

Introduction to multivariate analysis, ordination and PCA

Contents

- 1. Learning targets and introduction to multivariate analysis**
2. Overview ordination
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Learning targets

- Explain the specifics of multivariate analysis
- List ordination methods and select one based on research goal and question
- Explain the mathematical basis of and apply PCA
- Interpreting results from a PCA

Learning targets and study questions

- Explain the specifics of multivariate analysis
 - Why should you favour multivariate approaches for multivariate data?
 - Outline the differences to the univariate case when diagnosing multivariate outliers and normality.
- List ordination methods and select one based on research goal and question
 - Outline the aim of ordination
 - Distinguish constrained and unconstrained ordination
 - Which criteria should guide the selection of an ordination method?

Learning targets and study questions

- Explain the mathematical basis of and apply PCA
 - Explain the variance-covariance matrix.
 - What are eigenvalues and eigenvectors?
 - How do eigenvalues relate to the variance of a PC?
 - Outline criteria to determine the optimal number of PCs.
 - What is sparse PCA and how does it influence the evaluation of descriptor contribution to PCs?
- Interpreting results from a PCA
 - Explain biplots with respect to (a) correlation between variables and (b) relation between sites/species and variables.
 - How does scaling influence interpretation of a biplot?
 - Which objects from a PCA would be extracted as non-collinear predictors for a multiple regression analysis?

From univariate to multivariate statistics

	univariate	multivariate
Variables (vars.)	Single/multiple predictors, single response variable Y	Single/multiple predictors, multiple response vars. Y_1, \dots, Y_n
Distribution of response	One-dimensional	n -dimensional
Data format	Y is vector	Y_1, \dots, Y_n constitute matrix
Example	Species richness explained by environmental variables	Community explained by environmental variables

Multivariate data analysis: Introduction

Some advantages of multivariate over univariate methods for analysing multivariate data

- Not all research questions can be answered with univariate statistical methods
e.g. What are the most important environmental variables determining community composition?
- Multivariate methods allow for dimension reduction and visualisation of multidimensional data
e.g. Ordination, Cluster dendrogram
- Joint (multivariate) analysis can reduce noise and increase power when assessing statistical hypotheses

Multivariate approaches in R

Available methods and developments: CRAN Task View

CRAN Task View: Multivariate Statistics

Maintainer: Paul Hewson

Contact: Paul.Hewson at plymouth.ac.uk

Version: 2018-07-21

URL: <https://CRAN.R-project.org/view=Multivariate>

Base R contains most of the functionality for classical multivariate analysis, somewhere. There are a large number of packages on CRAN which extend this methodology, a brief overview is given below. Application-specific uses of multivariate statistics are described in relevant task views, for example whilst principal components are listed here, ordination is covered in the [Environmetrics](#) task view. Further information on supervised classification can be found in the [MachineLearning](#) task view, and unsupervised classification in the [Cluster](#) task view.

The packages in this view can be roughly structured into the following topics. If you think that some package is missing from the list, please let me know.

Visualising multivariate data

- *Graphical Procedures:* A range of base graphics (e.g. `pairs()` and `coplot()`) and [lattice](#) functions (e.g. `xyplot()` and `spplot()`) are useful for visualising pairwise arrays of 2-dimensional scatterplots, clouds and 3-dimensional densities. `scatterplot.matrix` in the [car](#) provides usefully enhanced pairwise scatterplots. Beyond this, [scatterplot3d](#) provides 3 dimensional scatterplots, [aplpack](#) provides bagplots and `spin3d()`, a function for rotating 3d clouds. [misc3d](#), dependent upon [rgl](#), provides animated functions within R useful for visualising densities. [YaleToolkit](#) provides a range of useful visualisation techniques for multivariate data. More specialised multivariate plots include the following: `faces()` in [aplpack](#) provides Chernoff's faces; `parcoord()` from [MASS](#) provides parallel coordinate plots; `stars()` in [graphics](#) provides a choice of star, radar and cobweb plots respectively. `mstree()` in [ade4](#) and `spantree()` in [vegan](#) provide minimum spanning tree functionality. [calibrate](#) supports biplot and scatterplot axis labelling. [geometry](#), which provides an interface to the `qhull` library, gives indices to the relevant points via `convexhulln()`. [ellipse](#) draws ellipses for two parameters, and provides `plotcorr()`, visual display of a correlation matrix. [denpro](#) provides level set trees for multivariate visualisation. Mosaic plots are available via `mosaicplot()` in [graphics](#) and `mosaic()` in [vcd](#) that also contains other visualization techniques for multivariate categorical data. [gclus](#) provides a number of cluster specific graphical enhancements for scatterplots and parallel coordinate plots See the links for a reference to GGobi. [rggobi](#) interfaces with GGobi. [xgobi](#) interfaces to the XGobi and XGvis programs which allow linked, dynamic multivariate plots as well as projection pursuit. Finally, [iplots](#) allows particularly powerful dynamic interactive graphics, of which interactive parallel coordinate plots and mosaic plots may be of great interest. Seriation methods are provided by [seriation](#) which can reorder matrices and dendrograms.
- *Data Preprocessing:* `summarize()` and `summary.formula()` in [Hmisc](#) assist with descriptive functions; from the same package `varclus()` offers variable clustering while `dataRep()` and `find.matches()` assist in exploring a given dataset in terms of representativeness and finding matches. Whilst `dist()` in base and `daisy()` in [cluster](#) provide a wide range of distance measures, [proxy](#) provides a framework for more distance measures, including measures between matrices. [simba](#) provides functions for dealing with presence / absence data including similarity matrices and reshaping.

Hypothesis testing

- [ICSNP](#) provides Hotellings T2 test as well as a range of non-parametric tests including location tests based on marginal ranks, spatial median and spatial signs computation, estimates of shape. Non-parametric two sample tests are also available from [cramer](#) and spatial sign and rank tests to investigate location, sphericity and independence are available in [SpatialNP](#).

Multivariate distributions

- *Descriptive measures:* `cov()` and `cor()` in stats will provide estimates of the covariance and correlation matrices respectively. [ICSNP](#) offers several descriptive measures such as `spatial.median()` which provides an estimate of the spatial median and further functions which provide estimates of scatter. Further robust methods are provided such as `cov.rob()` in [MASS](#) which provides robust estimates of the variance-covariance matrix by minimum volume ellipsoid, minimum covariance determinant or classical product-moment. [covRobust](#) provides robust covariance estimation via nearest neighbor variance estimation. [robustbase](#) provides robust covariance estimation via fast minimum covariance determinant with `covMCD()` and the Orthogonalized pairwise estimate of Gnanadesikan-Kettenring via `covOGK()`. Scalable robust methods are provided within [rcov](#) also using fast minimum covariance determinant with `covMcd()` as well as M-estimators with `covMest()`. [corpcor](#) provides shrinkage estimation of large scale covariance and (partial) correlation matrices.
- *Densities (estimation and simulation):* `mvnrm()` in [MASS](#) simulates from the multivariate normal distribution. [mvtnorm](#) also provides simulation as well as probability and quantile functions for both the multivariate t distribution and multivariate normal distributions as well as density functions for the multivariate normal distribution. [mnormt](#) provides multivariate normal and multivariate t density and distribution functions as well as random number simulation. [sn](#) provides density, distribution and random number generation for the multivariate skew normal and skew t distribution. [delt](#) provides a range of functions for estimating multivariate densities by CART and greedy methods. Comprehensive information on mixtures is given in the [Cluster](#) view, some density estimates and random numbers are provided by `rmvnorm.mixt()` and `dmvnorm.mixt()` in [ks](#), mixture fitting is also provided within [bayesm](#). Functions to simulate from the Wishart distribution are provided in a number of places, such as `rwishart()` in [bayesm](#) and `rwish()` in [MCMCpack](#) (the latter also has a density function `dwish()`). `bkd2d()` from [KernSmooth](#) and `kde2d()` from [MASS](#) provide binned and non-binned 2-dimensional kernel density estimation, [ks](#) also provides multivariate kernel smoothing as does [ash](#) and [GenKern](#). [prim](#) provides patient rule induction methods to attempt to find regions of high density in high dimensional multivariate data, [feature](#) also provides methods for determining feature significance in multivariate data (such as in relation to local modes).
- *Assessing normality:* [mvnormtest](#) provides a multivariate extension to the Shapiro-Wilks test, [myoutlier](#) provides multivariate outlier detection based on robust methods. [ICS](#) provides tests for multi-normality. `mvnorm.etest()` in [energy](#) provides an assessment of normality based on E statistics (energy); in the same package `k.sample()` assesses a number of samples for equal distributions. Tests for Wishart-distributed covariance matrices are given by `mauchly.test()` in stats.
- *Copulas:* [copula](#) provides routines for a range of (elliptical and archimedean) copulas including normal, t, Clayton, Frank, Gumbel, [fgac](#) provides generalised archimedean copula.

Multivariate outlier checking

Multivariate and univariate exploration/diagnosis similar for e.g. multicollinearity and transformation, but approaches differ for checking of e.g. outliers and distributional assumptions

Multivariate outliers

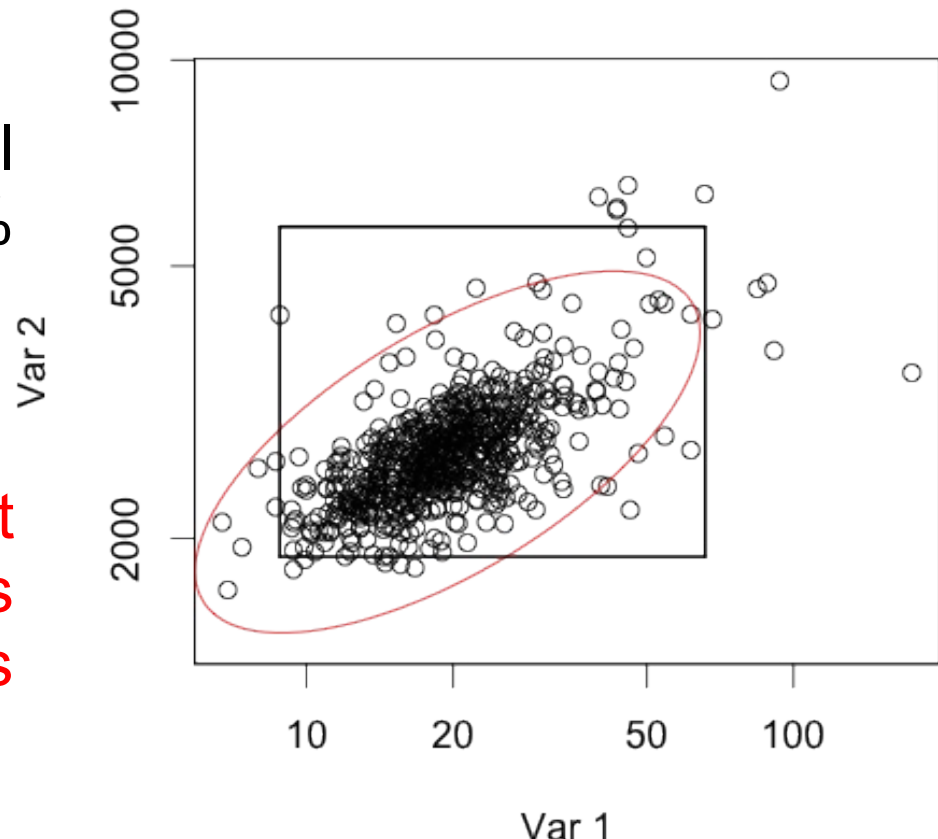
Use joint multivariate distribution of variables to find outliers (i.e. extreme points from centre of multivariate sample)

Outside of box:

Outliers if considering individual distributions of variables (99% quantile of each variable)

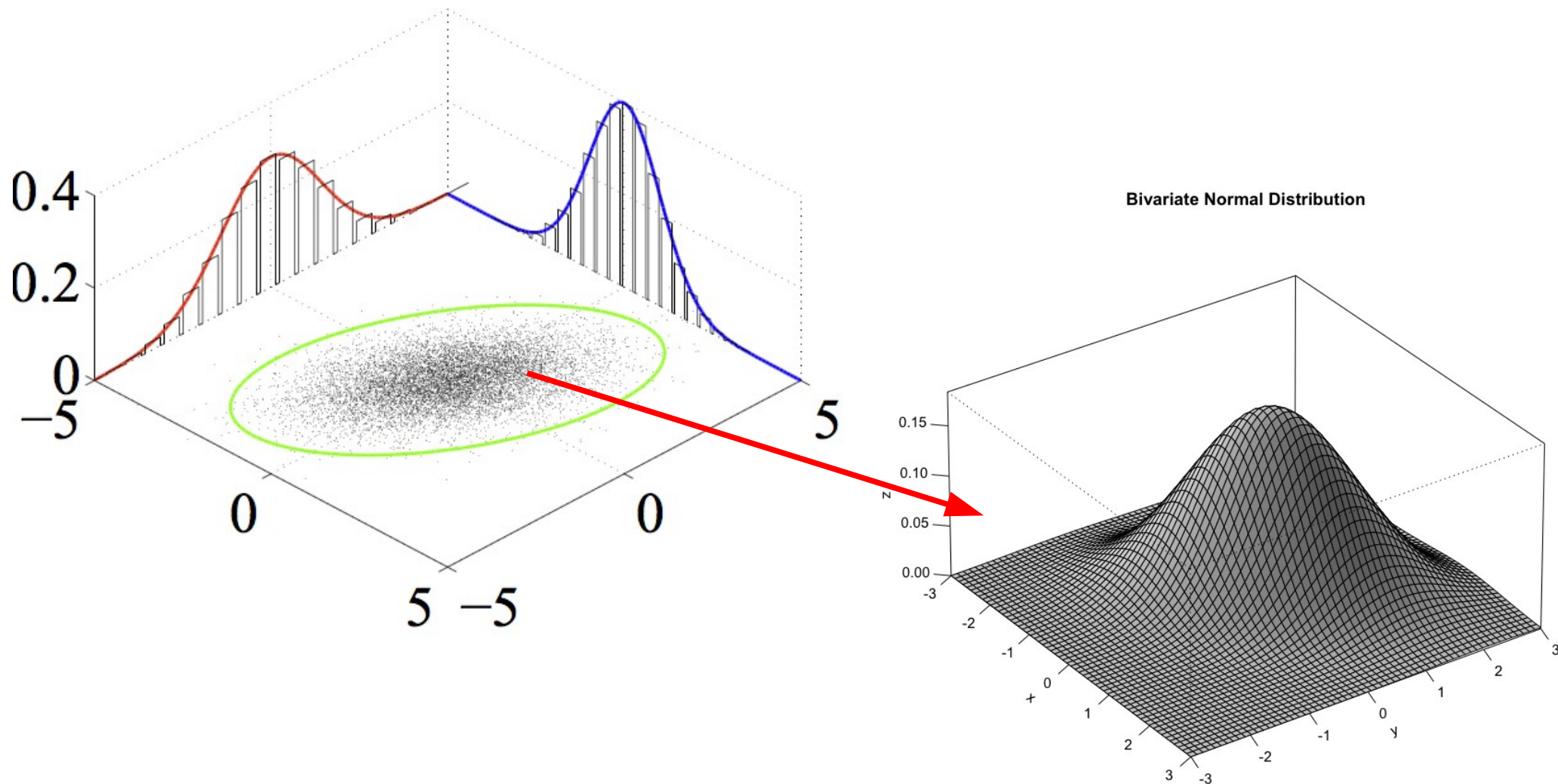
Outside of ellipse:

Outliers if considering joint distribution of both variables (99% quantile using Mahalanobis distance)



Multivariate normal distribution

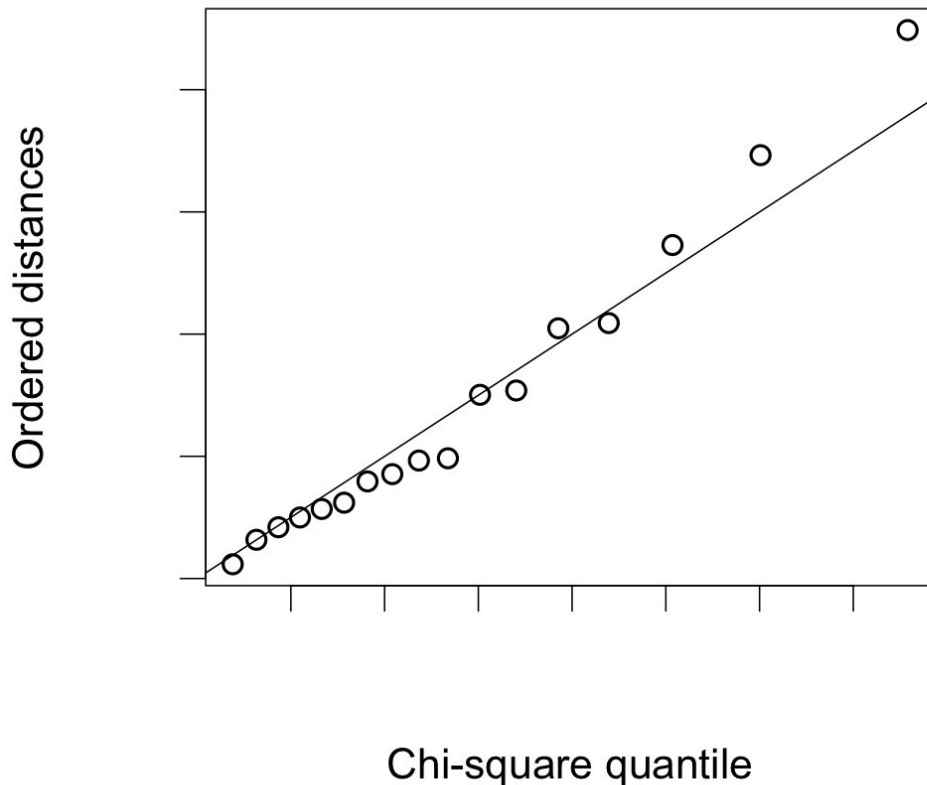
Multivariate normal distribution is evaluated instead of univariate normal distribution



