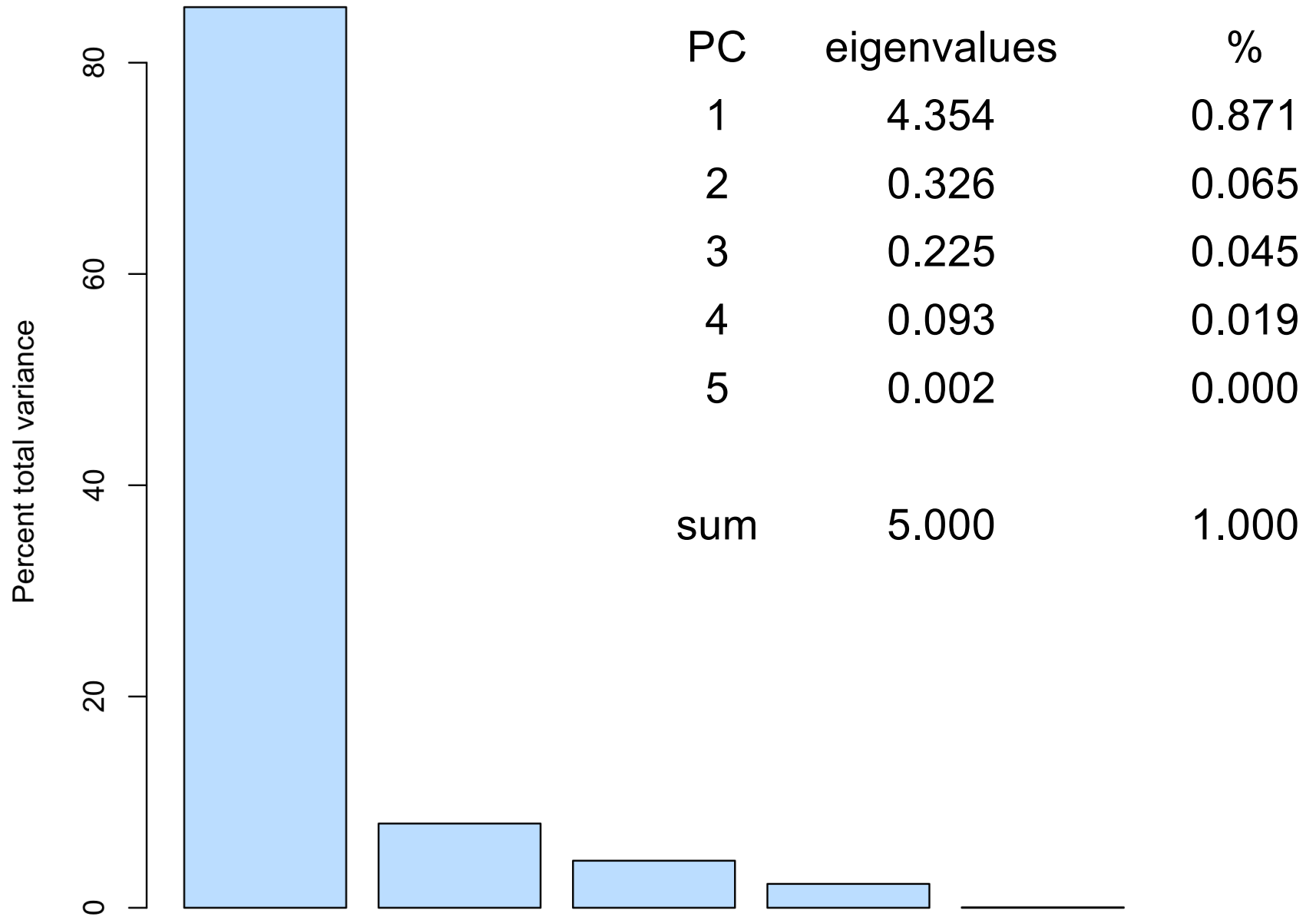
“one
dimension”

Eigenvalues:

PC	eigenvalues	%
1	4.354	0.871
2	0.326	0.065
3	0.225	0.045
4	0.093	0.019
5	0.002	0.000
sum	5.000	1.000

“Lower”
dimensionality
because it kept a
large proportion of
the variation in the
data in the first PC.

Plot of eigenvalue contributions



Eigenvalues

1.00	0.87	0.96	0.04	0.05
0.87	1.00	0.95	0.03	0.07
0.96	0.95	1.00	0.04	0.05
0.04	0.03	0.04	1.00	0.84
0.05	0.07	0.05	0.84	1.00

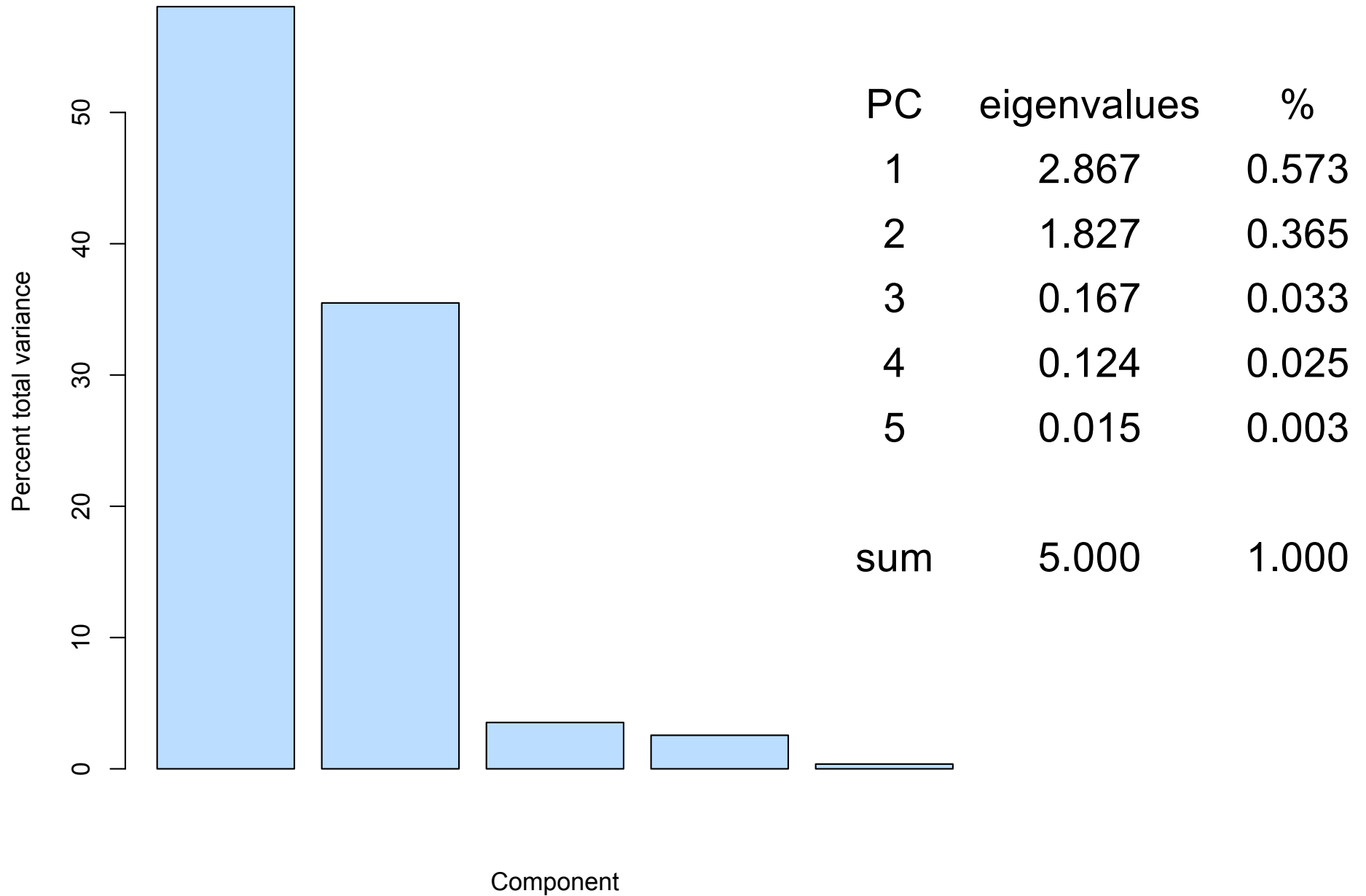
“two
dimensions”

Eigenvalues:

PC	eigenvalues	%
1	2.867	0.573
2	1.827	0.365
3	0.167	0.033
4	0.124	0.025
5	0.015	0.003
sum	5.000	1.000

“Higher” dimensionality
because two components
are needed to summarize
variation.

Plot of eigenvalues



Principal component analysis presents three important structures:

2 - **Eigenvectors**: Each principal component is a linear function with coefficients for each variable.

- Eigenvectors contain these coefficients. High values, positive or negative, represents high association with the component.

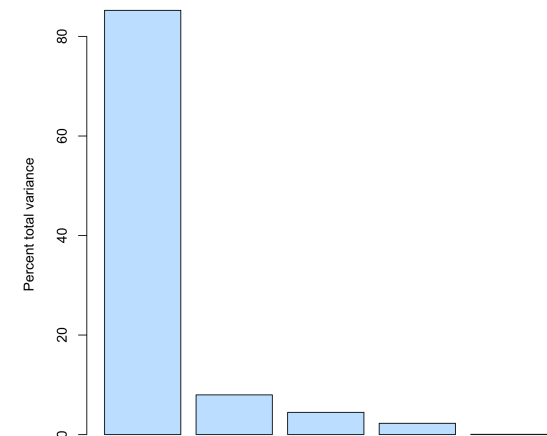
Correlation matrix

	X ₁	X ₂	X ₃	X ₄	X ₅
X ₁	1.00	0.80	0.90	0.78	0.87
X ₂	0.80	1.00	0.76	0.87	0.78
X ₃	0.90	0.76	1.00	0.78	0.89
X ₄	0.78	0.87	0.78	1.00	0.95
X ₅	0.87	0.78	0.89	0.95	1.00

“one
dimension”

Associated eigenvectors

	PC				
var	1	2	3	4	5
1	0.447	-0.436	0.330	-0.687	0.170
2	0.432	0.533	0.644	0.181	-0.288
3	0.445	-0.534	0.035	0.692	0.192
4	0.450	0.489	-0.413	-0.063	0.619
5	0.462	-0.039	-0.552	-0.109	-0.684



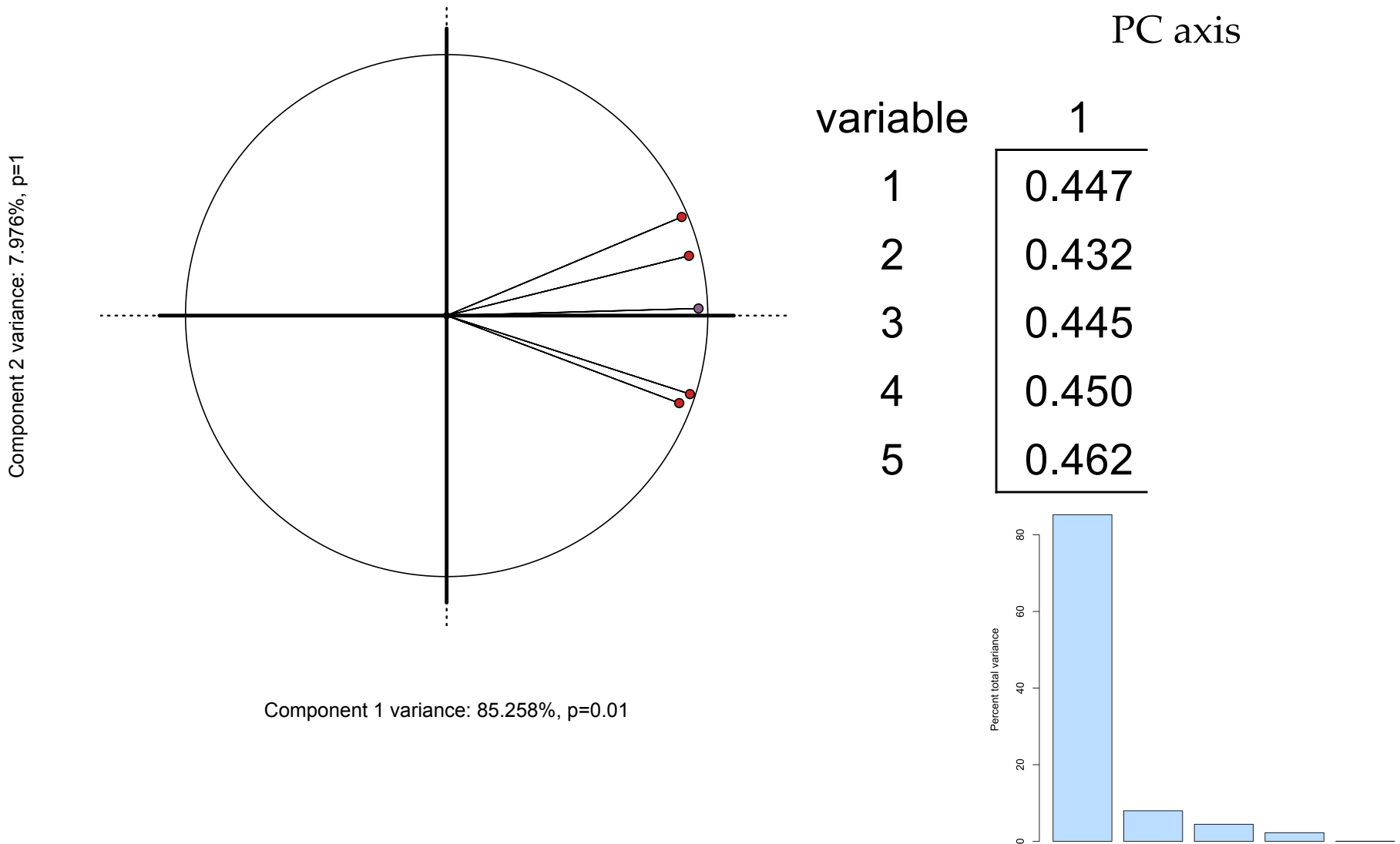
Eigenvectors can be seen as regression coefficients, where the component is the dependent variable. A “one dimension” matrix has only one interpretable principal component.

$$\text{PC-1} = 0.447X_1 + 0.432X_2 + 0.445X_3 + 0.450X_4 + 0.462X_5$$

↙ Unlike the numbers after =, this is not a subtraction but a hyphen stating that this is the first and second Principal Components (PC).