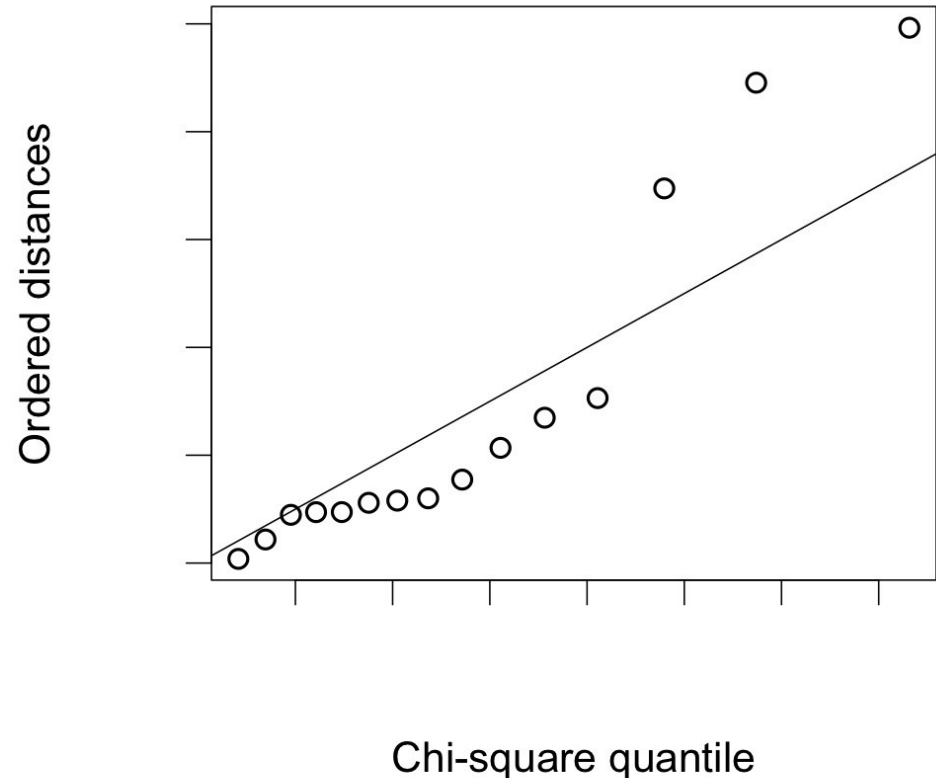
Multivariate normal distribution

Visual check of multivariate normality with QQ plots for sample Mahalanobis distances (to centroid) and theoretical quantiles from the χ^2 distribution

QQ plot with data sampled
from a multivariate normal distribution



QQ plot with data sampled
from an exponential distribution



Covariance matrix

For a data matrix \mathbf{Y} containing the variables y_1, \dots, y_p

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}'_1 \\ \mathbf{y}'_2 \\ \vdots \\ \mathbf{y}'_i \\ \vdots \\ \mathbf{y}'_n \end{pmatrix} = (\text{units}) \begin{matrix} & \begin{matrix} \text{(variables)} \\ 1 & 2 & \cdots & j & \cdots & p \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} & \begin{pmatrix} y_{11} & y_{12} & \cdots & y_{1j} & \cdots & y_{1p} \\ y_{21} & y_{22} & \cdots & y_{2j} & \cdots & y_{2p} \\ \vdots & \vdots & & \vdots & & \vdots \\ y_{i1} & y_{i2} & \cdots & y_{ij} & \cdots & y_{ip} \\ \vdots & \vdots & & \vdots & & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{nj} & \cdots & y_{np} \end{pmatrix} \end{matrix}$$

the sample covariance matrix for these variables is \mathbf{S} (Sigma)

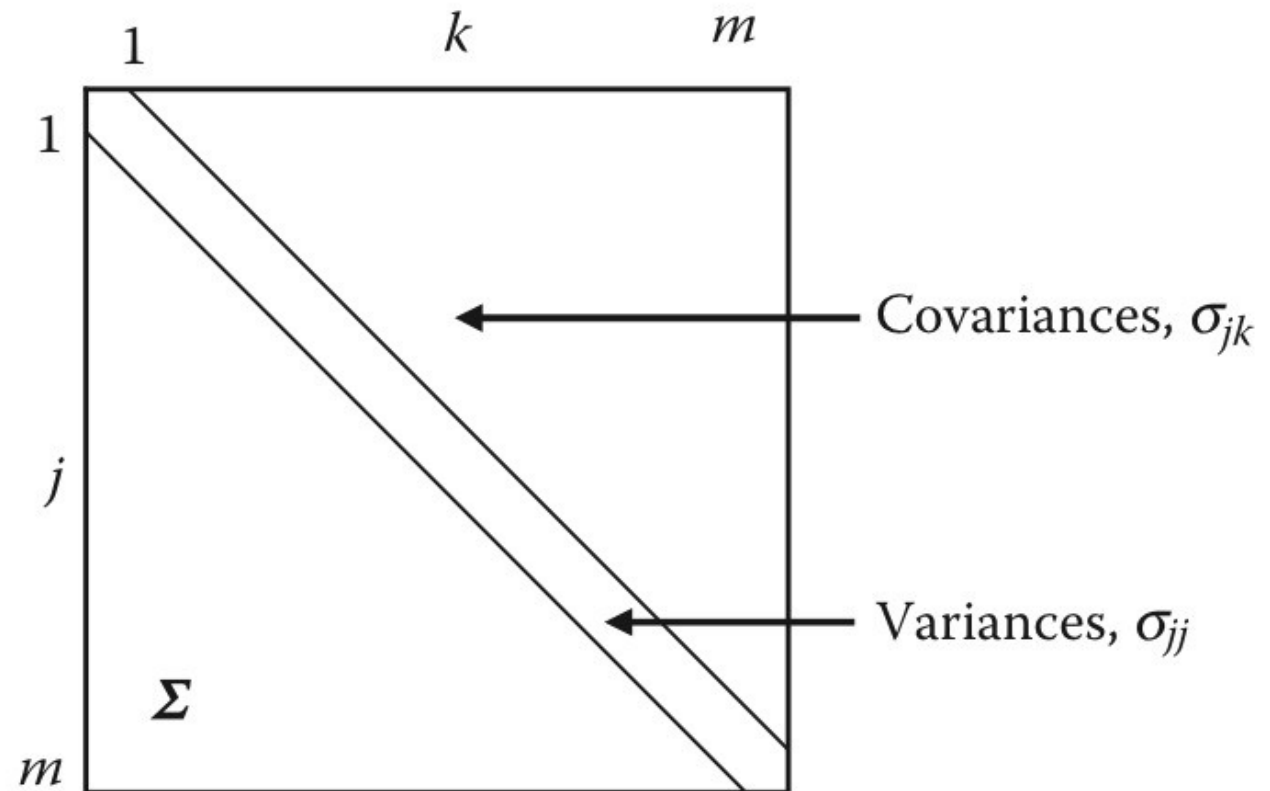
$$\mathbf{S} = (s_{jk}) = \begin{pmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{21} & s_{22} & \cdots & s_{2p} \\ \vdots & \vdots & & \vdots \\ s_{p1} & s_{p2} & \cdots & s_{pp} \end{pmatrix}$$

where
$$s_{jk} = \frac{1}{n-1} \sum_{i=1}^n (y_{i,j} - \bar{y}_j)(y_{i,k} - \bar{y}_k)$$

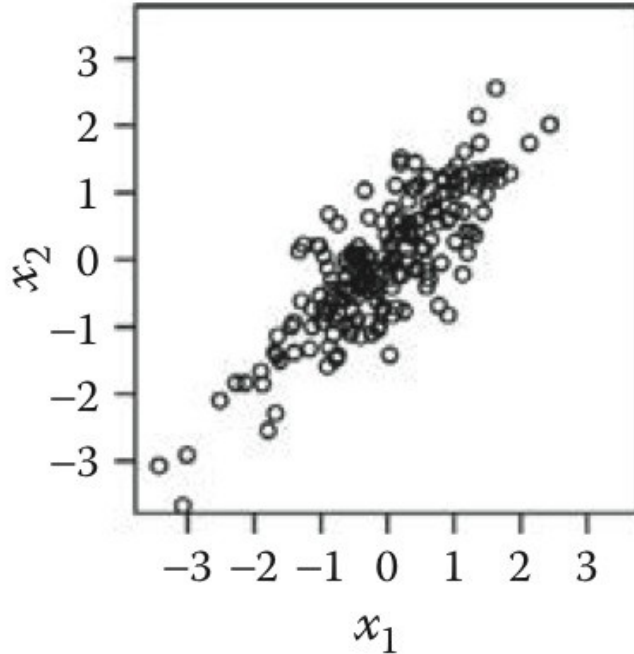
Covariance matrix

In the diagonal, the equation simplifies to:

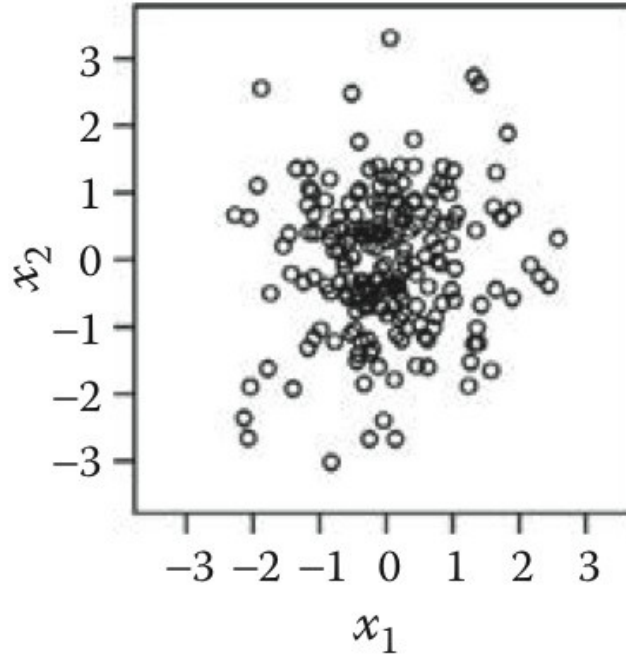
$$s_{jj} = \frac{1}{n-1} \sum_{i=1}^n (y_{i,j} - \bar{y}_j)^2$$



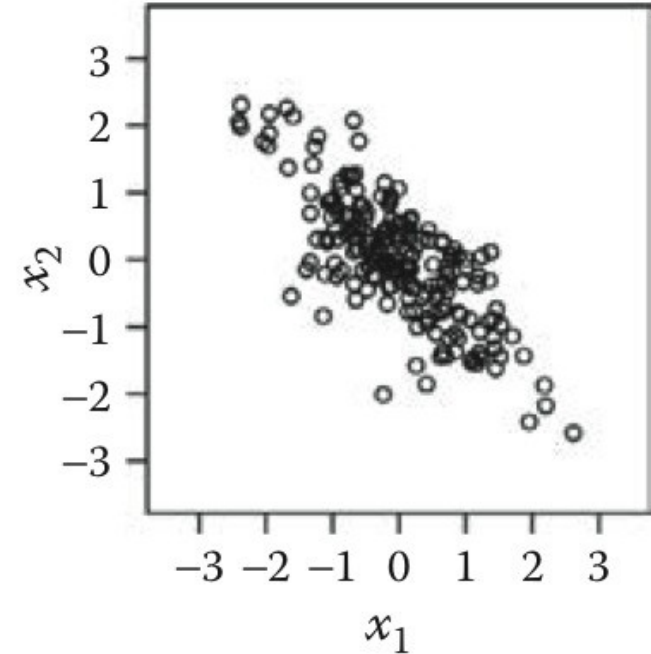
Covariance matrix for two variables



$$\Sigma_1 = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$



$$\Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma_3 = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}$$

Multivariate distance from mean to observation is measured with the Mahalanobis distance D_M :

$$D_M(x) = \sqrt{(x - \mu)^T \mathbf{S}^{-1} (x - \mu)}$$

→ D_M is distance of vector x from the mean vector μ weighed by their covariance (given in sample covariance matrix \mathbf{S})

Introduction to multivariate analysis, ordination and PCA

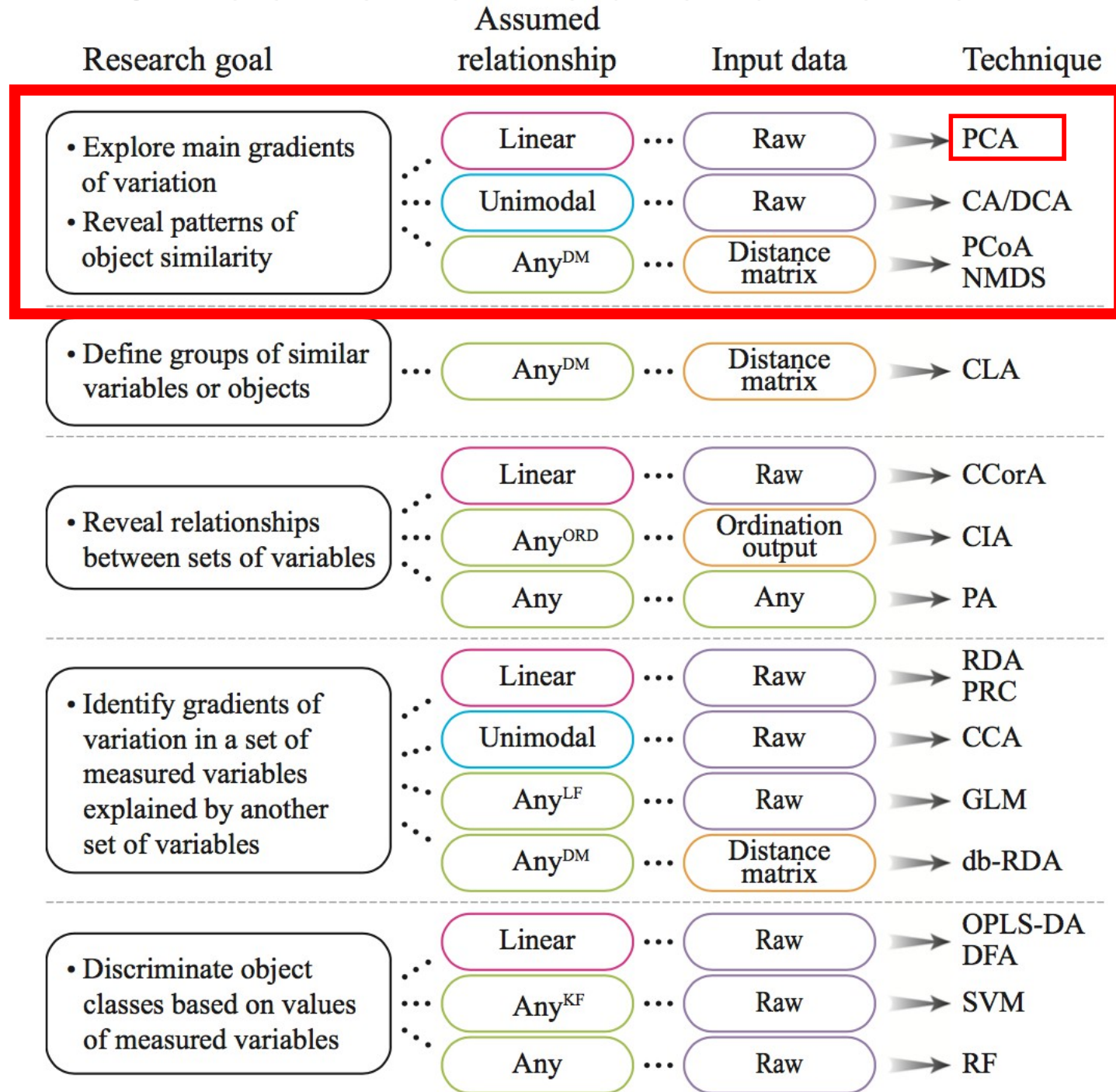
Contents

1. Learning targets and introduction to multivariate analysis
- 2. Overview ordination**
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Ordination: Introduction