	PC				
var	1	2	3	4	5
1	0.447	0.436	0.330	-0.687	0.170
2	0.432	-0.533	0.644	0.181	-0.288
3	0.445	0.534	0.035	0.692	0.192
4	0.450	-0.489	-0.413	-0.063	0.619
5	0.462	0.039	-0.552	-0.109	-0.684

Eigenvectors (simulated data with 1 dimension): only first axis (PC-1) should be interpreted



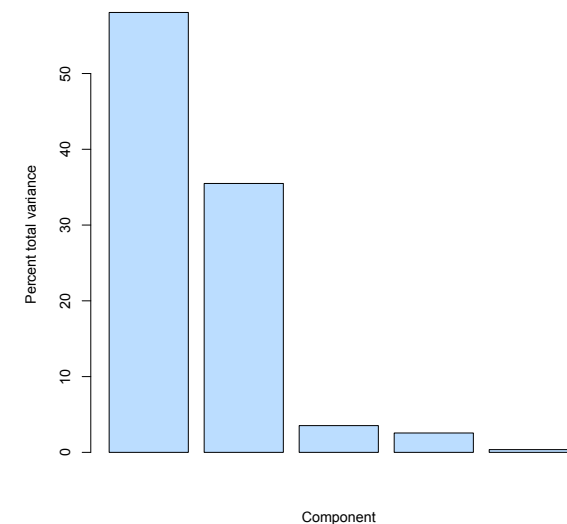
Correlation matrix

	X_1	X_2	X_3	X_4	X_5
X_1	1.00	0.87	0.96	0.04	0.05
X_2	0.87	1.00	0.95	0.03	0.07
X_3	0.96	0.95	1.00	0.04	0.05
X_4	0.04	0.03	0.04	1.00	0.84
X_5	0.05	0.07	0.05	0.84	1.00

“two dimensions”

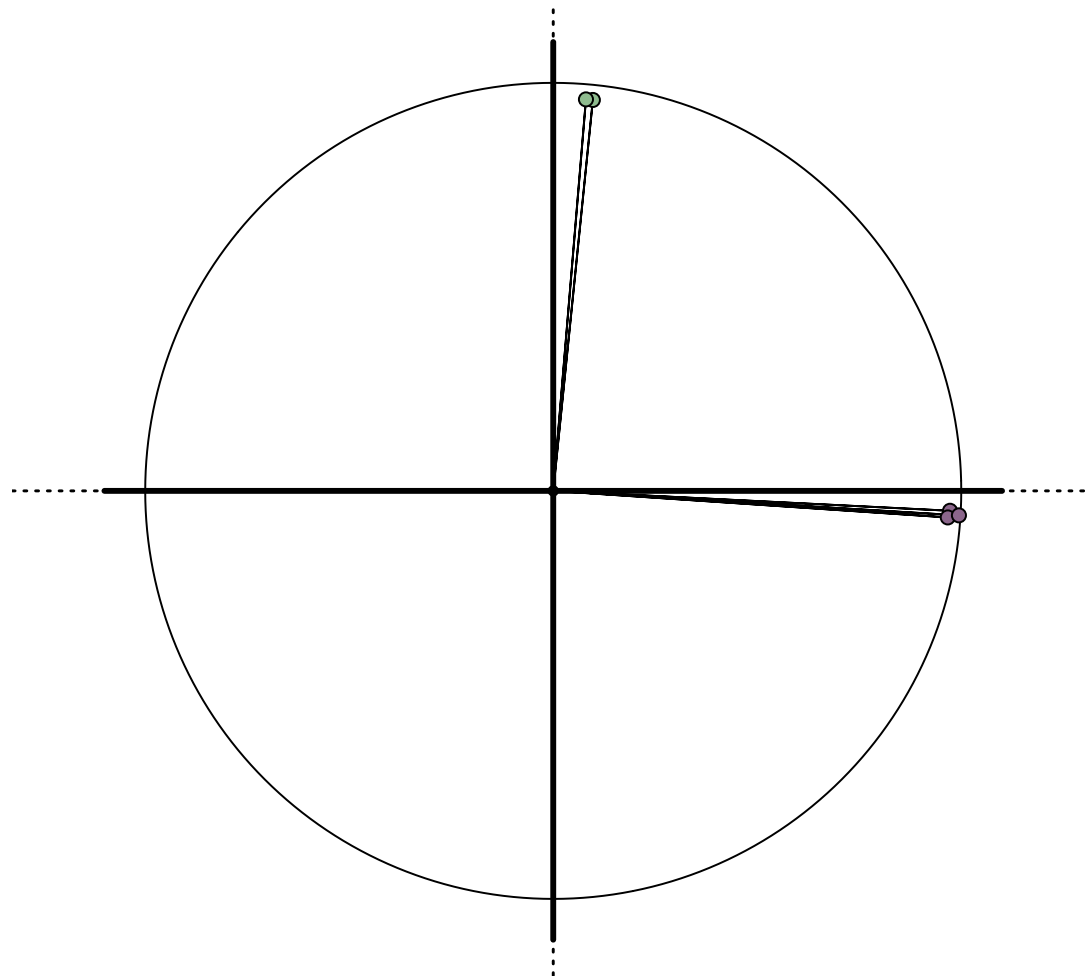
Associated eigenvectors (only interpret the first two components (PC))

	PC				
var	1	2	3	4	5
1	0.569	-0.064	0.249	-0.642	0.445
2	0.567	-0.060	-0.298	0.661	0.386
3	0.585	-0.067	0.061	-0.810	-0.806
4	0.072	0.704	0.651	0.273	0.039
5	0.085	0.702	-0.650	-0.277	-0.043



Eigenvectors (simulated data with 2 dimensions): only first two axis (PC-1 & PC-2) should be interpreted

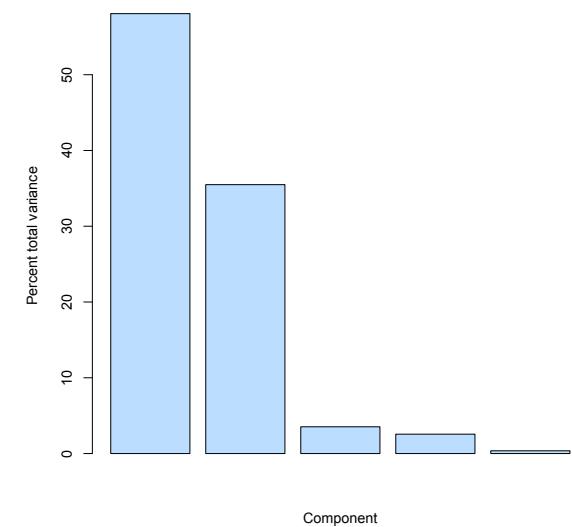
Component 2 variance: 36.974%, p=0.01



Component 1 variance: 57.703%, p=0.01

PC axis

var	1	2
1	0.569	-0.064
2	0.567	-0.060
3	0.585	-0.067
4	0.072	0.704
5	0.085	0.702



Principal component analysis presents three important structures:

3 – **Multivariate scores:** Since each component is a linear function of the variables, when multiplying the standardized variables (in the case of correlation matrices) by the eigenvector structure, a matrix containing the position of each observation in each principal component is produced.

The plot of these scores in the first few dimensions, represents the main patterns of variation among the original observations (more in the empirical example).

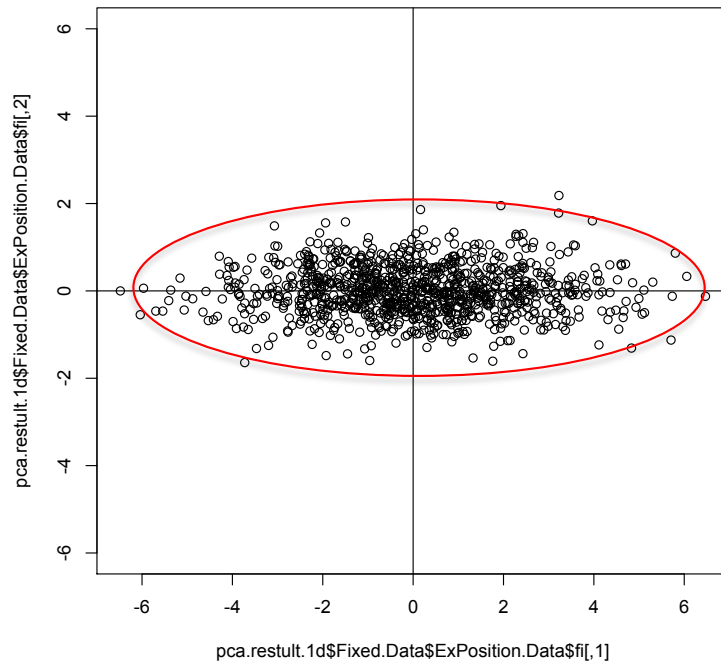
$$\text{PC-1} = 0.569X_1 + 0.567X_2 + 0.585X_3 + 0.072X_4 + 0.085X_5$$

$$\text{PC-2} = -0.064X_1 - 0.060X_2 - 0.067X_3 + 0.704X_4 + 0.702X_5$$

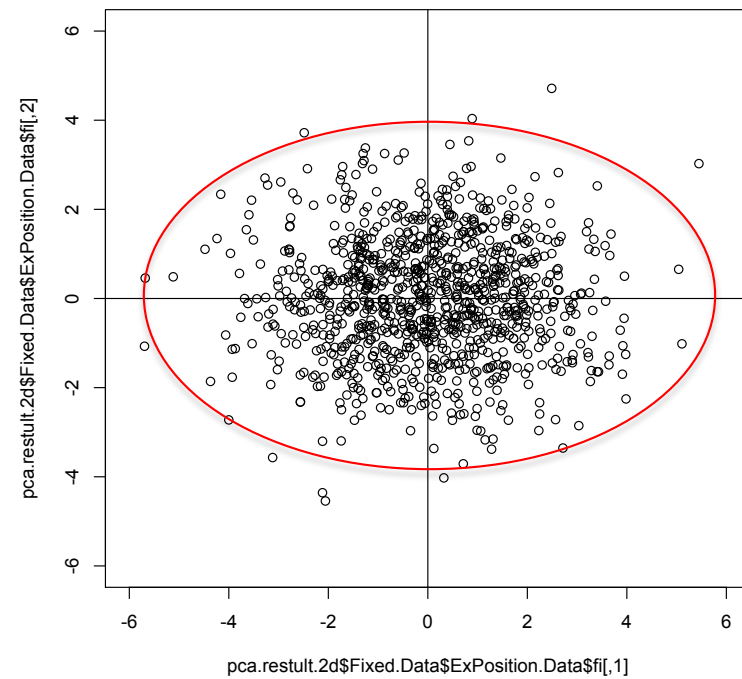
	PC axis	
var	1	2
1	0.569	-0.064
2	0.567	-0.060
3	0.585	-0.067
4	0.072	0.704
5	0.085	0.702