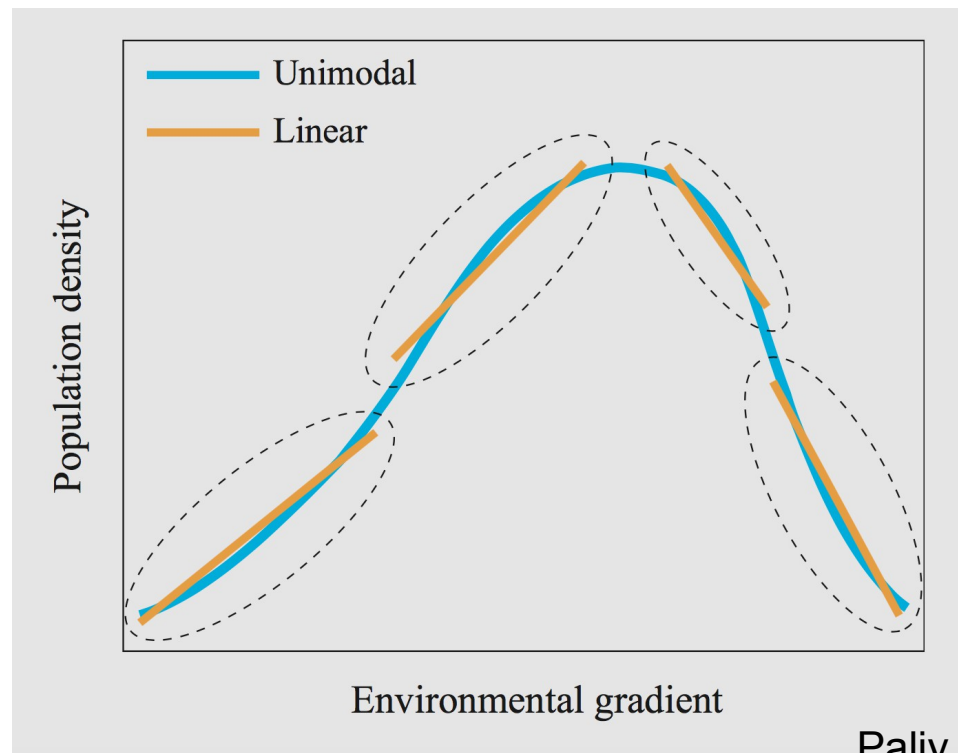
- Extraction of new axes from high dimensional data that sequentially maximise the variance
 - Dimension reduction (e.g. omission of axes that capture low amount of variance)
 - Aggregation of variables into gradients
 - Graphical representation in lower dimension
- Unconstrained ordination: extraction without consideration of variables outside of data set
- Constrained ordination: extraction of axes that are explained by variables of second data set

Unconstrained ordination



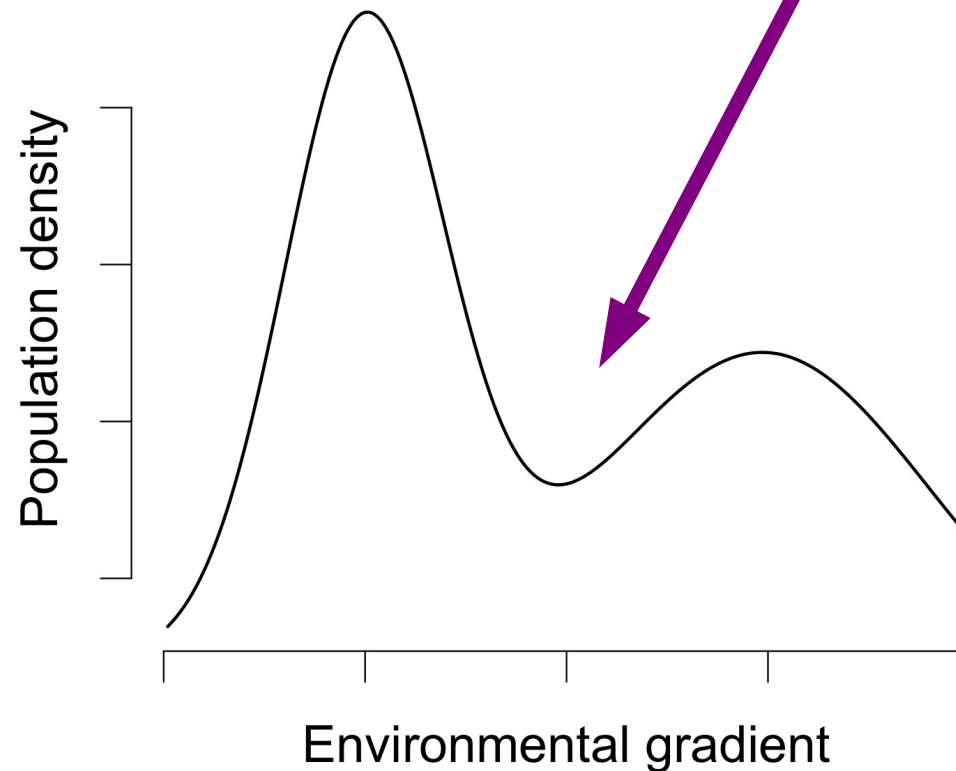
Ordination: Overview

Shape of response	Linear	Unimodal	Any
Unconstrained methods (examples)	PCA	CA	Distance-based: NMDS; GAM-based: UAO (U-VGAM)
Constrained methods (examples)	RDA	CCA	Distance-based: db-RDA; GAM-based: CAO (RR-VGAM)



Ordination: Overview

Shape of response	Linear	Unimodal	Any
Unconstrained methods (examples)	PCA	CA	Distance-based: NMDS; GAM-based: UAO (U-VGAM)
Constrained methods (examples)	RDA	CCA	Distance-based: db-RDA; GAM-based: CAO (RR-VGAM)



Introduction to multivariate analysis, ordination and PCA

Contents

1. Learning targets and introduction to multivariate analysis
2. Overview ordination
- 3. Introduction to PCA and mathematical background**
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Principal Component Analysis

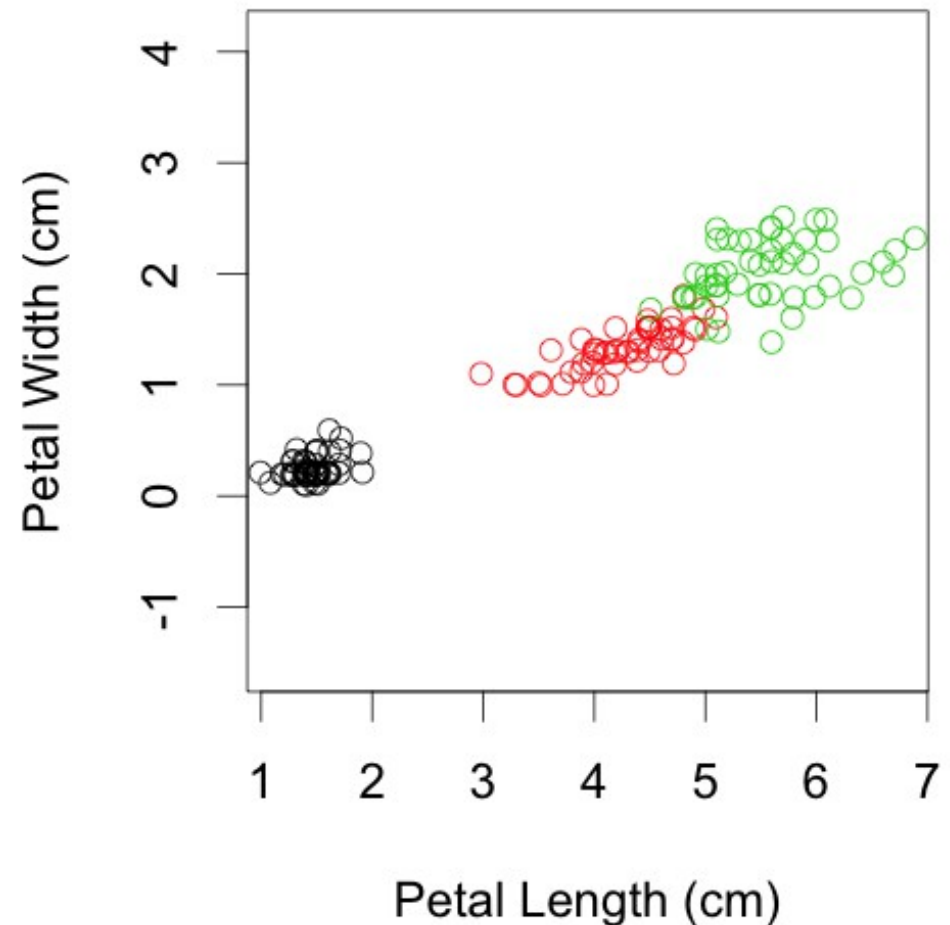
Example-based introduction

Iris data set: sepal length & width, petal length & width for 50 flowers from 3 species of Iris

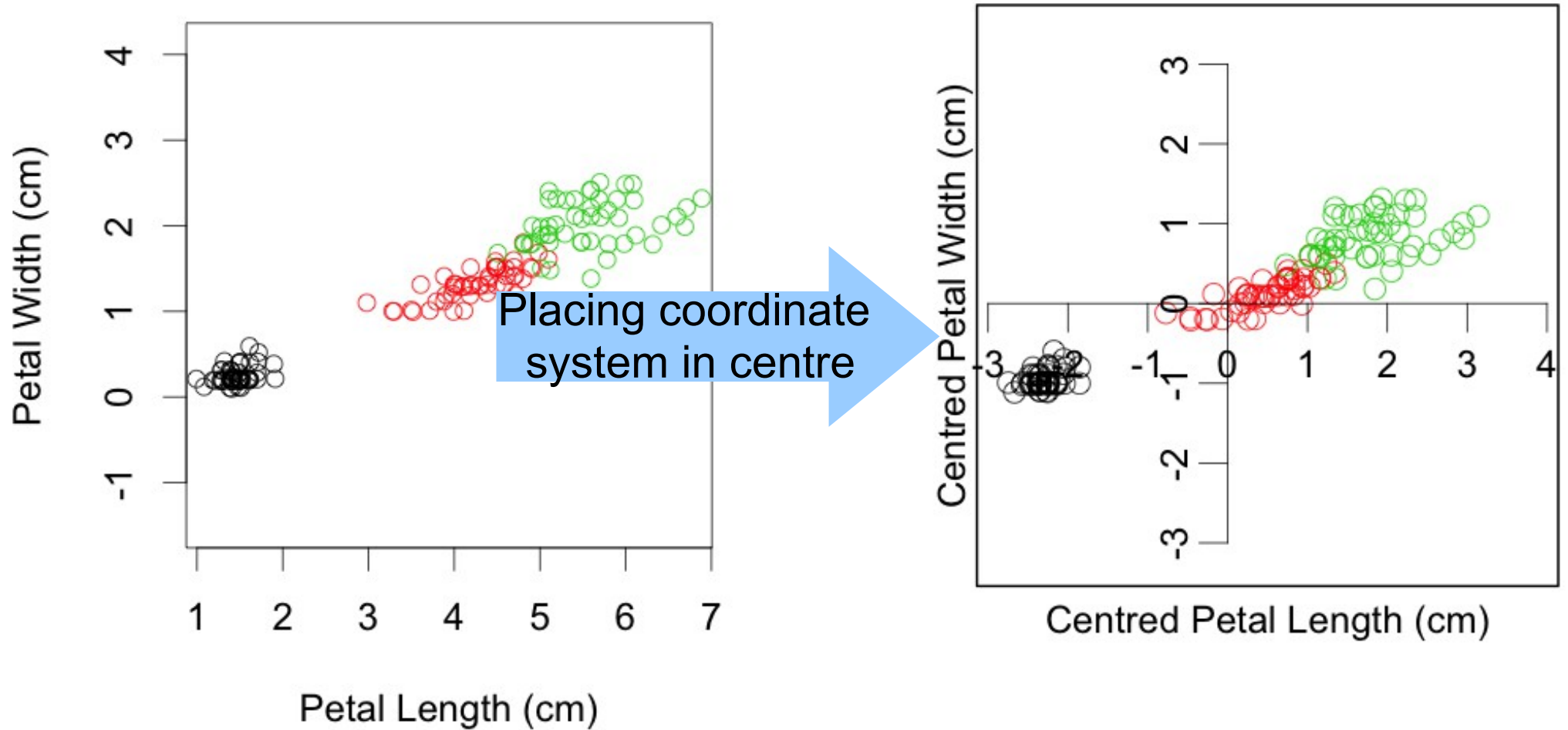


<http://de.wikipedia.org/wiki/Schwertlilien>

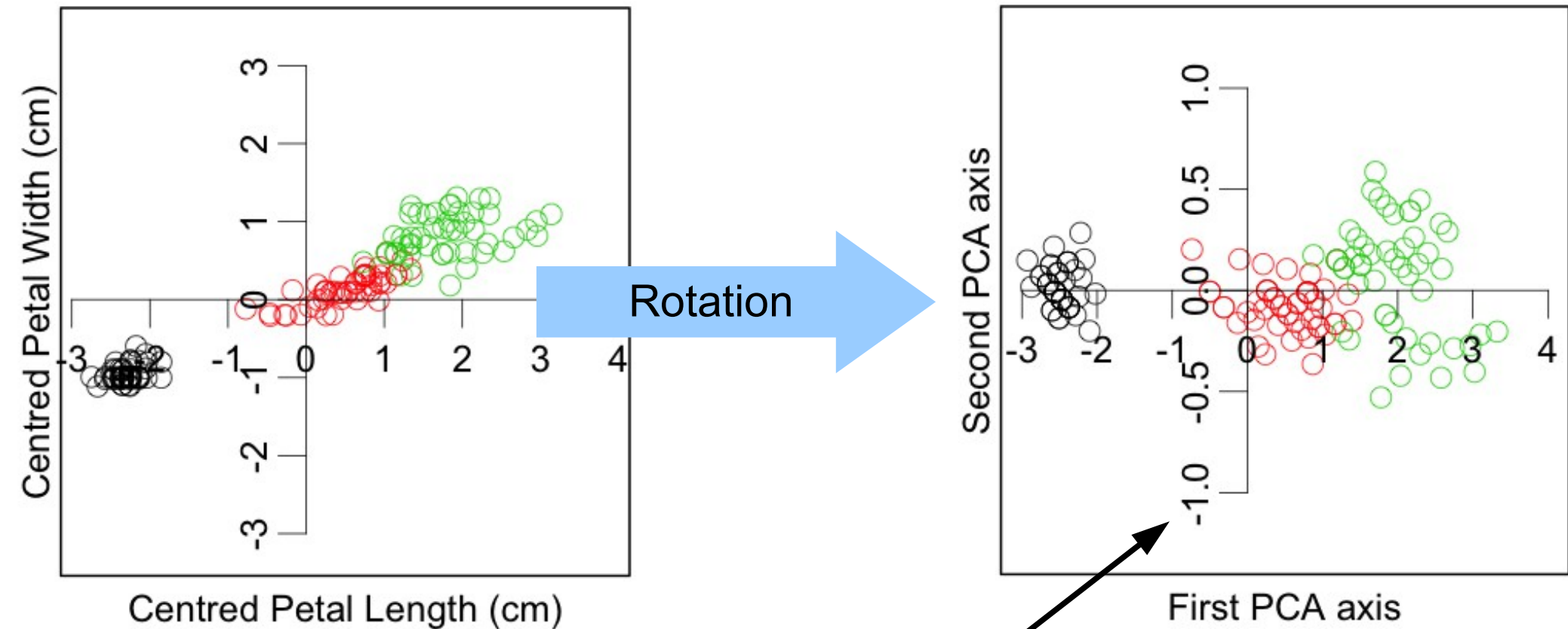
Aim: Represent as much variance as possible on first few axes



Introduction to PCA

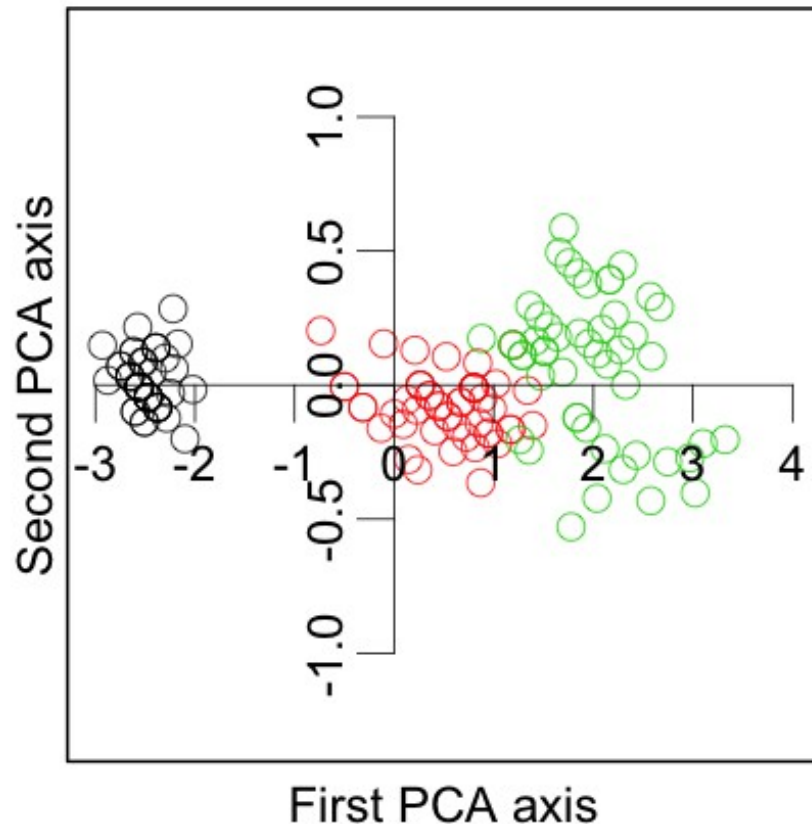


Introduction to PCA



Rescaling of axis based on variation along axis

Results of PCA



Variances

Petal Width 0.58

Petal Length 3.12

Results

Total Variation: 3.70

Importance of components:

	PC1	PC2
Eigenvalue	3.66	0.04
Proportion Explained	0.99	0.01
Cumulative Proportion	0.99	1.00

What is an Eigenvalue?

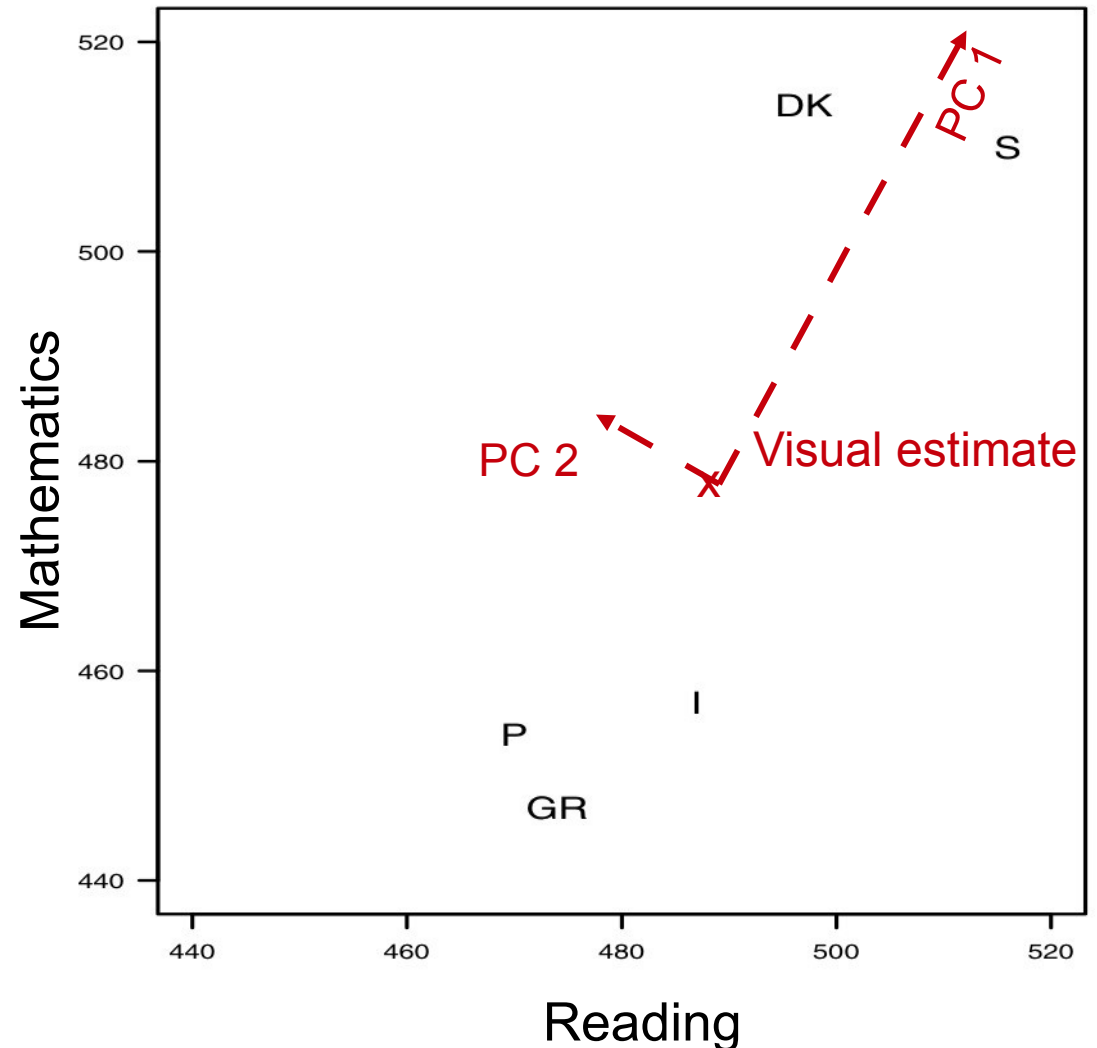
First axis explains 99% of variation!

Mathematical background of PCA

Country	Reading	Mathematics
DK	497	514
GR	474	447
I	487	457
P	470	454
S	516	510

Centering leads to

$$\tilde{X} = \begin{pmatrix} 8.2 & 37.6 \\ -14.8 & -29.4 \\ -1.8 & -19.4 \\ -18.8 & -22.4 \\ 27.2 & 33.6 \end{pmatrix}$$



Search for first axis with maximum variation!

Mathematical background of PCA

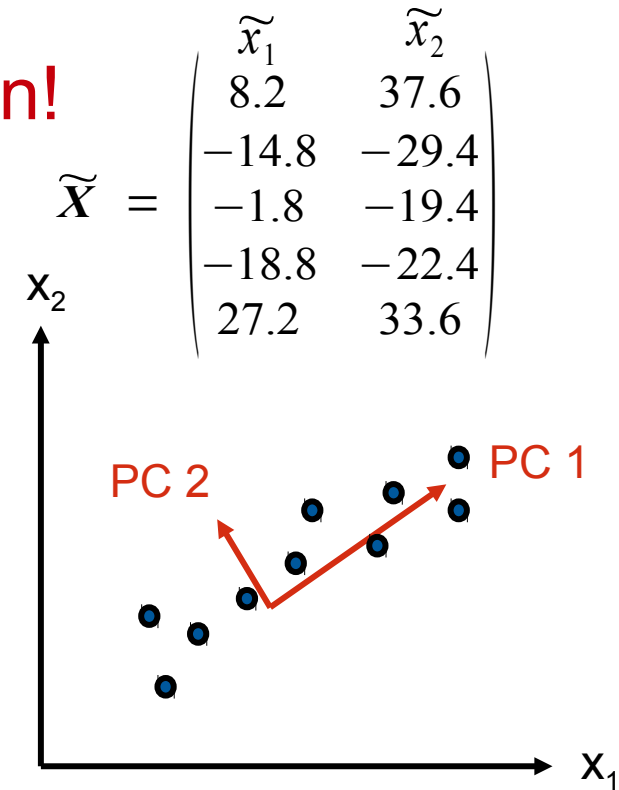
Search for first axis with maximum variation!

Scores on new axis given as

$$PC_1 = a_1 \tilde{x}_1 + a_2 \tilde{x}_2$$

and maximum variation means:

$$\max_{a_1, a_2} \text{Var}(a_1 \tilde{x}_1 + a_2 \tilde{x}_2)$$



Generalise problem: define a_1, a_2 as elements of vector a and likewise \tilde{x}_1, \tilde{x}_2 of matrix \tilde{X} : $\max_a \text{Var}(a \tilde{X})$

Trivial solution: choose high values for a_1, a_2, \dots, a_n

→ introduce condition: $a_1^2 + a_2^2 + \dots + a_n^2 = 1$

Mathematical background of PCA

Solve:

$$\max_a \text{Var}(a\widetilde{X}) \text{ with } a^T a = 1$$

This can be expressed as (see Handl 2010: 79)

$$\max_a (a^T \Sigma a) \text{ with } a^T a = 1$$

Covariance matrix

Using the Lagrange function yields:

$$L(a, \lambda) = a^T \Sigma a - \lambda (a^T a - 1)$$

$$\frac{\partial L(a, \lambda)}{\partial a} = 2 \Sigma a - 2 \lambda a \longrightarrow \text{Eigenvalue problem}$$

$\Sigma a = \lambda a$

$$\frac{\partial L(a, \lambda)}{\partial \lambda} = 1 - a^T a$$

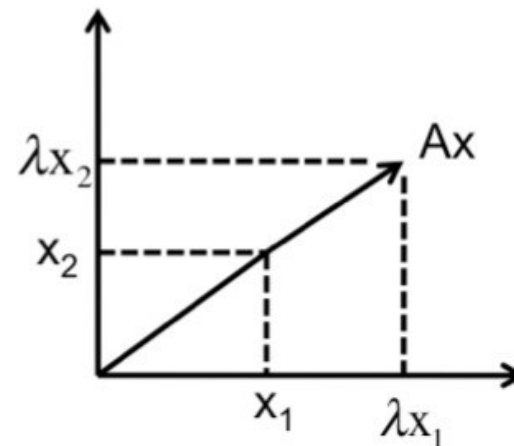
Mathematical basics II: Eigenvalues

Idea: Conversion of a matrix into a matrix with linear independent variables

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \dots & a_{22} & \dots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \xrightarrow{\text{Conversion}} \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_n \end{pmatrix}$$

Eigenvalue problem: $Ax = \lambda x$ Eigenvector
Eigenvalue

Eigenvectors form canonical basis and are only stretched or shrunk by λ when multiplied with A .



Mathematical basics II: Eigenvalues

$$A x = \lambda x \Leftrightarrow A x - \lambda x = 0$$

$$\Leftrightarrow (A - \lambda E) x = 0$$

$$\Leftrightarrow \begin{pmatrix} a_{11} - \lambda_1 & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} - \lambda_n \end{pmatrix} \begin{pmatrix} x_1 \\ \dots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ \dots \\ 0 \end{pmatrix} \quad \begin{array}{l} \text{Homogeneous} \\ \text{linear equation} \\ \text{system (HLS)} \end{array}$$

$$\text{Ignore trivial solution: } x = 0 \quad \Rightarrow \quad A - \lambda E = 0$$

The given HLS has only a non-trivial solution if the columns of $A - \lambda E$ are linearly dependent, which is the case if the determinant = 0.

$$\Rightarrow \det(A - \lambda E) = 0 \Leftrightarrow |A - \lambda E| = 0$$

Example I: Calculation of Eigenvalues

Sample Variance-Covariance matrix from Pisa example:

$$\mathbf{S} = \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix}$$

Following $|\mathbf{A} - \lambda \mathbf{E}| = 0$ we obtain:

$$\left| \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right| = 0 \Leftrightarrow$$

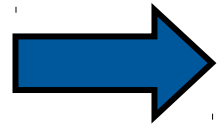
$$\left| \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \right| = 0 \Leftrightarrow$$

$$\begin{vmatrix} 345.7 - \lambda & 528.35 \\ 528.35 & 1071.30 - \lambda \end{vmatrix} = 0 \Leftrightarrow$$

$$(345.7 - \lambda)(1071.30 - \lambda) - 528.35^2 = 0$$

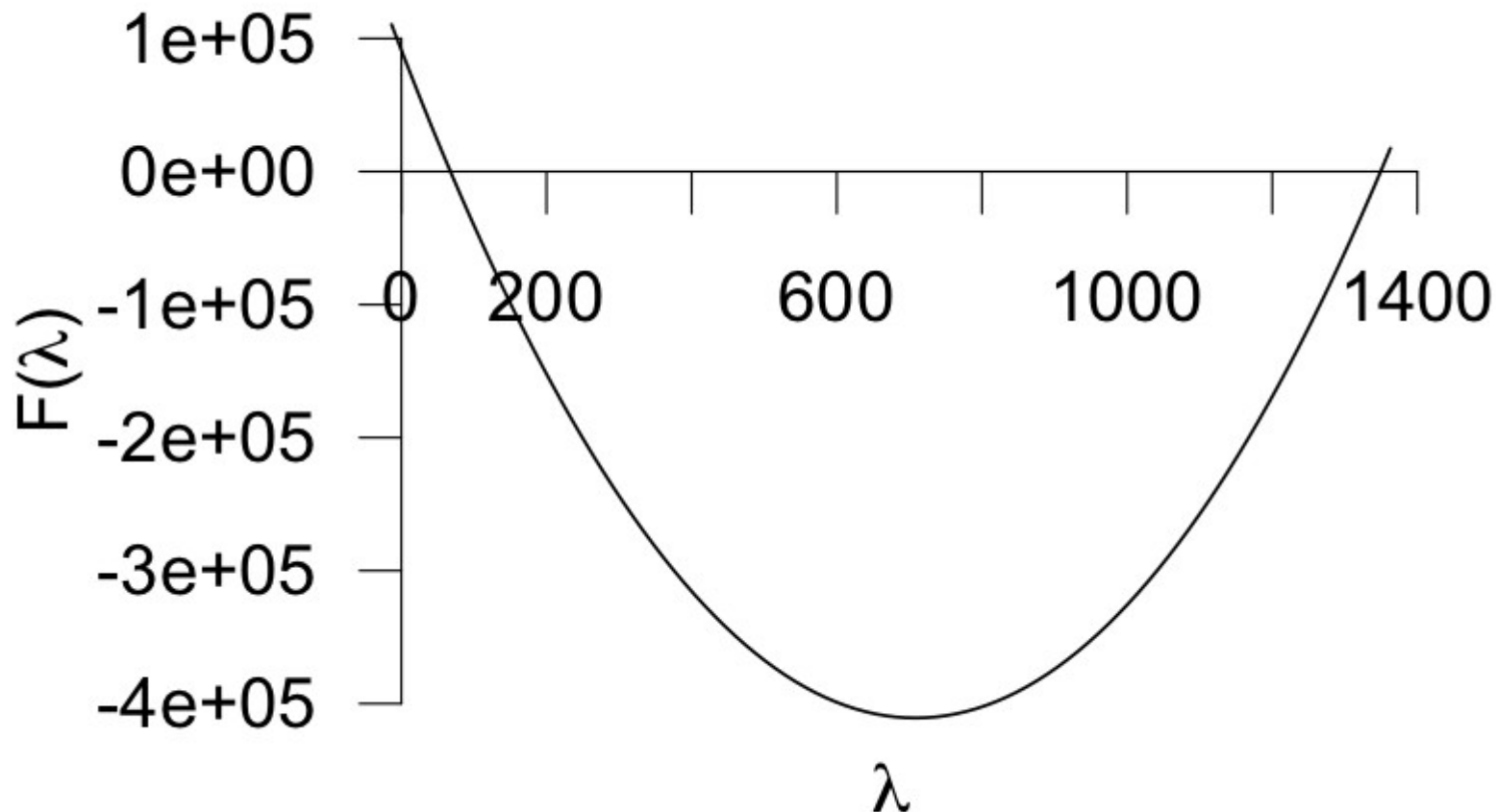
Example I: Calculation of Eigenvalues

$$(345.7 - \lambda)(1071.30 - \lambda) - 528.35^2 = 0$$




Characteristical polynom

$$\lambda^2 - 1417\lambda + 91194.69 = 0$$



Example II: Calculation of Eigenvalues and -vectors

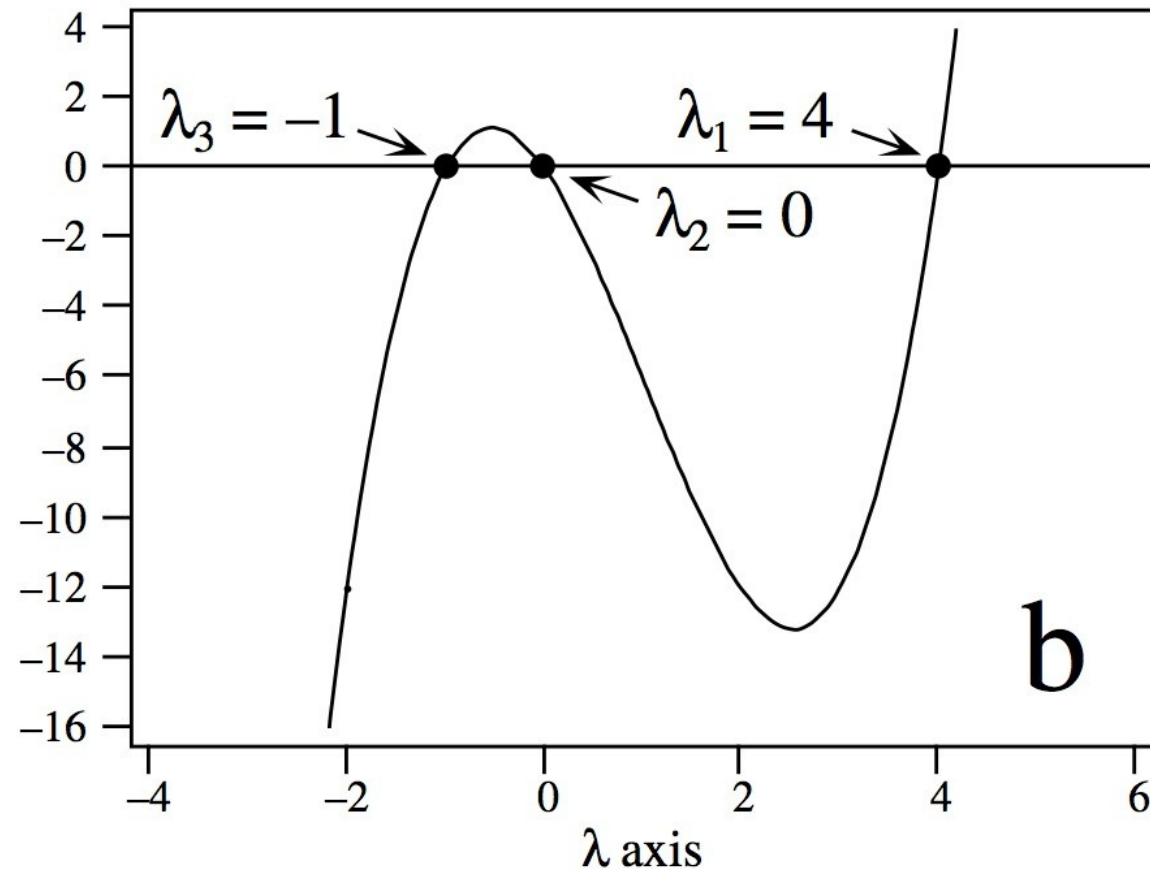
$$\begin{pmatrix} 1 & 3 & -1 \\ 0 & 1 & 2 \\ 1 & 4 & 1 \end{pmatrix}$$