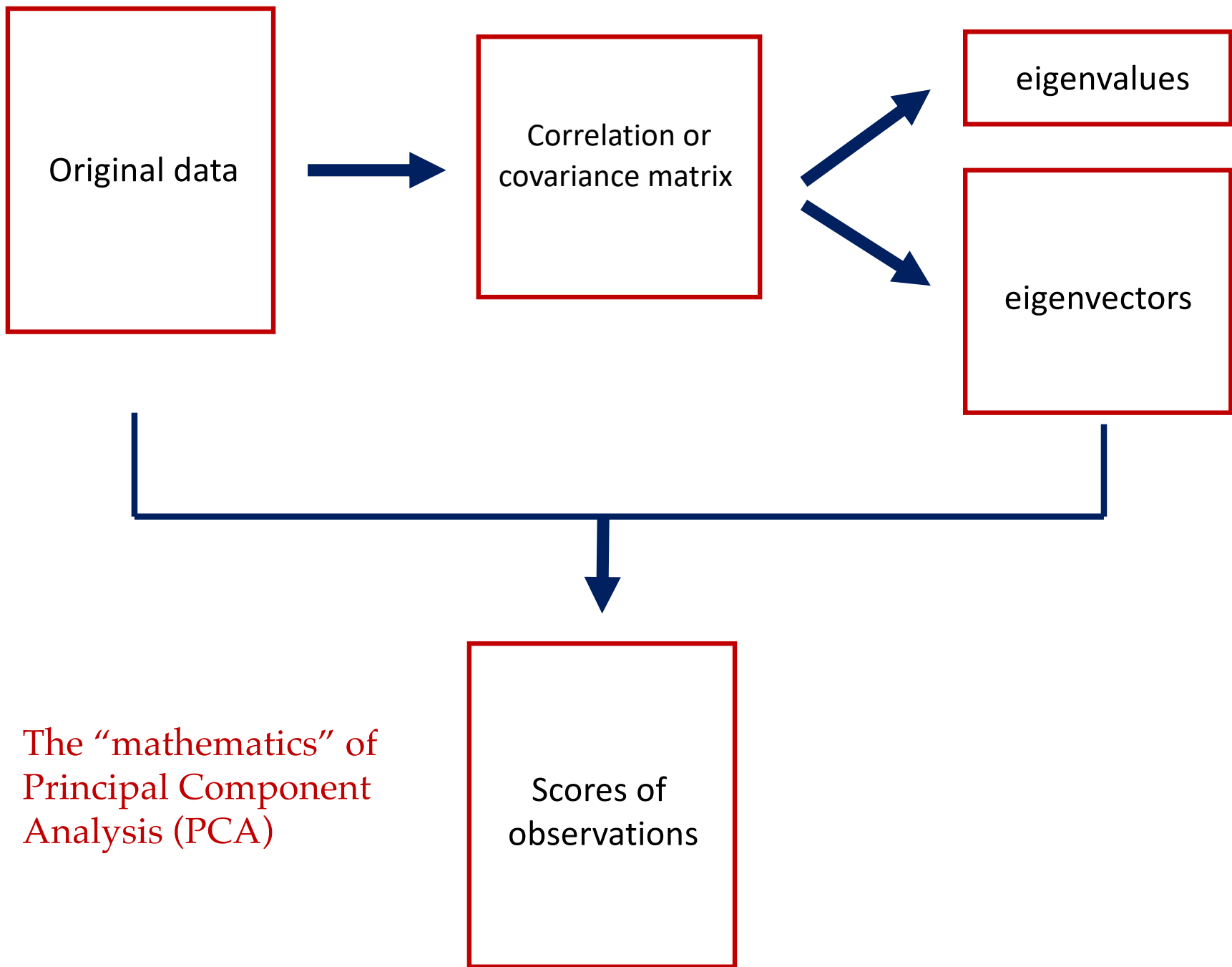
PCA Scores: one versus two dimensions

1D



2D





The "mathematics" of
Principal Component
Analysis (PCA)

Next lecture: How many PCA dimensions?
Inferential frameworks for determining number of axes to interpret and the significance of each variable on each axis (lots of work on this area).

1st) determine how many axes to interpret (i.e., how many PCs capture correlated variation in the data?).



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Computational Statistics & Data Analysis 49 (2005) 974–997

**COMPUTATIONAL
STATISTICS
& DATA ANALYSIS**

www.elsevier.com/locate/csda

How many principal components? stopping rules for determining the number of non-trivial axes revisited

Pedro R. Peres-Neto*, Donald A. Jackson, Keith M. Somers

Inferential frameworks for determining number of axes to interpret and the significance of each variable on each axis are usually not performed.

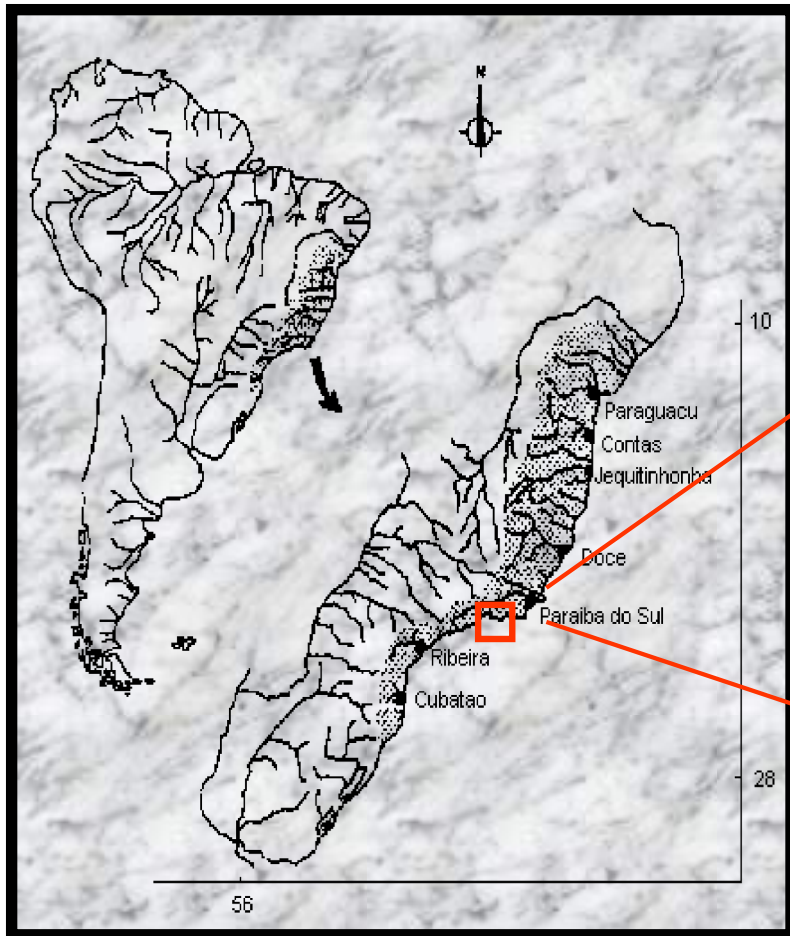
2nd) for each significant axis, determine which variable is significant on each of them.

Ecology, 84(9), 2003, pp. 2347–2363
© 2003 by the Ecological Society of America

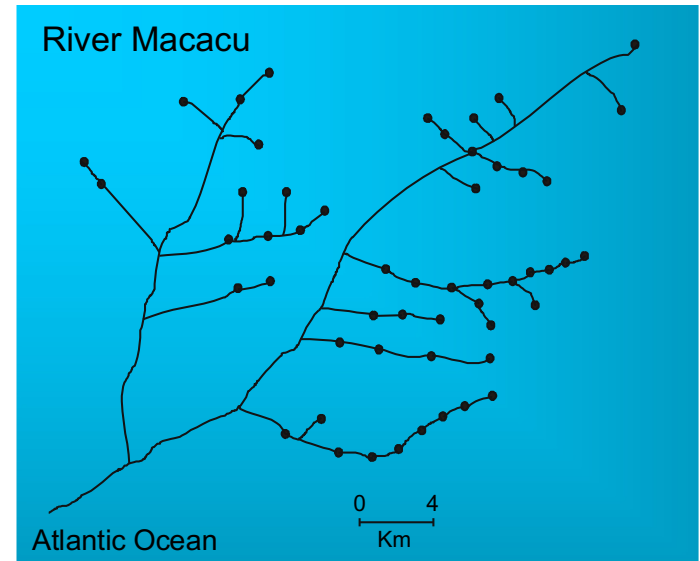
GIVING MEANINGFUL INTERPRETATION TO ORDINATION AXES:
ASSESSING LOADING SIGNIFICANCE IN
PRINCIPAL COMPONENT ANALYSIS

PEDRO R. PERES-NETO,¹ DONALD A. JACKSON, AND KEITH M. SOMERS

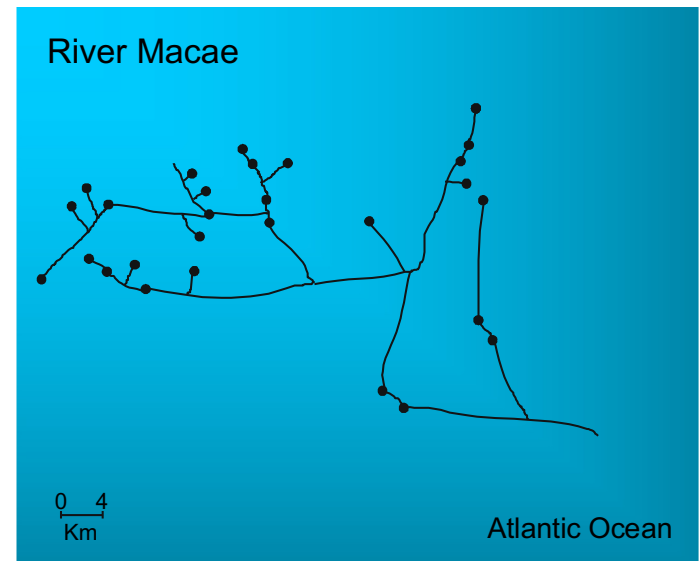
Principal component analysis: a complete example



53 sites



28 sites



What is the correlation structure and differences among streams in terms of their environmental features?

Depth

Depth variation

Current velocity

Current variation

Substrate composition: Boulder, rubble, gravel and sand

Substrate variation (variance in composition)

Stream width variation (irregularity)

Area

Altitude

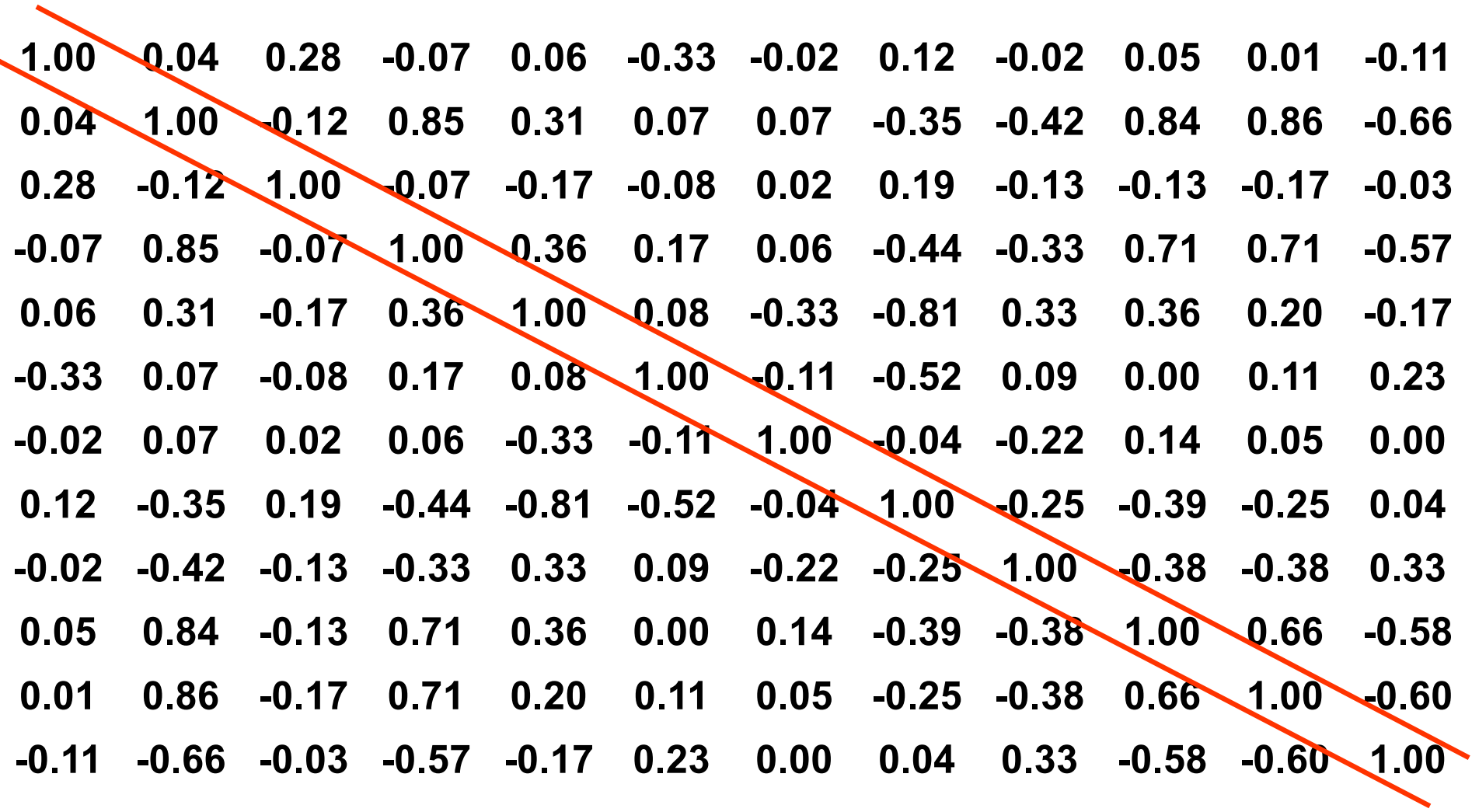
Oecologia (2004) 140: 352–360
DOI 10.1007/s00442-004-1578-3