EV calculation
following Sarrus

Characteristical polynom

$$\lambda^3 - 3\lambda^2 - 4\lambda = 0$$

**Eigenvalues λ : 4,
0 and -1**



Example II: Calculation of Eigenvalues and -vectors

Calculation of eigenvector for $\lambda = 4$

$$(A - \lambda E)x = 0$$

$$\left(\begin{pmatrix} 1 & 3 & -1 \\ 0 & 1 & 2 \\ 1 & 4 & 1 \end{pmatrix} - \lambda_1 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0 \Leftrightarrow \begin{pmatrix} 1-4 & 3 & -1 \\ 0 & 1-4 & 2 \\ 1 & 4 & 1-4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$\begin{array}{cccc} -3x_1 & 3x_2 & -x_3 & 0 \\ 0 & -3x_2 & 2x_3 & = 0 \\ x_1 & 4x_2 & -3x_3 & 0 \end{array} \Leftrightarrow \begin{array}{cccc} -3x_1 & 0 & x_3 & 0 \\ 0 & 1.5x_2 & 0 & = x_3 \\ x_1 & 4x_2 & -3x_3 & 0 \end{array}$$

Matrix is singular (no unique solution)

→ fix value of one variable e.g. $x_1 = 1$.

Example II: Calculation of Eigenvalues and -vectors

$$x_1=1 \Rightarrow \begin{pmatrix} -3 & 0 & 0 \\ 0 & 1.5x_2 & 0 \\ 1 & 4x_2 & -3x_3 \end{pmatrix} = \begin{pmatrix} -x_3 \\ x_3 \\ 0 \end{pmatrix} \Rightarrow x_1=1; x_2=2; x_3=3$$

Calculation of eigenvectors for all eigenvalues yields the following matrix of eigenvectors (or multiples of columns):

$$\begin{pmatrix} 1 & 7 & 2 \\ 2 & -2 & -1 \\ 3 & 1 & 1 \end{pmatrix}$$

Eigenvalues: 4; 0 and -1

Introduction to multivariate analysis, ordination and PCA

Contents

1. Learning targets and introduction to multivariate analysis
2. Overview ordination
3. Introduction to PCA and mathematical background
- 4. PCA results and interpretation**
5. PCA diagnosis and extensions, brief tutorial

Results of PCA II

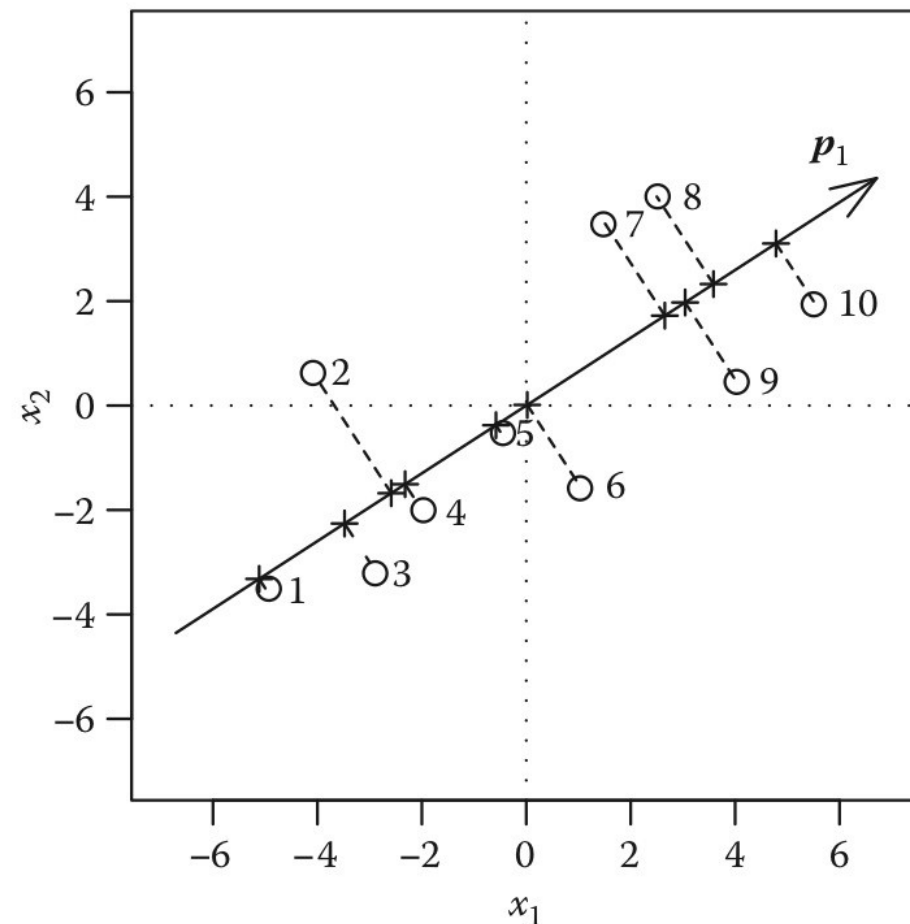
PC Scores

Demo Example for PCA with 10 Objects and Two Mean-Centered Variables x_1 and x_2

i	x_1	x_2	t_1	t_2
1	-5.0	-3.5	-6.10	-0.21
2	-4.0	0.5	-3.08	2.60
3	-3.0	-3.0	-4.15	-0.88
4	-2.0	-2.0	-2.77	-0.59
5	-0.5	-0.5	-0.69	-0.15
6	1.0	-1.5	0.02	-1.80
7	1.5	3.5	3.16	2.12
8	2.5	4.0	4.27	1.99
9	4.0	0.5	3.63	-1.76
10	5.5	2.0	5.70	-1.32
\bar{x}	0.00	0.00	0.00	0.00
v	12.22	6.72	16.22	2.72
$v\%$	64.52	35.48	85.64	14.36

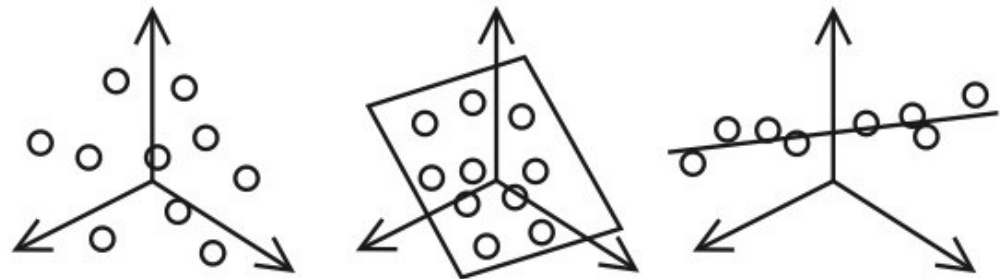
Note: i , Object number; t_1 and t_2 are the PCA scores of PC1 and PC2, respectively;
 \bar{x} , mean; v , variance; $v\%$, variance in percent of total variance.

PC scores result from multiplication of scores from initial axes with eigenvectors



Number of principal components

- Number of descriptors/explanatory variables determines number of eigenvalues and thus principal components
- Largest eigenvalue (and the corresponding eigenvector) explains highest share of total variance
- Aim is to represent the major variation with a few principal components → How many components are needed?



Number of variables

3

3

3

Number of relevant components
= intrinsic dimensionality

3

2

1

How many principal components needed?

Some criteria to evaluate the optimal number of axes:

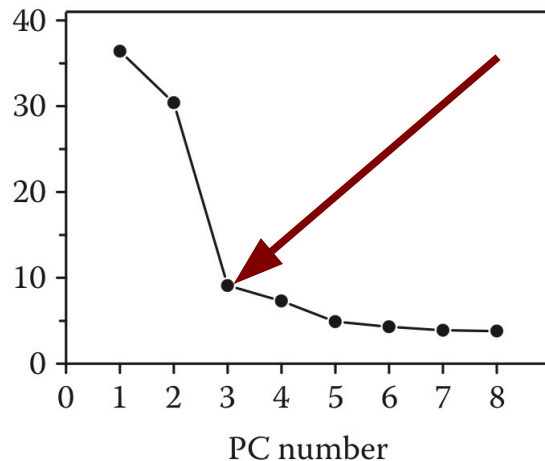
1. Sum criterion

$$\frac{\sum_{i=1}^r \lambda_i}{\sum_{j=1}^p \lambda_j} \geq \alpha$$

2. Broken-Stick criterion

$$\lambda_i > \frac{1}{p} \sum_{i=1}^p \frac{1}{i}$$

3. Scree plot



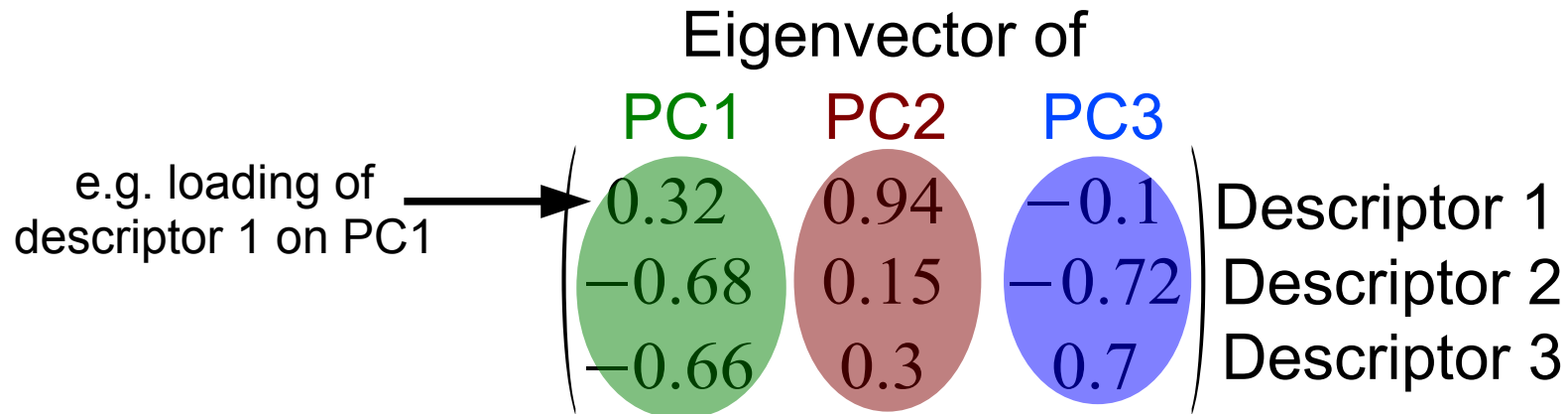
4. Cross-validation

$$\min_S \text{MSEP}(S) = \frac{1}{I K} \sum_{i=1}^I \sum_{k=1}^K \left(x_{i,k} - (\hat{x}_{i,k})^{(S)} \right)^2$$

For the matrix $X_{I \times K}$, we search the number of PC S minimizing the mean square error of prediction (MSEP).

Importance of descriptor for PC axes

- Elements of eigenvector matrix = 'loadings', indicate weight of original descriptor on PCs



- Easier to interpret: Correlation loadings $r_i = a_i \sqrt{\lambda_i}$
- Interpretation of descriptor importance complicated if many variables load on PC
- Sparse PCA – introduces penalty term (cf. LASSO):

$$\max_a \text{Var}(a\tilde{\mathbf{X}}) - \lambda \|a\| \quad \text{with} \quad a^T a = 1$$

Introduction to multivariate analysis, ordination and PCA

Contents

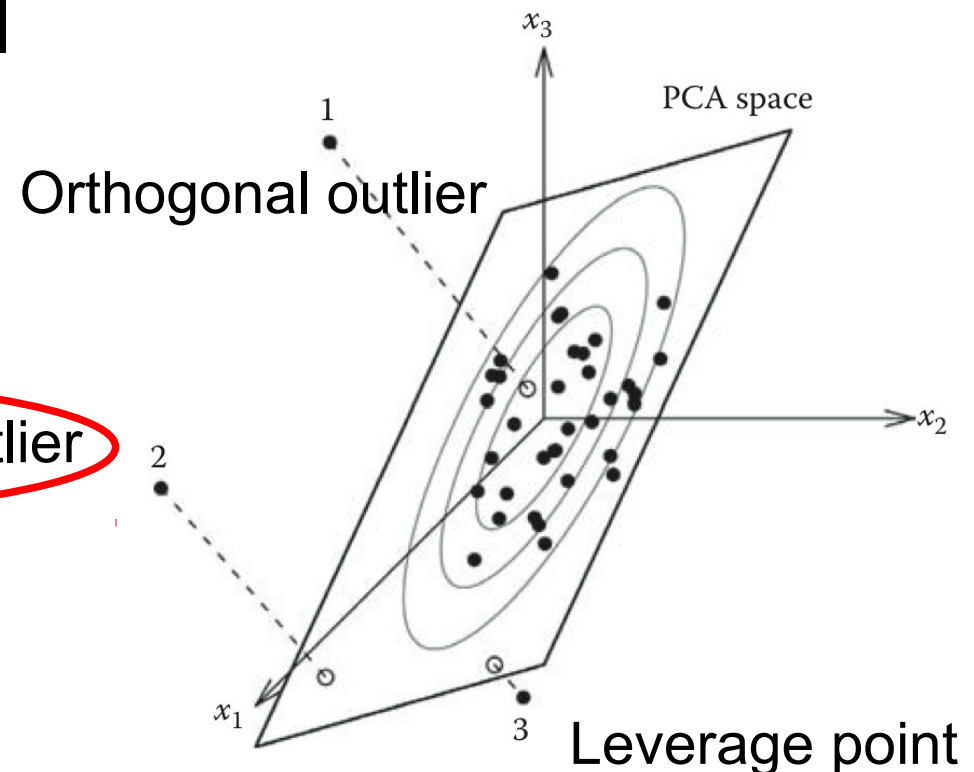
1. Learning targets and introduction to multivariate analysis
2. Overview ordination
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
- 5. PCA diagnosis and extensions, brief tutorial**

PCA assumptions and diagnosis

- Independence of observations (temporal and spatial independence)
- multivariate normality (depending on research question)
- no serious outliers, otherwise use robust PCA (see Varmuza & Filzmoser 2009: chapter 3.5)
- Computation of orthogonal distance (OD) and score distance (SD)

Leverage point and orthogonal outlier

Problematic



PCA assumptions and limitations

- Linear gradient of descriptors (rarely the case for species data, but often for environmental data)
- Euclidean distances inappropriate for species data with many double zeros (joint absences)
- Adding noise variables to data increases fraction of variance on first axis, but has no meaning
- Useful technique in exploratory analysis and for analysing environmental data (or other data with linear gradients)
- Best results for large n and high $n:p$ (p = descriptors)

Brief tutorial for PCA

1. Check if conditions for descriptors are met (quantitative, multivariate normality, linear)
2. Conduct PCA (or sparse PCA) on scaled descriptors unless they exhibit a similar variation and have been measured on similar scale
3. Check for outliers
4. Determine number of principal components
5. How informative are first two PCs?
6. Which descriptors contribute most to PCs?
7. Visualise and interpret

Principal component regression

- Extract (unscaled) PC scores to use PCs as descriptors in multiple regression analysis
- PCs are orthogonal → fix for multicollinearity in regression analysis
- In low $n:p$ situations, the few last PCs are often removed to reduce number of predictors in regression
→ Can be problematic because low variance of PC does not imply low explanatory power for response variable
→ not necessarily a fix for low $n:p$ ratios

Introduction to multivariate analysis, ordination and PCA

Ralf B. Schäfer

These slides and notes complement the lecture with exercises “Tools for complex data analysis” for ecotoxicologists and environmental scientists. Do not hesitate to contact me if you have any comments or you find any errors (slides, slide notes, or code): schaefer-ralf@uni-landau.de

While I made notes below the slides, some aspects are only mentioned in the R demonstration associated with the lecture.

Introduction to multivariate analysis, ordination and PCA

Contents

- 1. Learning targets and introduction to multivariate analysis**
2. Overview ordination
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Learning targets

- Explain the specifics of multivariate analysis
- List ordination methods and select one based on research goal and question
- Explain the mathematical basis of and apply PCA
- Interpreting results from a PCA

Learning targets and study questions

- Explain the specifics of multivariate analysis
 - Why should you favour multivariate approaches for multivariate data?
 - Outline the differences to the univariate case when diagnosing multivariate outliers and normality.
- List ordination methods and select one based on research goal and question
 - Outline the aim of ordination
 - Distinguish constrained and unconstrained ordination
 - Which criteria should guide the selection of an ordination method?

Learning targets and study questions

- Explain the mathematical basis of and apply PCA
 - Explain the variance-covariance matrix.
 - What are eigenvalues and eigenvectors?
 - How do eigenvalues relate to the variance of a PC?
 - Outline criteria to determine the optimal number of PCs.
 - What is sparse PCA and how does it influence the evaluation of descriptor contribution to PCs?
- Interpreting results from a PCA
 - Explain biplots with respect to (a) correlation between variables and (b) relation between sites/species and variables.
 - How does scaling influence interpretation of a biplot?
 - Which objects from a PCA would be extracted as non-collinear predictors for a multiple regression analysis?

From univariate to multivariate statistics

	univariate	multivariate
Variables (vars.)	Single/multiple predictors, single response variable Y	Single/multiple predictors, multiple response vars. Y_1, \dots, Y_n
Distribution of response	One-dimensional	n -dimensional
Data format	Y is vector	Y_1, \dots, Y_n constitute matrix
Example	Species richness explained by environmental variables	Community explained by environmental variables

6

The table describes the difference between univariate and multivariate statistics with a focus on the response (i.e. dependent) variable(s). Sometimes linear (GLMs) with multiple predictors are also considered as multivariate (e.g. Lloyd 2010: Spatial data analysis. Oxford Univ. Press: Oxford), though they rather represent multivariable models (see Hidalgo & Goodman 2013).

Cited reference:
Hidalgo B. & Goodman M. (2013) Multivariate or Multivariable Regression? *American Journal of Public Health* 103, 39–40. Freely available under:
<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3518362/>

Multivariate data analysis: Introduction

Some advantages of multivariate over univariate methods for analysing multivariate data

- Not all research questions can be answered with univariate statistical methods
e.g. *What are the most important environmental variables determining community composition?*
- Multivariate methods allow for dimension reduction and visualisation of multidimensional data
e.g. *Ordination, Cluster dendrogram*
- Joint (multivariate) analysis can reduce noise and increase power when assessing statistical hypotheses

7

In univariate analyses, only the relationship between single taxa and environmental variables can be examined, whereas multivariate analyses allow for the analysis of how environmental variables act on a community of organisms.

Furthermore, response variables typically incorporate some random variation and therefore can display an association with noise variables for example in multiple regression analysis (Flack & Chang 1987). Moreover, the larger the number of noise variables, the higher the probability of high correlation with meaningful predictors of a response, which also leads to associations between the response and noise variables. The simultaneous consideration of several response variables as in multivariate analysis reduces the influence of noise variables and can increase the power when assessing hypotheses. For example, multivariate methods are used in climatology for the reconstruction of the global temperature trend, where single time lines would be insufficient to discover and establish the global warming trend (Ammann & Wahl 2007).

Ammann, E. R. & Wahl C. M. 2007: The importance of the geophysical context in statistical evaluations of climate reconstruction procedures. *Climate Change* 85 (1-2): 71-88
Flack, V. F. & Chang, P. C., Frequency of Selecting Noise Variables in Subset Regression-Analysis - a Simulation Study. *American Statistician* 1987: 84-86

Multivariate approaches in R

Available methods and developments: CRAN Task View

CRAN Task View: Multivariate Statistics

Maintainer: Paul Hewson
Contact: Paul.Hewson at plymouth.ac.uk
Version: 2018-07-21
URL: <https://CRAN.R-project.org/view=Multivariate>

Base R contains most of the functionality for classical multivariate analysis, somewhere. There are a large number of packages on CRAN which extend this methodology, a brief overview is given below. Application-specific uses of multivariate statistics are described in relevant task views, for example whilst principal components are listed here, ordination is covered in the [Environmetrics](#) task view. Further information on supervised classification can be found in the [MachineLearning](#) task view, and unsupervised classification in the [Cluster](#) task view.

The packages in this view can be roughly structured into the following topics. If you think that some package is missing from the list, please let me know.

Visualising multivariate data

- Graphical Procedures:* A range of base graphics (e.g. `pairs()` and `coplot()`) and [lattice](#) functions (e.g. `xyplot()` and `spplot()`) are useful for visualising pairwise arrays of 2-dimensional scatterplots, clouds and 3-dimensional densities. `scatterplot.matrix` in the [car](#) provides usefully enhanced pairwise scatterplots. Beyond this, [scatterplot3d](#) provides 3 dimensional scatterplots, [aplpack](#) provides bagplots and `spin3R()`, a function for rotating 3d clouds. [misc3d](#), dependent upon [rgl](#), provides animated functions within R useful for visualising densities. [YaleToolkit](#) provides a range of useful visualisation techniques for multivariate data. More specialised multivariate plots include the following: `faces()` in [aplpack](#) provides Chernoff's faces; `parcoord()` from [MASS](#) provides parallel coordinate plots; `stars()` in [graphics](#) provides a choice of star, radar and cobweb plots respectively. `mstree()` in [ade4](#) and `spantree()` in [vegan](#) provide minimum spanning tree functionality. [calibrate](#) supports biplot and scatterplot axis labelling. [geometry](#), which provides an interface to the `qhull` library, gives indices to the relevant points via `convexhulln()`. [ellipse](#) draws ellipses for two parameters, and provides `plotcorr()`, visual display of a correlation matrix. [denpro](#) provides level set trees for multivariate visualisation. Mosaic plots are available via `mosaicplot()` in [graphics](#) and `mosaic()` in [vcd](#) that also contains other visualization techniques for multivariate categorical data. [gclus](#) provides a number of cluster specific graphical enhancements for scatterplots and parallel coordinate plots See the links for a reference to GGobi. [rggobi](#) interfaces with GGobi. [xgobi](#) interfaces to the XGobi and XGvis programs which allow linked, dynamic multivariate plots as well as projection pursuit. Finally, [iplots](#) allows particularly powerful dynamic interactive graphics, of which interactive parallel co-ordinate plots and mosaic plots may be of great interest. Seriation methods are provided by [seriation](#) which can reorder matrices and dendrograms.
- Data Preprocessing:* `summarize()` and `summary.formula()` in [Hmisc](#) assist with descriptive functions; from the same package `varclus()` offers variable clustering while `dataRep()` and `find.matches()` assist in exploring a given dataset in terms of representativeness and finding matches. Whilst `dist()` in base and `daisy()` in [cluster](#) provide a wide range of distance measures, [proxy](#) provides a framework for more distance measures, including measures between matrices. [simba](#) provides functions for dealing with presence / absence data including similarity matrices and reshaping.

Hypothesis testing

- [ICSNP](#) provides Hotellings T2 test as well as a range of non-parametric tests including location tests based on marginal ranks, spatial median and spatial signs computation, estimates of shape. Non-parametric two sample tests are also available from [cramer](#) and spatial sign and rank tests to investigate location, sphericity and independence are available in [SpatialNP](#).

Multivariate distributions

- Descriptive measures:* `cov()` and `cor()` in stats will provide estimates of the covariance and correlation matrices respectively. [ICSNP](#) offers several descriptive measures such as `spatial.median()` which provides an estimate of the spatial median and further functions which provide estimates of scatter. Further robust methods are provided such as `cov.rob()` in [MASS](#) which provides robust estimates of the variance-covariance matrix by minimum volume ellipsoid, minimum covariance determinant or classical product-moment. [covRobust](#) provides robust covariance estimation via nearest neighbor variance estimation. [robustbase](#) provides robust covariance estimation via fast minimum covariance determinant with `covMCD()` and the Orthogonalized pairwise estimate of Gnanadesikan-Kettenring via `covOGK()`. Scalable robust methods are provided within [rccov](#) also using fast minimum covariance determinant with `covMed()` as well as M-estimators with `covMest()`. [corpcor](#) provides shrinkage estimation of large scale covariance and (partial) correlation matrices.
- Densities (estimation and simulation):* `mvnorm()` in [MASS](#) simulates from the multivariate normal distribution. [mvtnorm](#) also provides simulation as well as probability and quantile functions for both the multivariate t distribution and multivariate normal distributions as well as density functions for the multivariate normal distribution. [mnormt](#) provides multivariate normal and multivariate t density and distribution functions as well as random number simulation. [sn](#) provides density, distribution and random number generation for the multivariate skew normal and skew t distribution. [delt](#) provides a range of functions for estimating multivariate densities by CART and greedy methods. Comprehensive information on mixtures is given in the [Cluster](#) view, some density estimates and random numbers are provided by `rmvnorm.mixt()` and `dmvnorm.mixt()` in [ks](#), mixture fitting is also provided within [bayesm](#). Functions to simulate from the Wishart distribution are provided in a number of places, such as `rwishart()` in [bayesm](#) and `rwish()` in [MCMCpack](#) (the latter also has a density function `dwish()`). `bkde2D()` from [KernSmooth](#) and `kde2d()` from [MASS](#) provide binned and non-binned 2-dimensional kernel density estimation, [ks](#) also provides multivariate kernel smoothing as does [ash](#) and [GenKern](#). [prim](#) provides patient rule induction methods to attempt to find regions of high density in high dimensional multivariate data. [feature](#) also provides methods for determining feature significance in multivariate data (such as in relation to local modes).
- Assessing normality:* [mvnormtest](#) provides a multivariate extension to the Shapiro-Wilks test, [mvoutlier](#) provides multivariate outlier detection based on robust methods. [ICS](#) provides tests for multi-normality. `mvnorm.etest()` in [energy](#) provides an assessment of normality based on E statistics (energy); in the same package `k.sample()` assesses a number of samples for equal distributions. Tests for Wishart-distributed covariance matrices are given by `mauchly.test()` in stats.
- Copulas:* [copula](#) provides routines for a range of (elliptical and archimedean) copulas including normal, t, Clayton, Frank, Gumbel, [fgac](#) provides generalised archimedean copula.

Additional methods that are still under development can be found on R Forge:
https://r-forge.r-project.org/softwaremap/tag_cloud.php?tag=multivariate

Multivariate outlier checking

Multivariate and univariate exploration/diagnosis similar for e.g. multicollinearity and transformation, but approaches differ for checking of e.g. outliers and distributional assumptions

Multivariate outliers

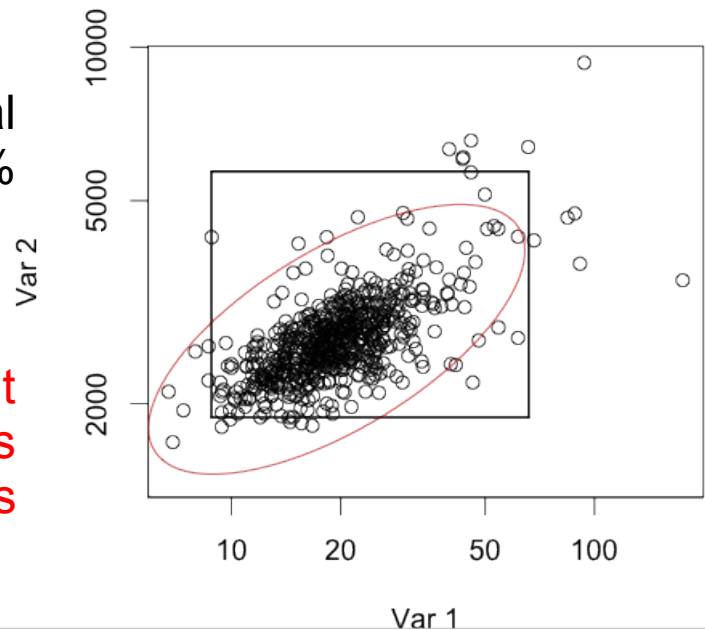
Use joint multivariate distribution of variables to find outliers (i.e. extreme points from centre of multivariate sample)

Outside of box:

Outliers if considering individual distributions of variables (99% quantile of each variable)

Outside of ellipse:

Outliers if considering joint distribution of both variables (99% quantile using Mahalanobis distance)

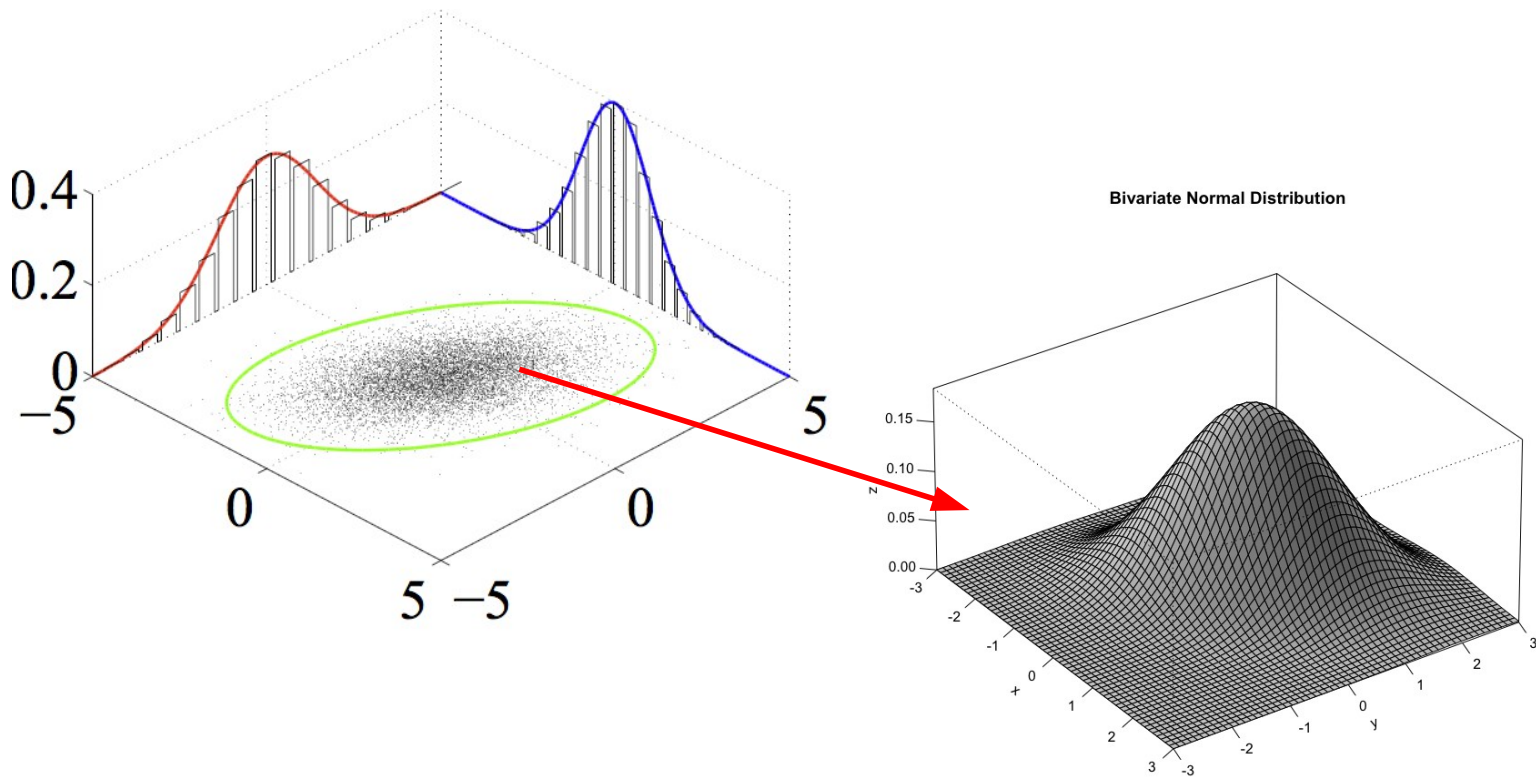


9

In univariate analysis, we focused on model outliers related to a single response, whereas in multivariate analysis we need to consider multiple responses simultaneously. Nevertheless, exploration in the context of univariate analyses such as multiple regression can also consider the joint distribution to identify predictor outliers. Hence, the boundaries are rather fluid between univariate and multivariate settings, and the general approach is often very similar.