COMMUNITY ECOLOGY

Pedro R. Peres-Neto

Patterns in the co-occurrence of fish species in streams: the role of site suitability, morphology and phylogeny versus species interactions

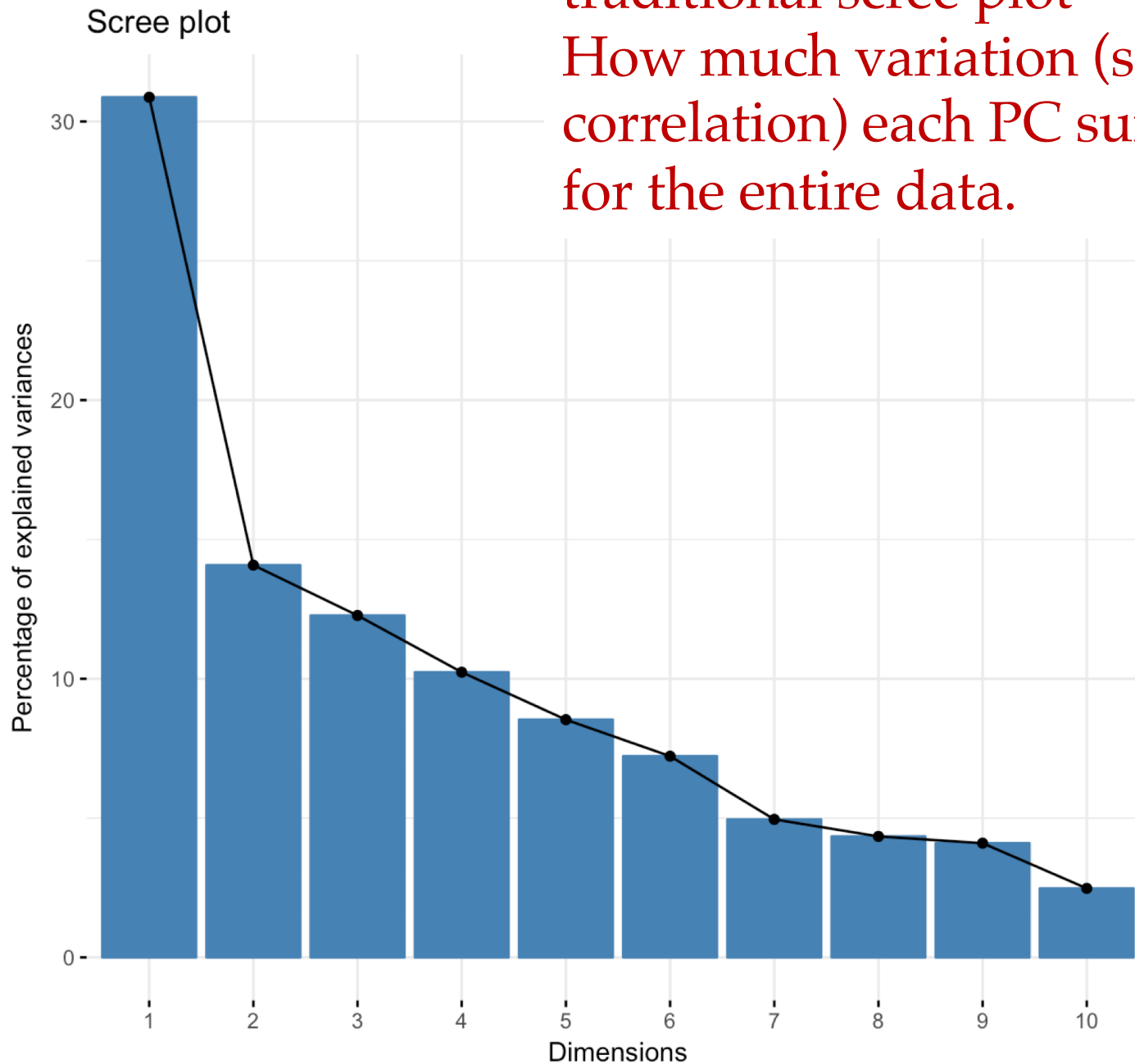
Correlation matrix



A 12x12 correlation matrix is displayed. The diagonal elements are all 1.00. The matrix is symmetric. A red diagonal line is drawn across the matrix, starting from the top-left corner and ending at the bottom-right corner.

1.00	0.04	0.28	-0.07	0.06	-0.33	-0.02	0.12	-0.02	0.05	0.01	-0.11
0.04	1.00	-0.12	0.85	0.31	0.07	0.07	-0.35	-0.42	0.84	0.86	-0.66
0.28	-0.12	1.00	-0.07	-0.17	-0.08	0.02	0.19	-0.13	-0.13	-0.17	-0.03
-0.07	0.85	-0.07	1.00	0.36	0.17	0.06	-0.44	-0.33	0.71	0.71	-0.57
0.06	0.31	-0.17	0.36	1.00	0.08	-0.33	-0.81	0.33	0.36	0.20	-0.17
-0.33	0.07	-0.08	0.17	0.08	1.00	-0.11	-0.52	0.09	0.00	0.11	0.23
-0.02	0.07	0.02	0.06	-0.33	-0.11	1.00	-0.04	-0.22	0.14	0.05	0.00
0.12	-0.35	0.19	-0.44	-0.81	-0.52	-0.04	1.00	-0.25	-0.39	-0.25	0.04
-0.02	-0.42	-0.13	-0.33	0.33	0.09	-0.22	-0.25	1.00	-0.38	-0.38	0.33
0.05	0.84	-0.13	0.71	0.36	0.00	0.14	-0.39	-0.38	1.00	0.66	-0.58
0.01	0.86	-0.17	0.71	0.20	0.11	0.05	-0.25	-0.38	0.66	1.00	-0.60
-0.11	-0.66	-0.03	-0.57	-0.17	0.23	0.00	0.04	0.33	-0.58	-0.60	1.00

Eigenvalue contribution – the traditional scree plot - How much variation (sd and correlation) each PC summarizes for the entire data.



Eigenvector structure (2 first dimensions)

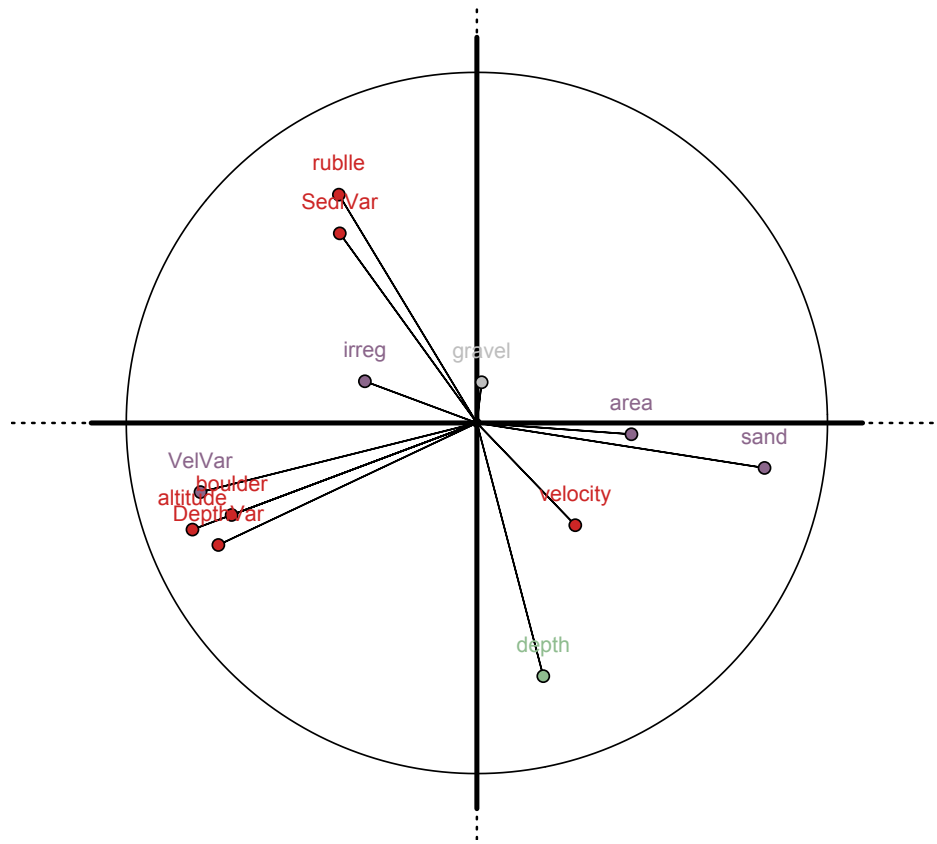
PC-1 PC-2

depth	0.098416371	-0.55557259
DepthVar	-0.383072589	-0.26772556
velocity	0.145820452	-0.22434910
VelVar	-0.409585483	-0.15169873
boulder	-0.363399847	-0.20189977
rubble	-0.204526467	0.50098773
gravel	0.007091107	0.08935752
sand	0.426264131	-0.09866678
altitude	-0.421467330	-0.23396335
area	0.229031867	-0.02477526
irreg	-0.165951470	0.09149688
SediVar	-0.203159109	0.41607768

Eigenvector plot :

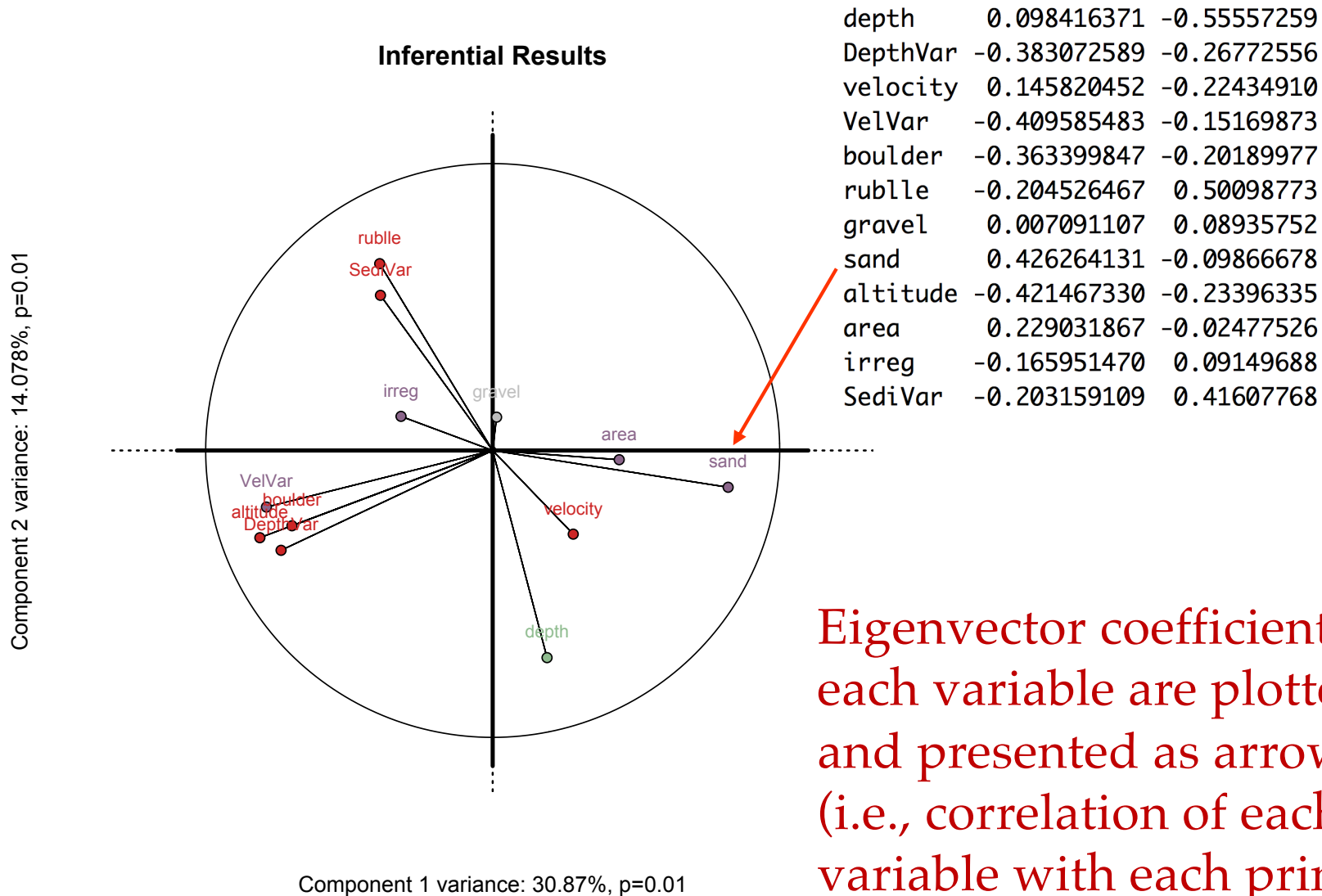
Inferential Results

Component 2 variance: 14.078%, p=0.01



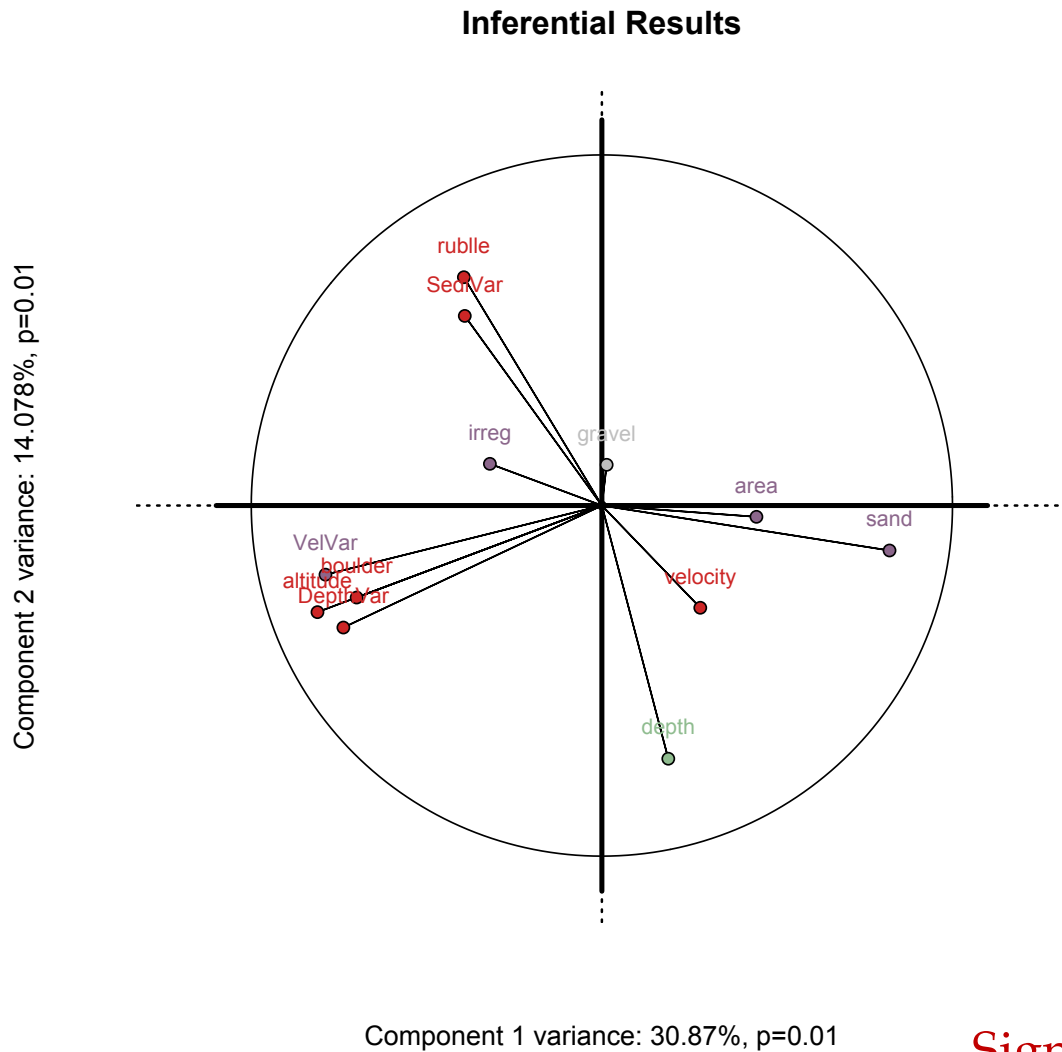
Component 1 variance: 30.87%, p=0.01

Eigenvector plot :

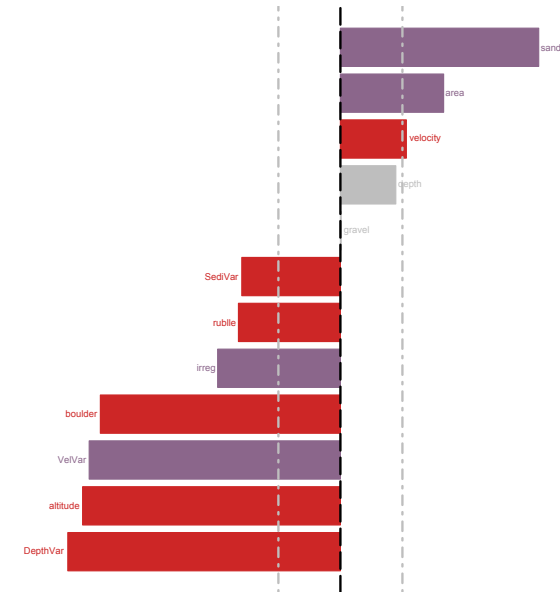


Eigenvector coefficients for each variable are plotted and presented as arrows (i.e., correlation of each variable with each principal component).

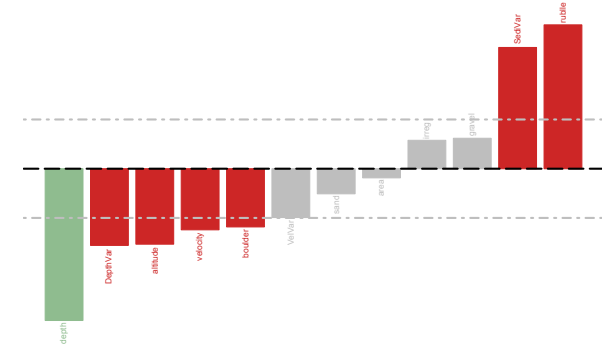
Eigenvector plot:



Bootstrap Ratios Component: 1



Bootstrap Ratios Component: 2

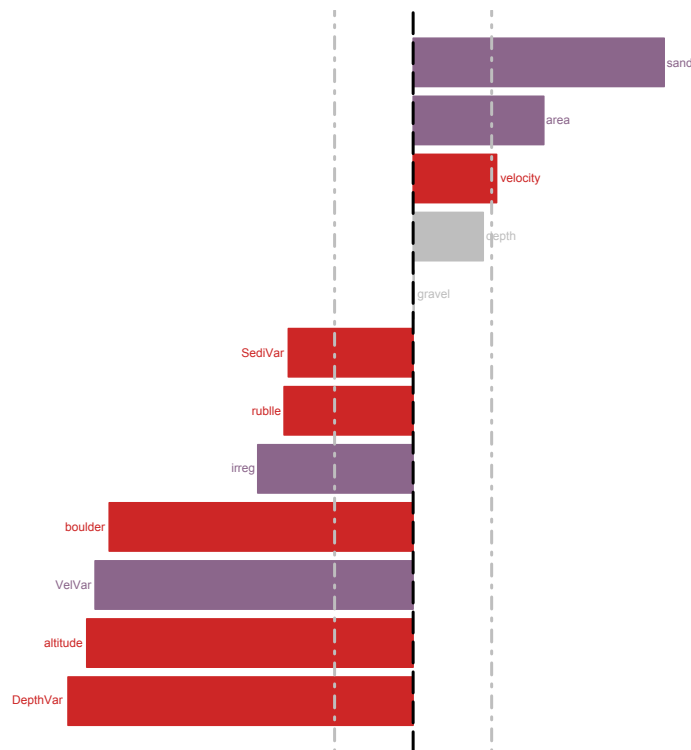


Significant variable contributions are determined if the eigenvector coefficient for the variables are beyond the confidence interval.

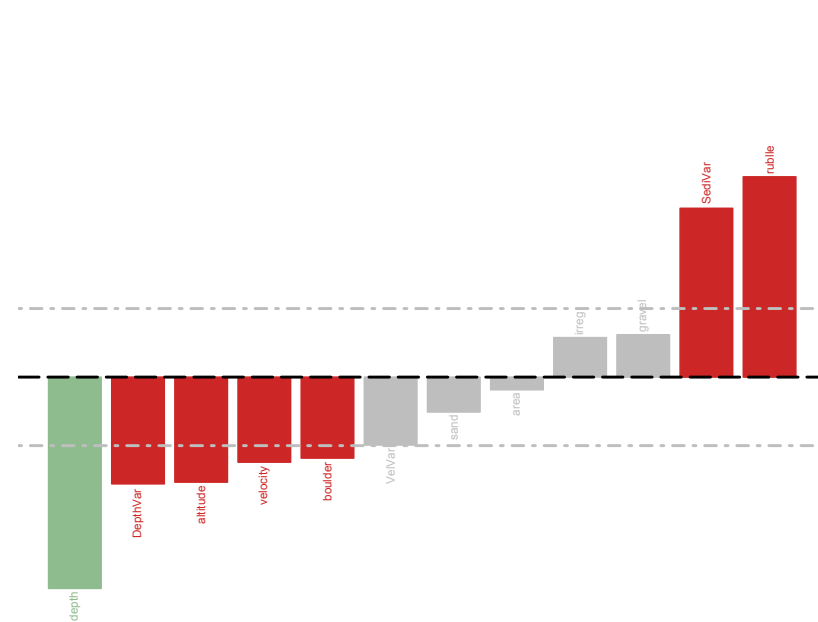
GIVING MEANINGFUL INTERPRETATION TO ORDINATION AXES: ASSESSING LOADING SIGNIFICANCE IN PRINCIPAL COMPONENT ANALYSIS