Multivariate normal distribution

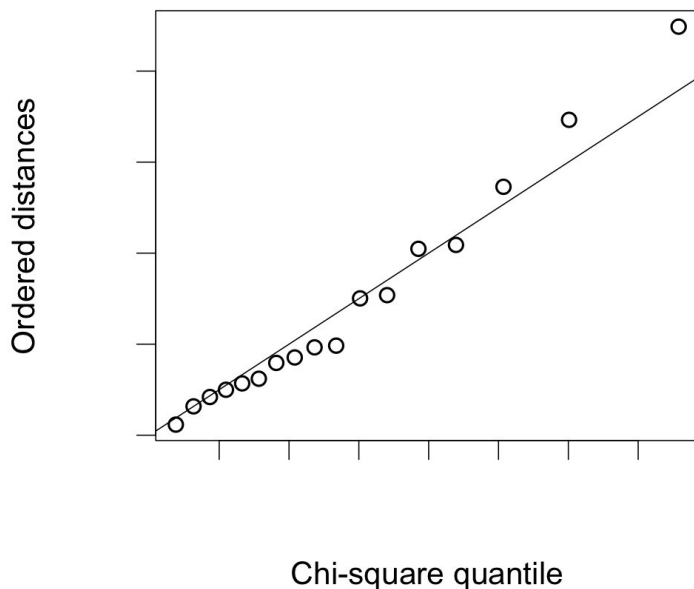
Multivariate normal distribution is evaluated instead of univariate normal distribution



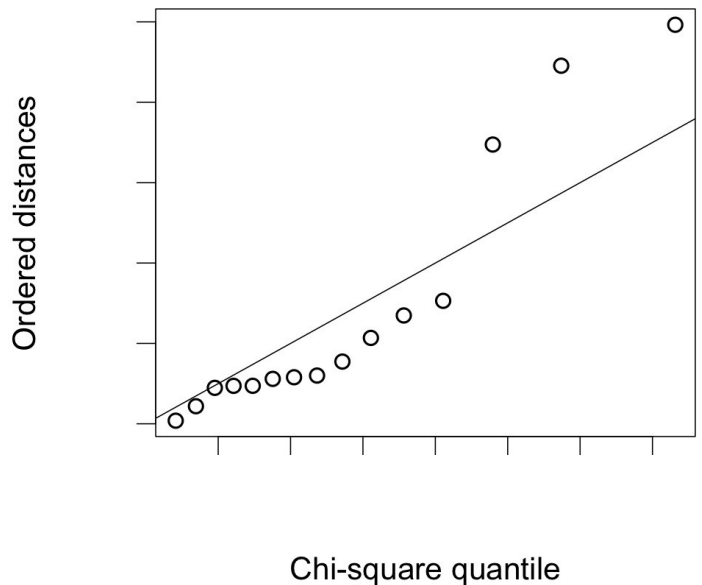
Multivariate normal distribution

Visual check of multivariate normality with QQ plots for sample Mahalanobis distances (to centroid) and theoretical quantiles from the χ^2 distribution

QQ plot with data sampled from a multivariate normal distribution



QQ plot with data sampled from an exponential distribution



11

Regarding the calculation of the QQ plot: The distance to the multivariate centroid of the sample is calculated for each empirical multivariate observation x_i and weighed by the inverse of the sample covariance matrix (Mahalanobis distance (D_M); further details on the D_M are provided in a few slides). The D_M s are then ordered and compared to the quantiles of a beta or χ^2 distribution. The χ^2 distribution can be misleading for a low ratio of observations to variables (<25), see Small (1978) for details. In this case the beta distribution may yield more reliable results.

Several hypothesis tests are available to check for multivariate normality (see CRAN Task View). The criticism on hypothesis testing for normality outlined for the univariate case largely applies to these.

Small, N. J. H. 1978: Plotting squared radii. *Biometrika* 65 (3): 657-658

Covariance matrix

For a data matrix \mathbf{Y} containing the variables y_1, \dots, y_p

$$\mathbf{Y} = \begin{pmatrix} y'_1 \\ y'_2 \\ \vdots \\ y'_i \\ \vdots \\ y'_n \end{pmatrix} = (\text{units}) \begin{matrix} & \begin{matrix} \text{(variables)} \\ 1 & 2 & \dots & j & \dots & p \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} & \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1j} & \dots & y_{1p} \\ y_{21} & y_{22} & \dots & y_{2j} & \dots & y_{2p} \\ \vdots & \vdots & & \vdots & & \vdots \\ y_{i1} & y_{i2} & \dots & y_{ij} & \dots & y_{ip} \\ \vdots & \vdots & & \vdots & & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nj} & \dots & y_{np} \end{pmatrix} \end{matrix}$$

the sample covariance matrix for these variables is \mathbf{S} (Sigma)

$$\mathbf{S} = (s_{jk}) = \begin{pmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ s_{21} & s_{22} & \dots & s_{2p} \\ \vdots & \vdots & & \vdots \\ s_{p1} & s_{p2} & \dots & s_{pp} \end{pmatrix}$$

$$\text{where } s_{jk} = \frac{1}{n-1} \sum_{i=1}^n (y_{i,j} - \bar{y}_j)(y_{i,k} - \bar{y}_k)$$

12

taken from Rencher (2012): *Methods of multivariate Analysis*. 55f

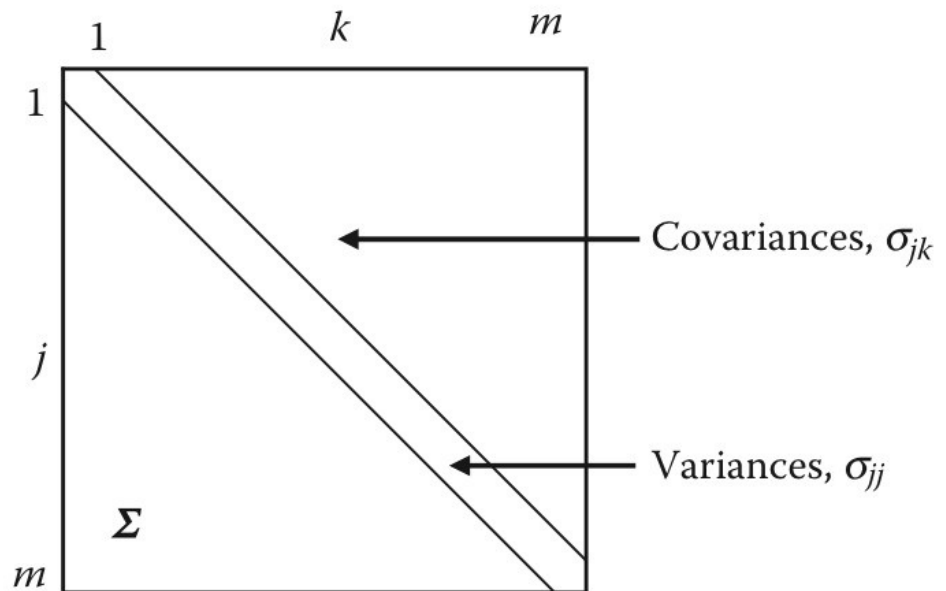
The formula for the covariance of two variables y_j and y_k (i.e. $s_{j,k}$) should be familiar to you from univariate statistics: the covariance is calculated in the context of bivariate correlation (the bivariate correlation for two variables x, y is the covariance divided by the product of the two standard deviations of the variables):

$$r_{x,y} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}}$$

Covariance matrix

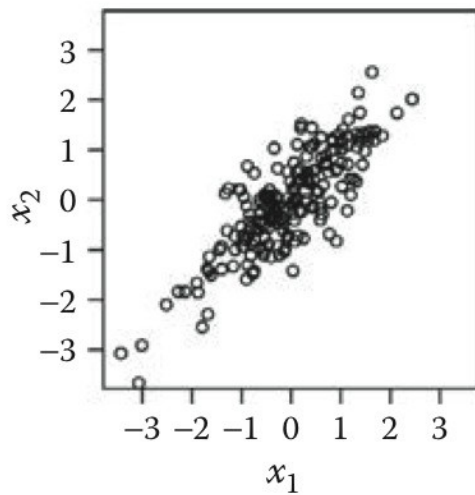
In the diagonal, the equation simplifies to:

$$s_{jj} = \frac{1}{n-1} \sum_{i=1}^n (y_{i,j} - \bar{y}_j)^2$$

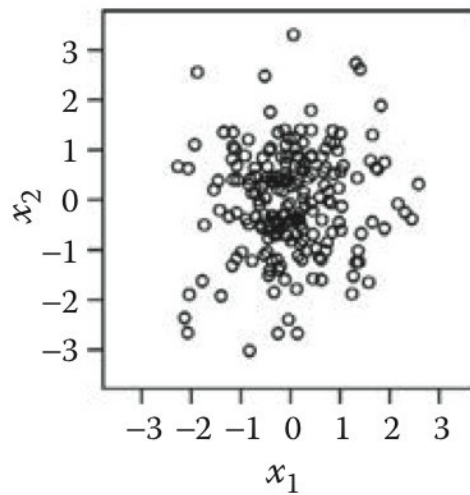


Varmuza K. & Filzmoser P. (2009) Introduction to multivariate statistical analysis in chemometrics. CRC Press/Taylor & Francis, Boca Raton, Fla.

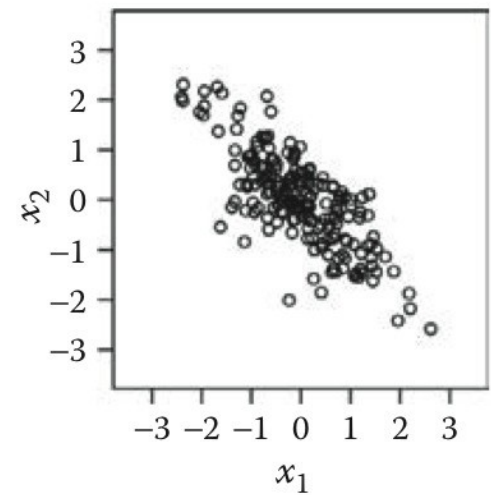
Covariance matrix for two variables



$$\Sigma_1 = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$



$$\Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



$$\Sigma_3 = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}$$

Multivariate distance from mean to observation is measured with the Mahalanobis distance D_M :

$$D_M(x) = \sqrt{(x - \mu)^T \mathbf{S}^{-1} (x - \mu)}$$

→ D_M is distance of vector x from the mean vector μ weighed by their covariance (given in sample covariance matrix \mathbf{S})

14

Varmuza & Filzmoser 2009

\mathbf{S} is the sample covariance matrix, which represents an estimate of the true covariance matrix Σ .

D_M can also be calculated for other vectors than μ , for example to quantify the multivariate distance between any two vectors of the same length. It represents the multivariate extension of the z -score, which is calculated by subtracting the mean of x (given as μ) from each observation i of x , and dividing the result by the standard deviation of x :

$$z_i = \frac{x_i - \mu}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}}$$

The division by the standard deviation and the multiplication with the inverse of the sample covariance matrix have the same purpose: to standardise the difference to the sample mean/centroid by the variation in the respective variable(s).

Introduction to multivariate analysis, ordination and PCA

Contents

1. Learning targets and introduction to multivariate analysis
- 2. Overview ordination**
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Ordination: Introduction

- Extraction of new axes from high dimensional data that sequentially maximise the variance
 - Dimension reduction (e.g. omission of axes that capture low amount of variance)
 - Aggregation of variables into gradients
 - Graphical representation in lower dimension
- Unconstrained ordination: extraction without consideration of variables outside of data set
- Constrained ordination: extraction of axes that are explained by variables of second data set

16

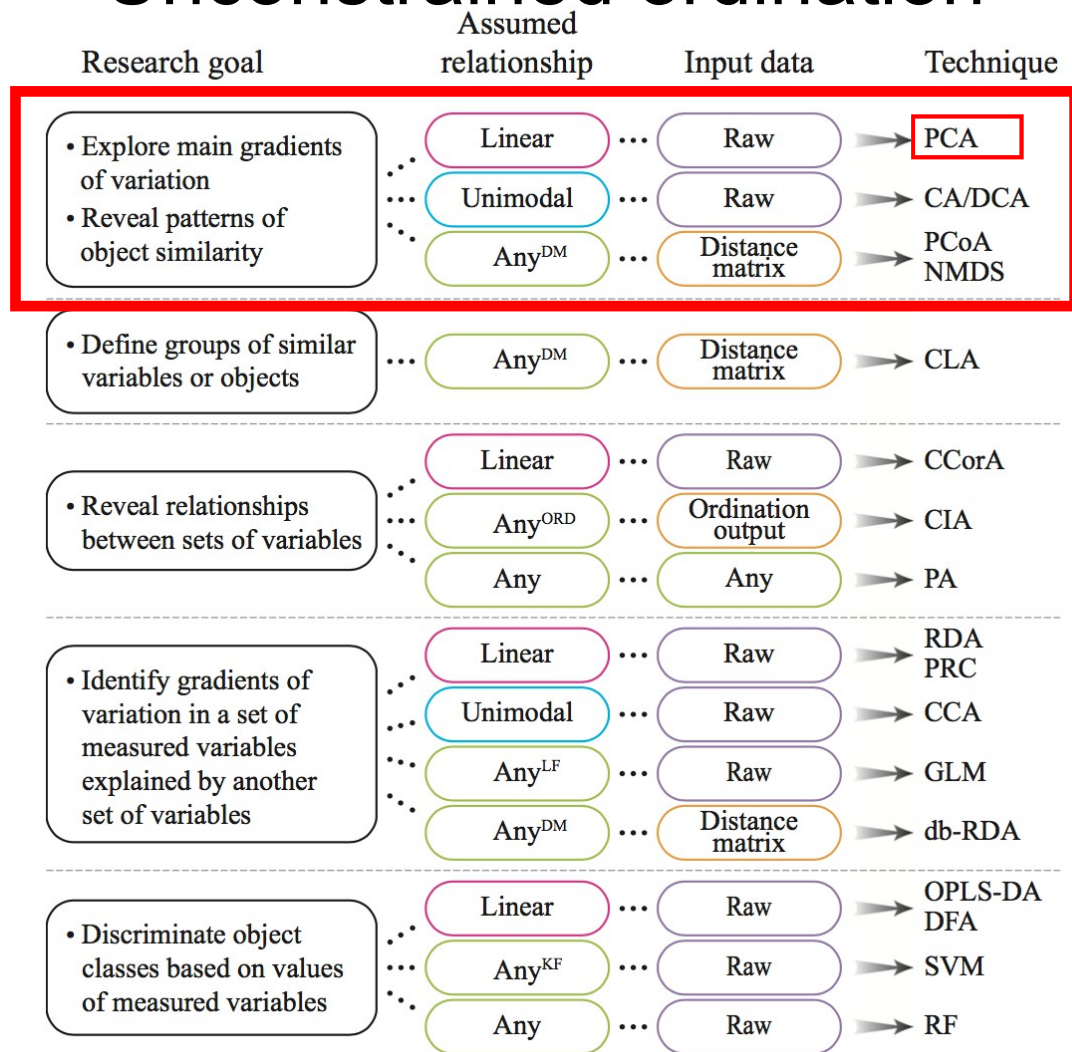
Very readable introductions into multivariate methods are:
Ramette, A. (2007) Multivariate analyses in microbial ecology. *FEMS Microbiology Ecology*, 62, 142-160. Freely available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2121141/pdf/fem0062-0142.pdf>
Paliy O. & Shankar V. (2016) Application of multivariate statistical techniques in microbial ecology. *Molecular Ecology* 25, 1032–1057. <http://onlinelibrary.wiley.com/doi/10.1111/mec.13536/abstract>

with a strong focus on ecotoxicology and a bit outdated:
van den Brink, P. J.; van den Brink, N. W.; Ter Braak, C. J.F. (2003): Multivariate analysis of ecotoxicological data using ordination: demonstrations of utility on the basis of various examples. *Australasian Journal of Ecotoxicology* 9, 141–156. Freely available at:
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.77.403&rep=rep1&type=pdf>

and finally a paper discussing different multivariate techniques in the context of diversity research:
Anderson, M. J.; Crist, T. O.; Chase, J. M.; Vellend, M.; Inouye, B. D.; Freestone, A. L. et al. (2011): Navigating the multiple meanings of diversity: a roadmap for the practicing ecologist. *Ecology Letters* 14, 19–28. <http://onlinelibrary.wiley.com/doi/10.1111/j.1461-0248.2010.01552.x/full>

For novel developments, see:
Warton D.I., Foster S.D., De'ath G., Stoklosa J. & Dunstan P.K. (2015) Model-based thinking for community ecology. *Plant Ecology* 216, 669–682. Freely available at:
https://www.researchgate.net/publication/276481409_Model-based_thinking_for_community_ecology
Warton D.I., Blanchet F.G., O'Hara R.B., Ovaskainen O., Taskinen S., Walker S.C., et al. (2015) So Many Variables: Joint Modeling in Community Ecology. *Trends in Ecology & Evolution*. 30, 766-779
<https://www.sciencedirect.com/science/article/pii/S0169534715002402>

Unconstrained ordination



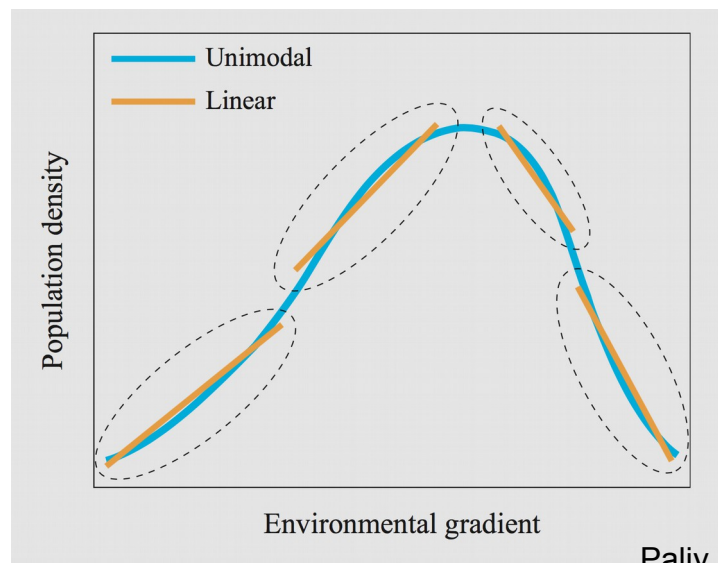
17

Paliy & Shankar 2016 *Mol. Ecol.* 25: 1032

Paliy O. & Shankar V. (2016) Application of multivariate statistical techniques in microbial ecology. *Molecular Ecology* 25, 1032–1057.