PEDRO R. PERES-NETO,¹ DONALD A. JACKSON, AND KEITH M. SOMERS

Bootstrap Ratios Component: 1

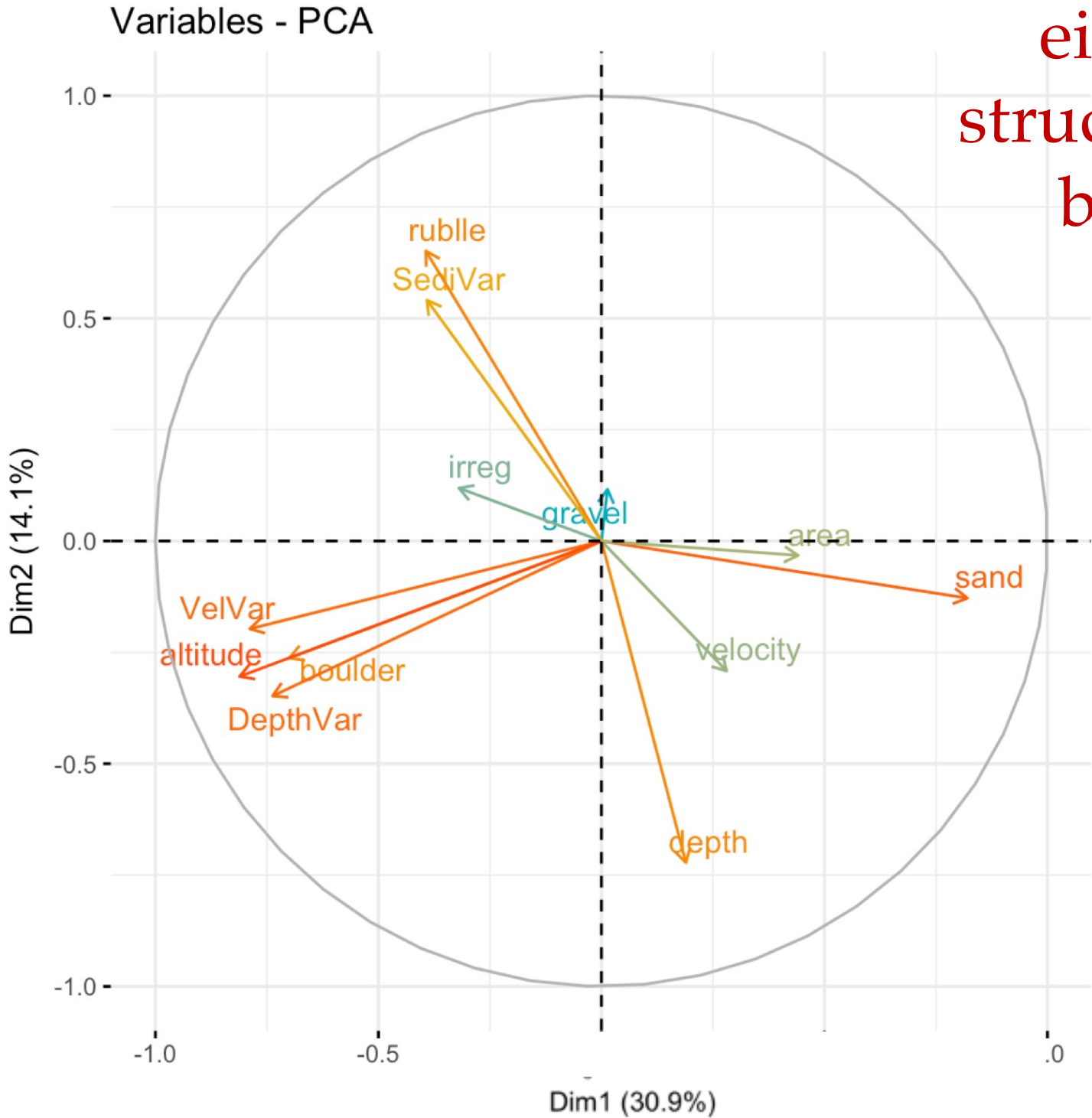


Bootstrap Ratios Component: 2



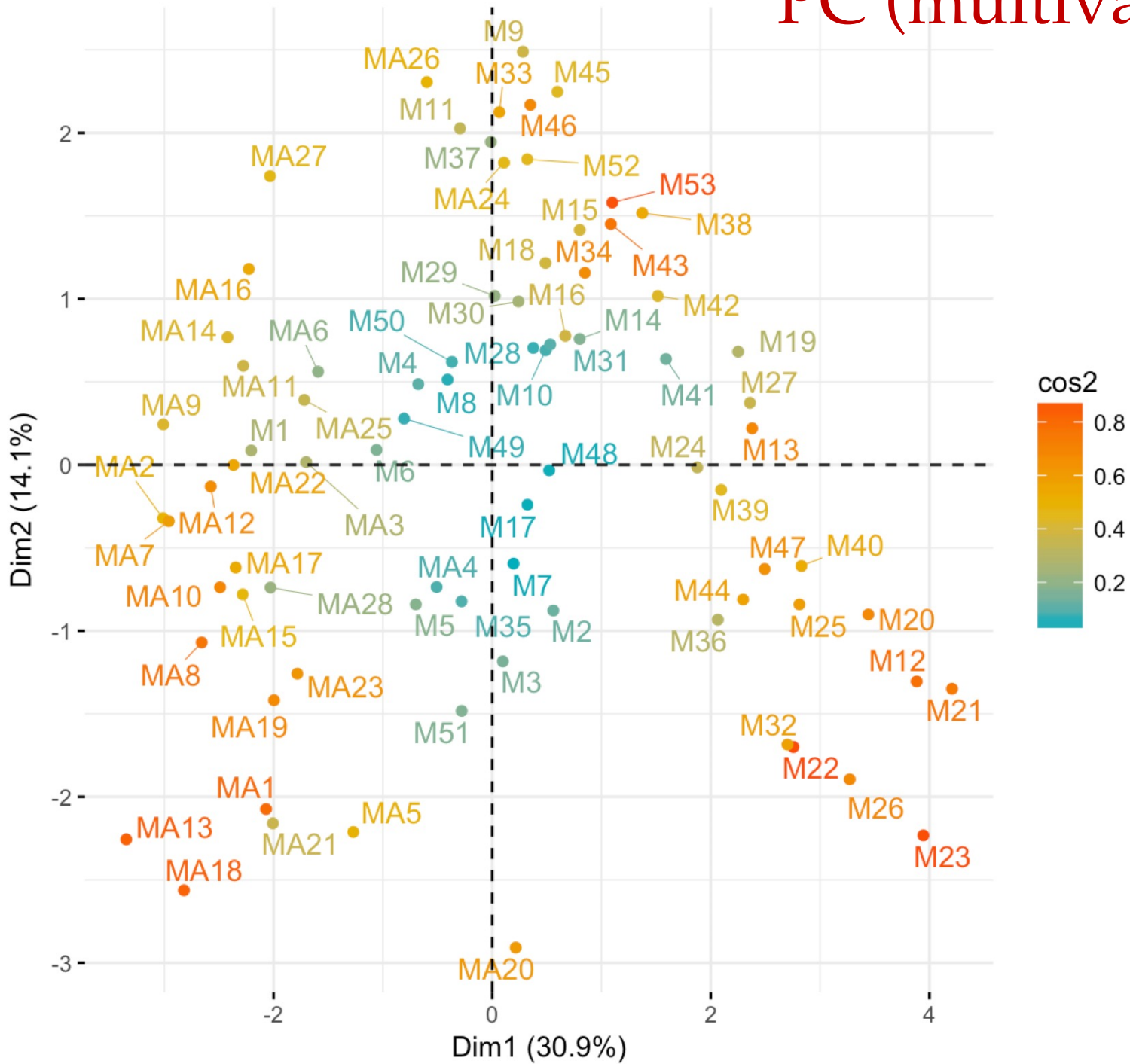
Significant variable contributions are determined if the eigenvector coefficient for the variables are beyond the confidence interval.

eigenvector
structure – much
better plot



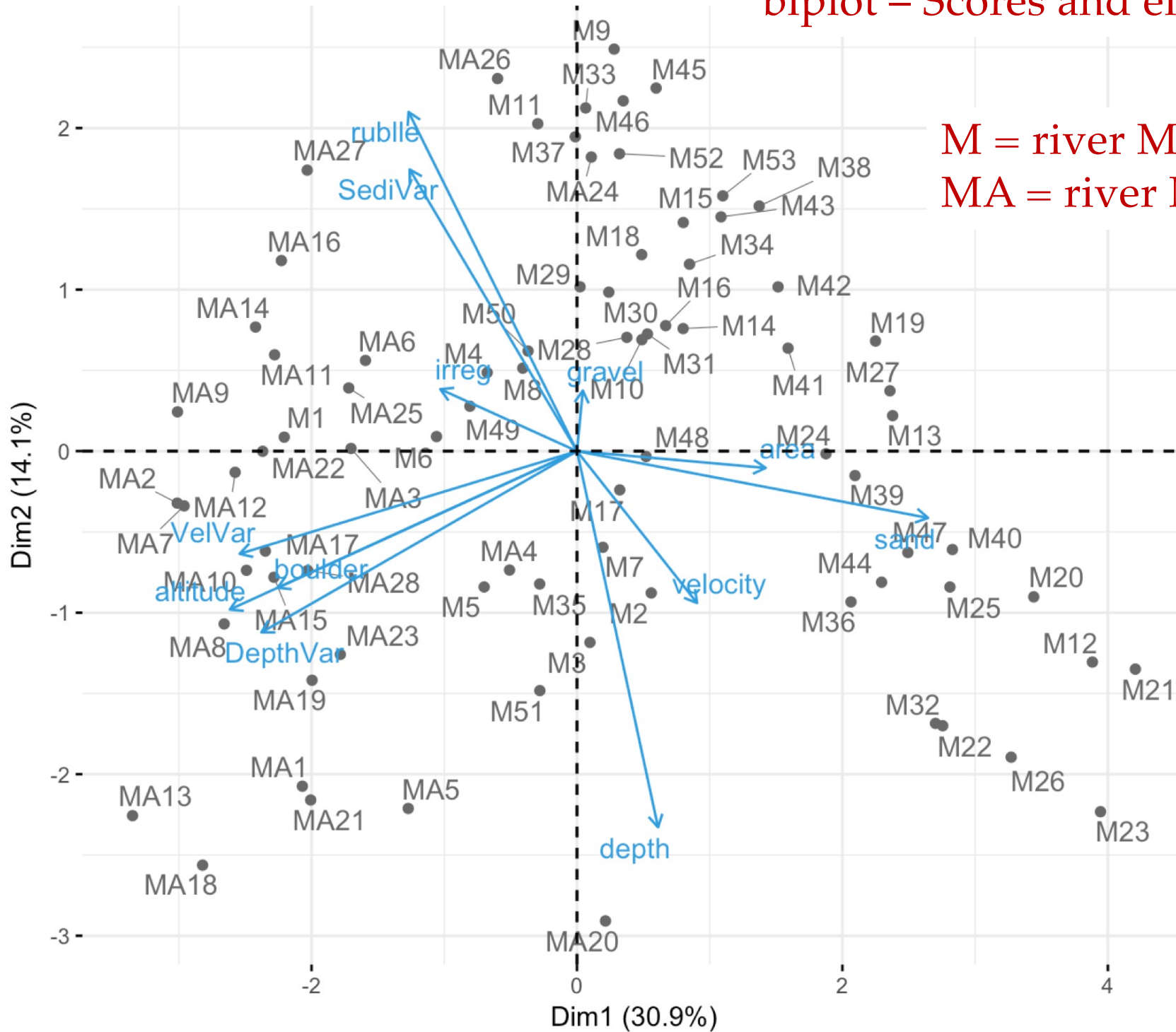
Individuals - PCA

PC (multivariate) scores

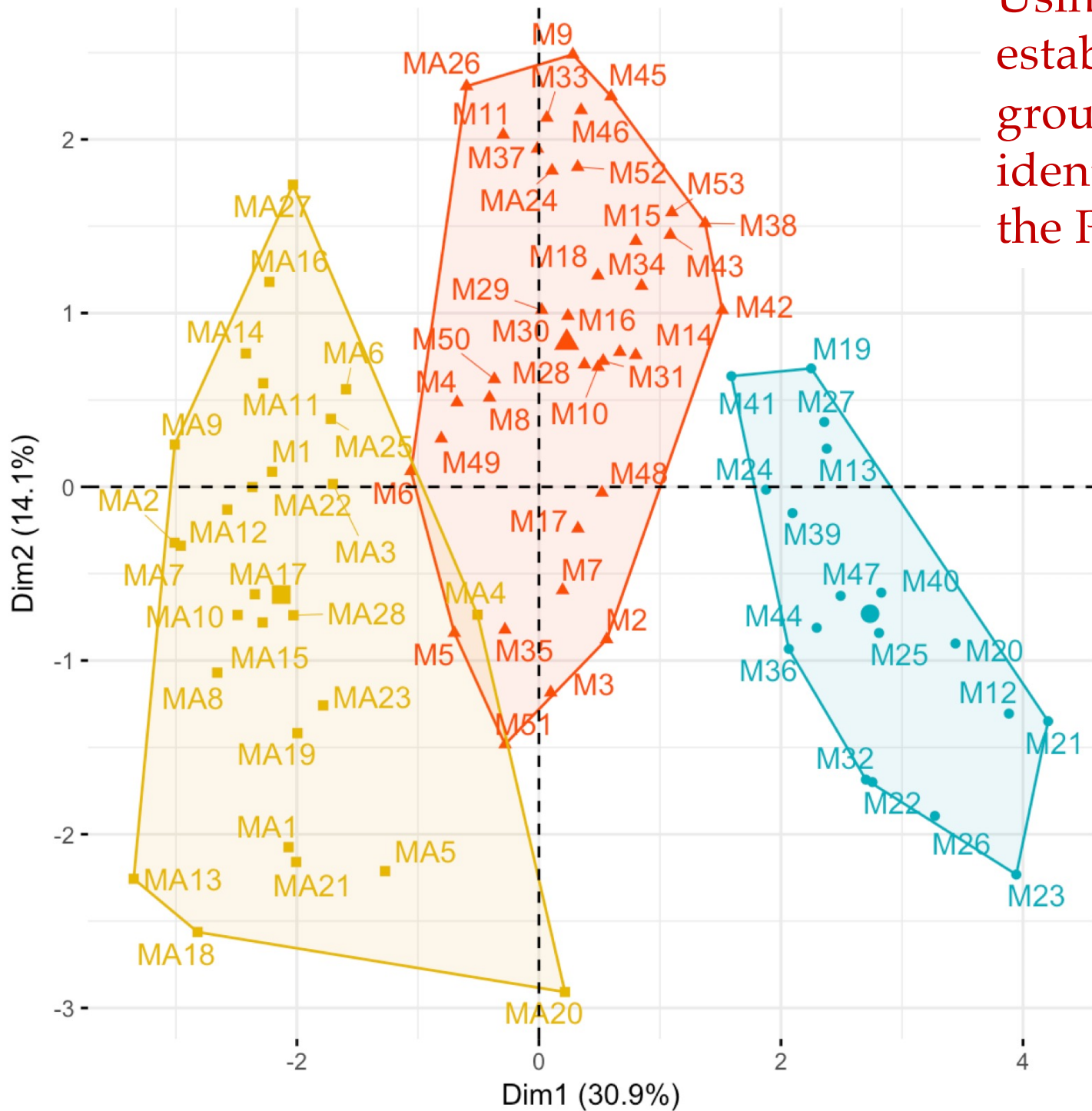


PCA - Biplot

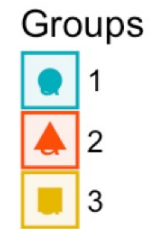
biplot – Scores and eigenvectors



Individuals - PCA



Using K-means to establish number of groups and then identify them the PCA plot



Mapping the environment of our planet – a very “small” example



Long	Lat	avg_prec	avg_ET	avg_VI	avg_Alt	range_Alt	avg_temp	seas_temp	seas_prec
-70.5	-55.344	89.49	222.167	421.958	370.806	2160	38.501	13.97282505	258.8423462
-69.5	-55.344	68.95	241.5	482.354	472.088	2470	32.631	15.6611433	282.420105
-68.5	-55.344	50.23	229	599.458	348.21	1258	38.392	16.68809319	316.756958
-67.5	-55.344	37.9	222.333	623.583	222.572	1047	44.807	22.50125504	323.7763367
-66.5	-55.344	38.94	170.167	498.25	176.965	833	45.774	23.47930336	300.9299011
-70.5	-54.046	47.71	222.167	421.958	174.06	763	54.352	16.45620728	315.1140137
-69.5	-54.046	36.37	241.5	482.354	186.163	786	53.772	17.72191429	331.6278076
-68.5	-54.046	29.06	229	599.458	83.993	342	56.596	20.50442696	354.6487122
-67.5	-54.046	28.69	222.333	623.583	42.762	224	53.089	22.89826965	380.9208679
-73.5	-52.788	149.71	447.167	571.81	287.777	1590	54.122	11.40344238	291.9087524
-72.5	-52.788	54.23	415.333	778.905	267.908	1190	53.38	14.44744396	341.0378723
-71.5	-52.788	26.4	315	867.833	214.992	691	57.626	14.78347588	374.8889771
-70.5	-52.788	18.47	285.167	742.786	148.188	355	61.686	20.54120636	392.6990967
-69.5	-52.788	21.24	158.833	697.405	69.194	263	65.235	23.69093513	386.2383728
-73.5	-51.564	108.16	447.167	562.875	720.814	2785	40.757	11.40344238	291.9087524
-72.5	-51.564	33.43	415.333	761.104	605.359	2000	46.611	14.44744396	341.0378723
-71.5	-51.564	15.04	315	655.229	416.289	861	57.116	14.78347588	374.8889771
-70.5	-51.564	13.15	285.167	622.688	247.667	418	66.142	20.54120636	392.6990967
-69.5	-51.564	15.56	158.833	607.313	161.634	367	69.428	23.69093513	386.2383728
-74.5	-50.373	266.93	512.833	540.292	427.273	2164	60.755	12.20302963	278.6205139
-73.5	-50.373	108.13	353.333	562.875	1127.405	3405	29.841	14.35677052	314.4966431
-72.5	-50.373	31.43	296.333	761.104	568.539	1756	61.505	25.55801201	418.8746643
-71.5	-50.373	13.61	275	655.229	500.653	1080	65.597	26.75129128	476.8866882
-70.5	-50.373	11.79	207.833	622.688	340.328	775	74.516	22.04994202	503.4150391
-69.5	-50.373	13.93	198.167	607.313	171.078	549	83.211	20.19462204	514.8027954
-68.5	-50.373	16.74	134.333	448.188	96.374	463	87.676	22.06239319	502.4994507
-74.5	-49.21	252.75	512.833	547.905	383.44	1567	65.613	8.876086235	275.4307251
-73.5	-49.21	110.05	353.333	497.929	1054.021	3435	36.699	12.14703465	320.0220642
-72.5	-49.21	39.53	296.333	589.905	931.375	1961	48.015	20.98670387	388.9425964
-71.5	-49.21	16.23	275	582.929	768.218	1447	61.56	27.47363091	454.5921631

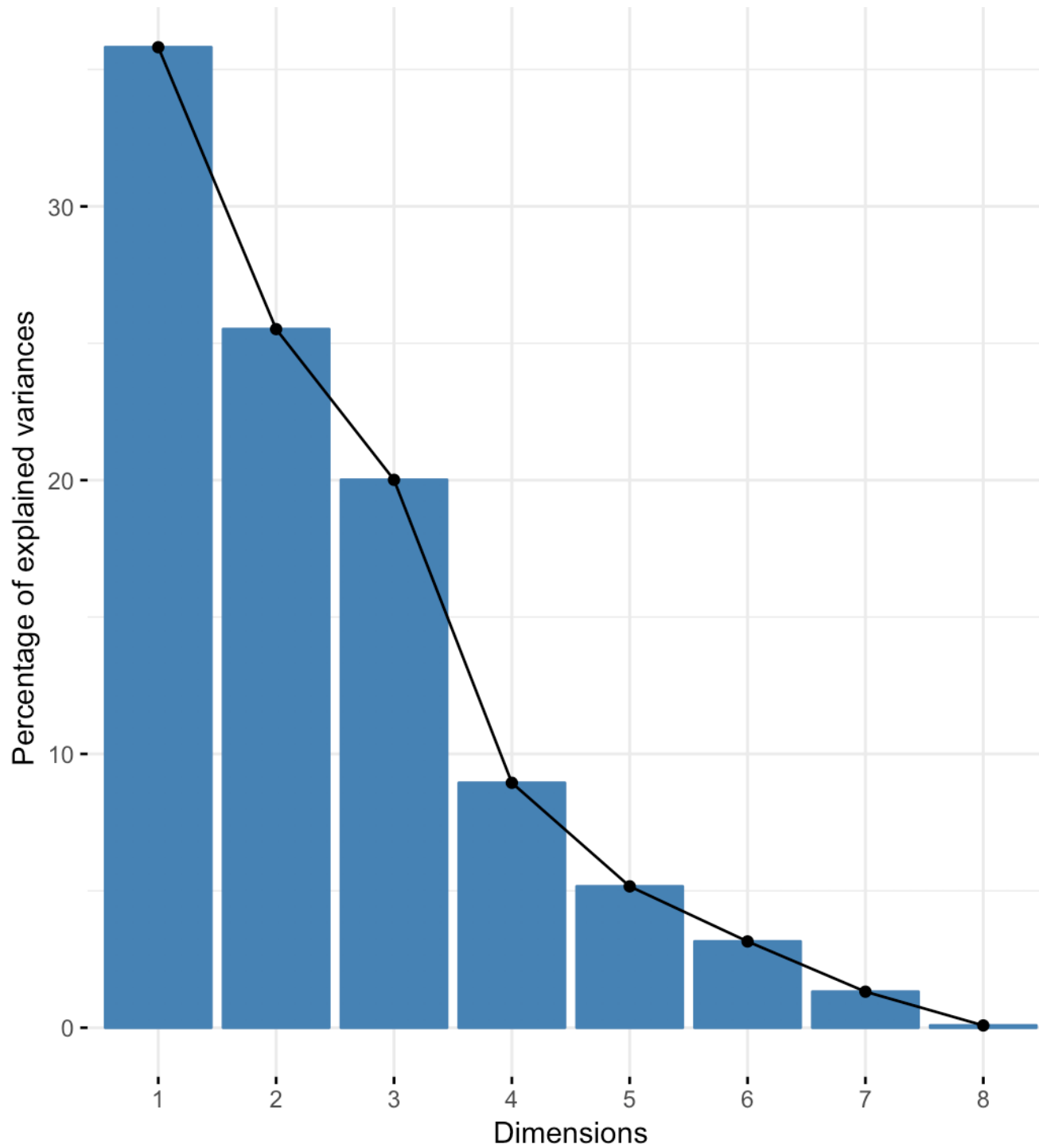


14909 geographic cells (110Km by 110Km)

Long	Lat	avg_prec	avg_ET	avg_VI	avg_Alt	range_Alt	avg_temp	seas_temp	seas_prec
-70.5	-55.344	89.49	222.167	421.958	370.806	2160	38.501	13.97282505	258.8423462
-69.5	-55.344	68.95	241.5	482.354	472.088	2470	32.631	15.6611433	282.420105
-68.5	-55.344	50.23	229	599.458	348.21	1258	38.392	16.68809319	316.756958
-67.5	-55.344	37.9	222.333	623.583	222.572	1047	44.807	22.50125504	323.7763367
-66.5	-55.344	38.94	170.167	498.25	176.965	833	45.774	23.47930336	300.9299011
-70.5	-54.046	47.71	222.167	421.958	174.06	763	54.352	16.45620728	315.1140137
-69.5	-54.046	36.37	241.5	482.354	186.163	786	53.772	17.72191429	331.6278076
-68.5	-54.046	29.06	229	599.458	83.993	342	56.596	20.50442696	354.6487122
-67.5	-54.046	28.69	222.333	623.583	42.762	224	53.089	22.89826965	380.9208679
-73.5	-52.788	149.71	447.167	571.81	287.777	1590	54.122	11.40344238	291.9087524
-72.5	-52.788	54.23	415.333	778.905	267.908	1190	53.38	14.44744396	341.0378723
-71.5	-52.788	26.4	315	867.833	214.992	691	57.626	14.78347588	374.8889771
-70.5	-52.788	18.47	285.167	742.786	148.188	355	61.686	20.54120636	392.6990967
-69.5	-52.788	21.24	158.833	697.405	69.194	263	65.235	23.69093513	386.2383728
-73.5	-51.564	108.16	447.167	562.875	720.814	2785	40.757	11.40344238	291.9087524
-72.5	-51.564	33.43	415.333	761.104	605.359	2000	46.611	14.44744396	341.0378723
-71.5	-51.564	15.04	315	655.229	416.289	861	57.116	14.78347588	374.8889771
-70.5	-51.564	13.15	285.167	622.688	247.667	418	66.142	20.54120636	392.6990967
-69.5	-51.564	15.56	158.833	607.313	161.634	367	69.428	23.69093513	386.2383728
-74.5	-50.373	266.93	512.833	540.292	427.273	2164	60.755	12.20302963	278.6205139
-73.5	-50.373	108.13	353.333	562.875	1127.405	3405	29.841	14.35677052	314.4966431
-72.5	-50.373	31.43	296.333	761.104	568.539	1756	61.505	25.55801201	418.8746643
-71.5	-50.373	13.61	275	655.229	500.653	1080	65.597	26.75129128	476.866882
-70.5	-50.373	11.79	207.833	622.688	340.328	775	74.516	22.04994202	503.4150391
-69.5	-50.373	13.93	198.167	607.313	171.078	549	83.211	20.19462204	514.8027954
-68.5	-50.373	16.74	134.333	448.188	96.374	463	87.676	22.06239319	502.4994507
-74.5	-49.21	252.75	512.833	547.905	383.44	1567	65.613	8.876086235	275.4307251
-73.5	-49.21	110.05	353.333	497.929	1054.021	3435	36.699	12.14703465	320.0220642
-72.5	-49.21	39.53	296.333	589.905	931.375	1961	48.015	20.98670387	388.9425964
-71.5	-49.21	16.23	275	582.929	768.218	1447	61.56	27.47363091	454.5921631



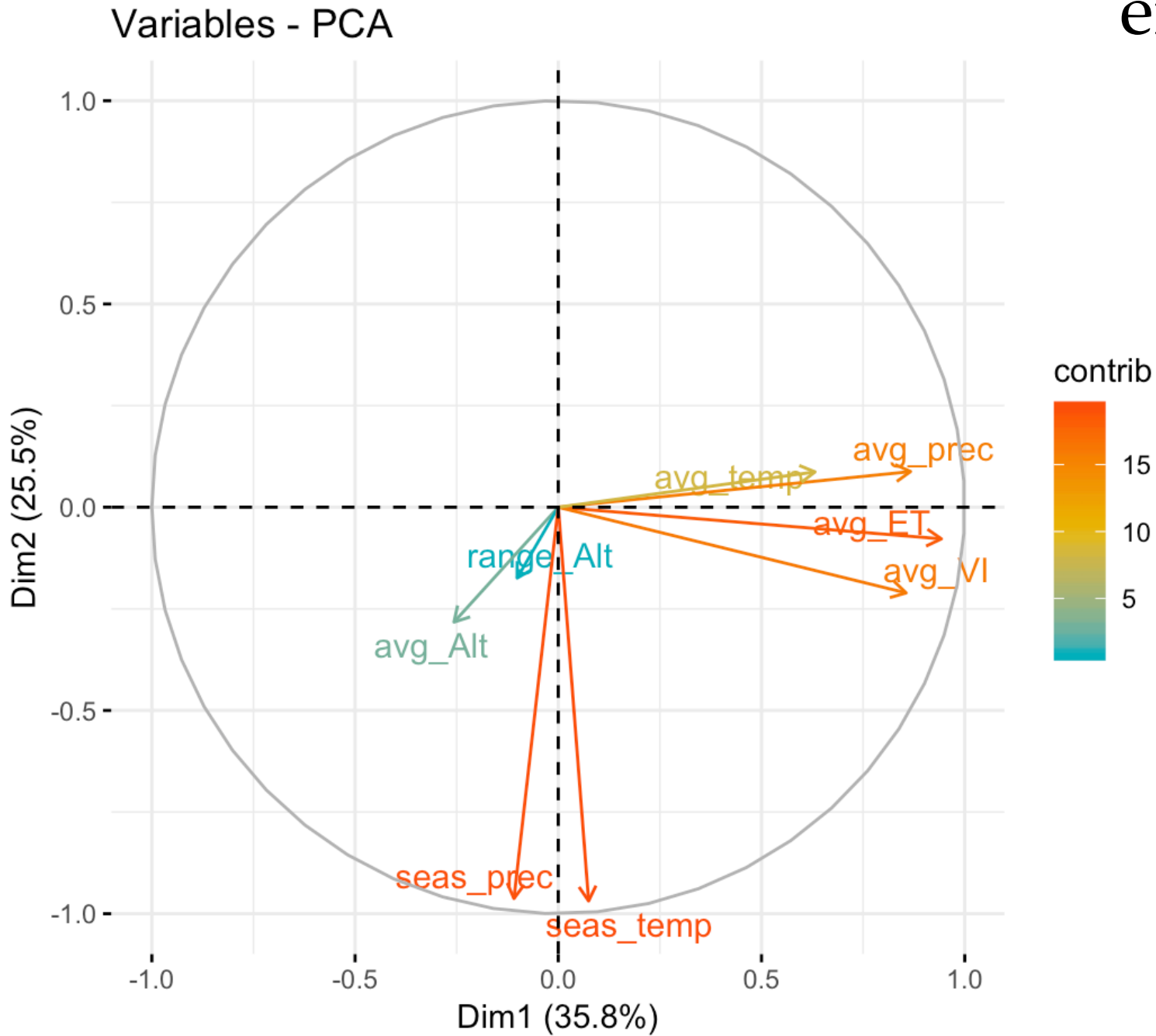
- 1) Latitude (Lat) and Longitude (Long) at the centre of geographic cell.
- 2) Average precipitation (last 40 years; avg_prec)
- 3) Average actual evapotranspiration (avg_ET, a proxy of productivity)
- 4) Average vegetation index (avg_VI)
- 5) Mean altitude (avg_Alt)
- 6) Maximum altitude minus minimum altitude (altitudinal range; range_Alt)
- 7) Average temperature (avg_temp)
- 8) Seasonal temperature (annual range in temperature; seas_temp)
- 9) Seasonal precipitation (annual range in precipitation; seas_prec)

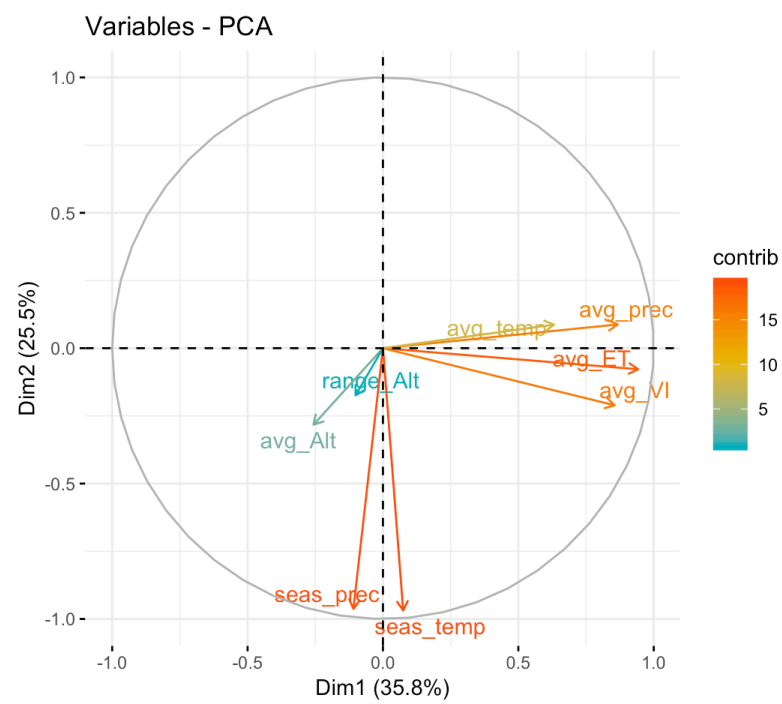
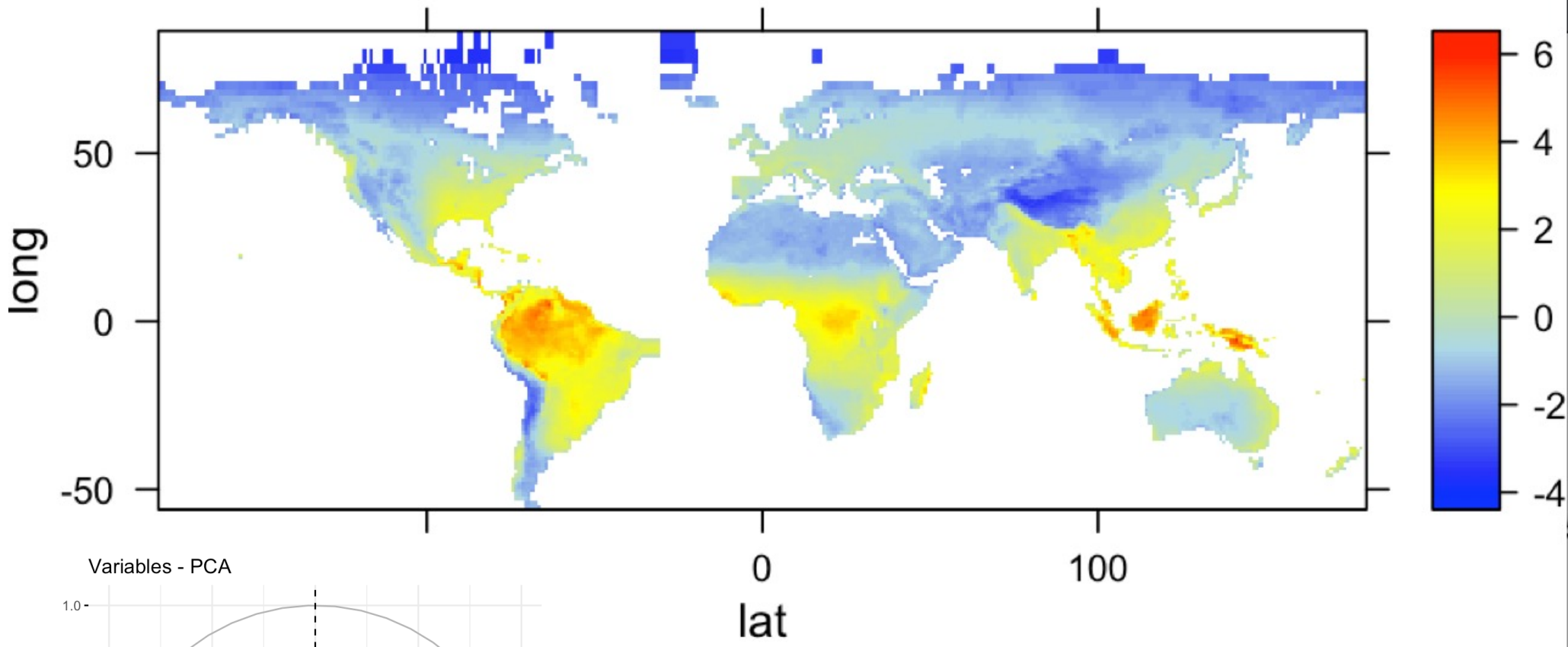


Eigenvalue
contribution –
the traditional
scree plot



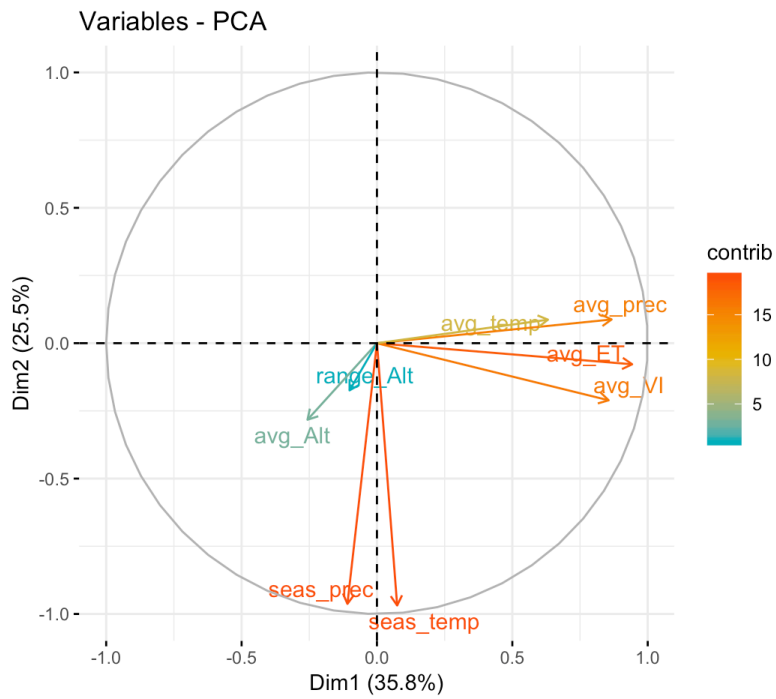
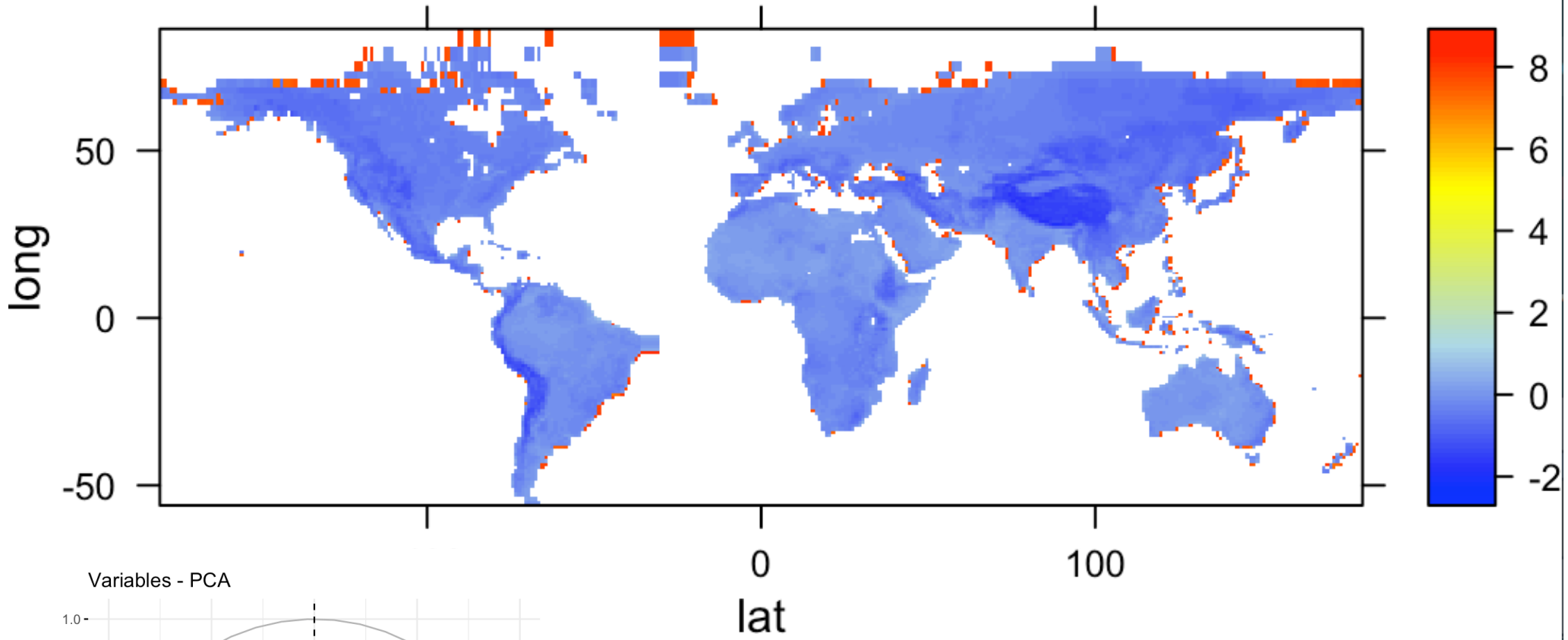
eigenvector structure





Map of PC-1





Map of PC-2