Ordination: Overview

Shape of response	Linear	Unimodal	Any
Unconstrained methods (examples)	PCA	CA	Distance-based: NMDS; GAM-based: UAO (U-VGAM)
Constrained methods (examples)	RDA	CCA	Distance-based: db-RDA; GAM-based: CAO (RR-VGAM)



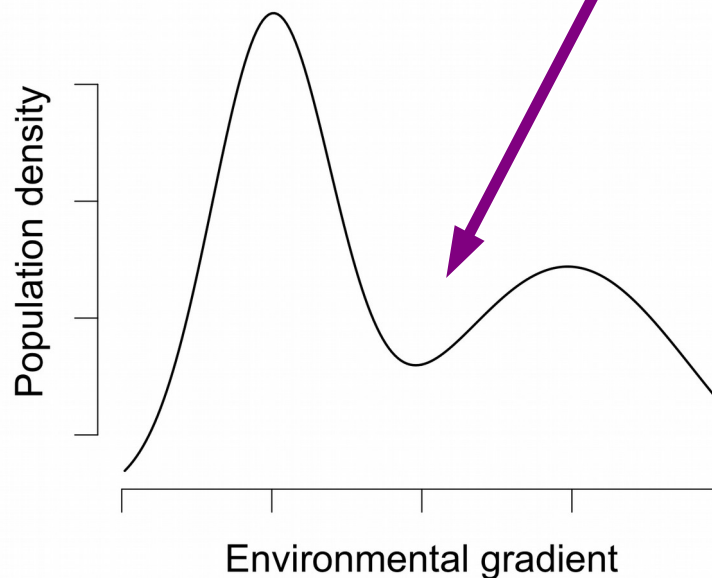
18

Paliy & Shankar 2016 *Mol Ecol*:1032

PCA: Principal Component Analysis
 RDA: Redundancy Discriminant Analysis
 CA: Correspondence Analysis
 CCA: Canonical Correspondence Analysis
 NMDS: Non-Metric Multidimensional Scaling
 db-RDA: distance-based RDA
 UAO: Unconstrained Additive Ordination
 CAO: Constrained Additive Ordination
 GAM: Generalised Additive Model
 VGAM: Vectorised Generalised Additive Model
 U-VGAM: Unconstrained-VGAM
 RR-VGAM: Reduced Rank-VGAM

Ordination: Overview

Shape of response	Linear	Unimodal	Any
Unconstrained methods (examples)	PCA	CA	Distance-based: NMDS; GAM-based: UAO (U-VGAM)
Constrained methods (examples)	RDA	CCA	Distance-based: db-RDA; GAM-based: CAO (RR-VGAM)



19

PCA: Principal Component Analysis
 RDA: Redundancy Discriminant Analysis
 CA: Correspondence Analysis
 CCA: Canonical Correspondence Analysis
 NMDS: Non-Metric Multidimensional Scaling
 db-RDA: distance-based RDA
 UAO: Unconstrained Additive Ordination
 CAO: Constrained Additive Ordination
 GAM: Generalised Additive Model
 VGAM: Vectorised Generalised Additive Model
 U-VGAM: Unconstrained-VGAM
 RR-VGAM: Reduced Rank-VGAM

Introduction to multivariate analysis, ordination and PCA

Contents

1. Learning targets and introduction to multivariate analysis
2. Overview ordination
- 3. Introduction to PCA and mathematical background**
4. PCA results and interpretation
5. PCA diagnosis and extensions, brief tutorial

Principal Component Analysis

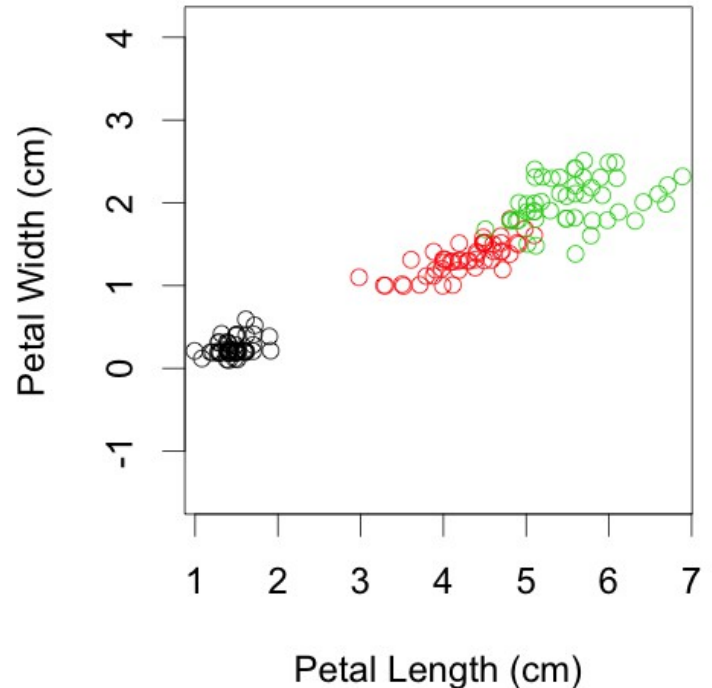
Example-based introduction

Iris data set: sepal length & width, petal length & width for 50 flowers from 3 species of Iris



<http://de.wikipedia.org/wiki/Schwertlilien>

Aim: Represent as much variance as possible on first few axes

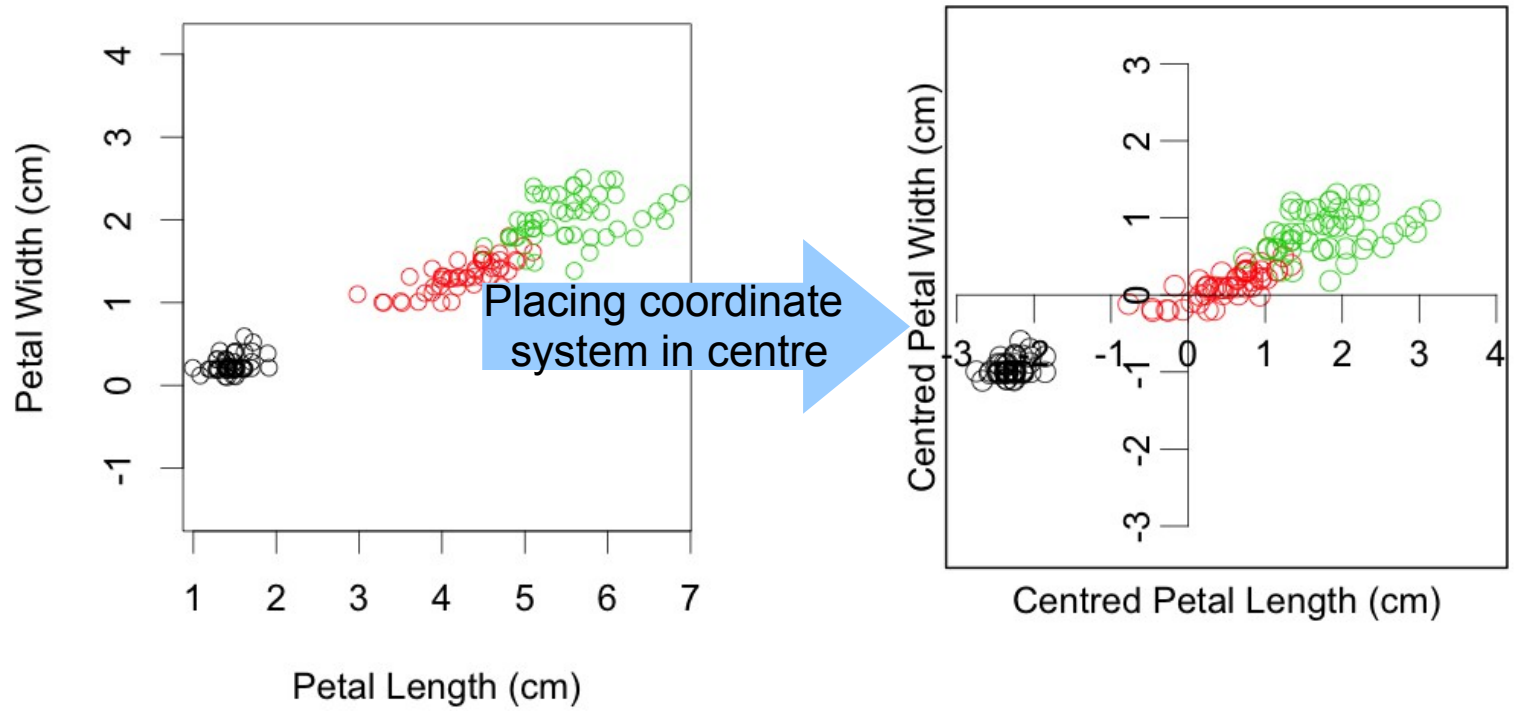


21

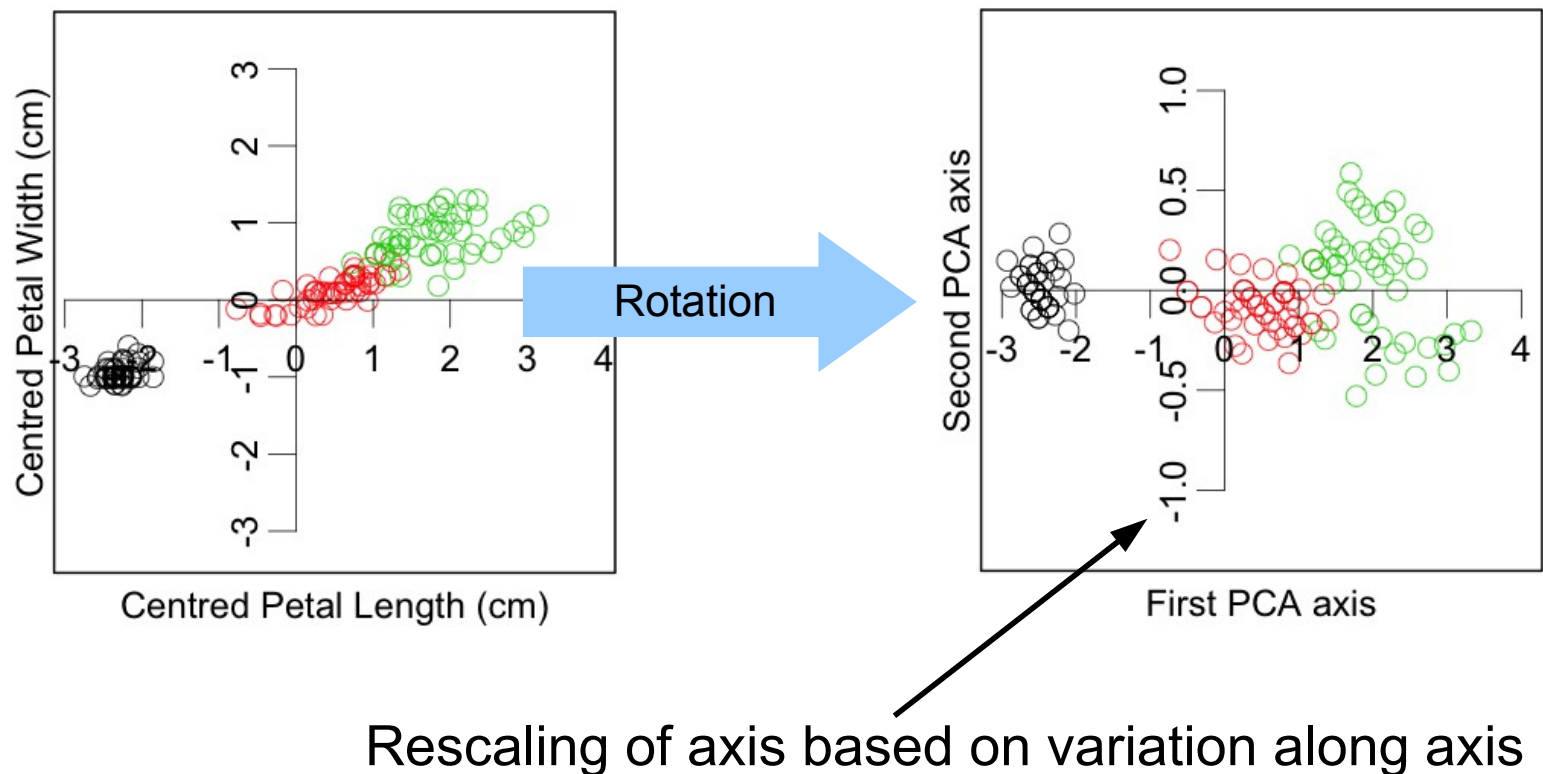
The example has been taken from:
<http://blogs.nature.com/boboh/2012/01/17/pca-and-pcoa-explained>
(link seems not to work anymore)

PCA can be regarded as a coordinate system with linearly independent axes that is placed in the centre of the data points. The axes are then rotated until the first axis explains the maximum variance, the second axis the highest remaining variance and so on. The number of axes always equals the number of variables. PCA is motivated by the expectation that the first few axes explain the major part of the total variation. The data are centred during analysis (the mean is subtracted) to facilitate interpretation.

Introduction to PCA



Introduction to PCA



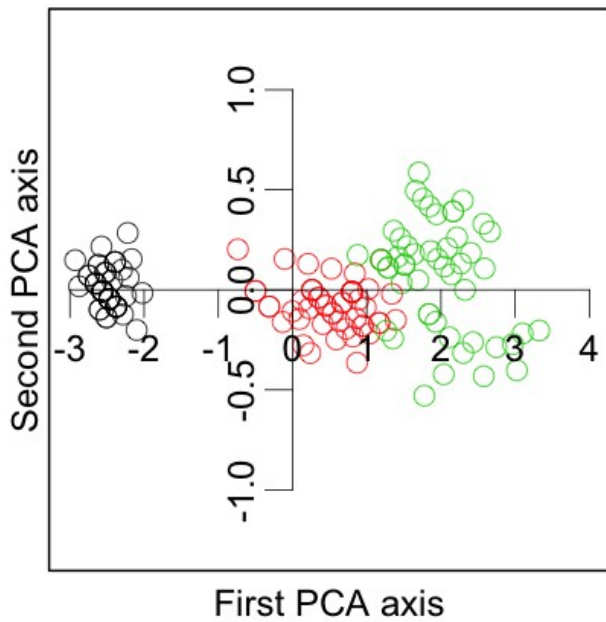
23

The term scaling can be used with two different meanings in the context of PCA and ordination in general. First, it refers to the question whether predictors are standardised to unit variance before analysis. Second, the PCA results can be scaled either to display the distances between objects or to display the relationship between descriptors (see Borcard 2011: 121 and Jolliffe (2002): 90ff for details). Here, we use the term for the second meaning.

Borcard D., Gillet F. & Legendre P. (2011) Numerical ecology with R, Springer, New York.

Jolliffe I.T. (2002) Principal component analysis, 2nd ed. Springer, New York.

Results of PCA



Variances

Petal Width 0.58

Petal Length 3.12

Results

Total Variation: 3.70

Importance of components:

	PC1	PC2
Eigenvalue	3.66	0.04
Proportion Explained	0.99	0.01
Cumulative Proportion	0.99	1.00

What is an Eigenvalue?

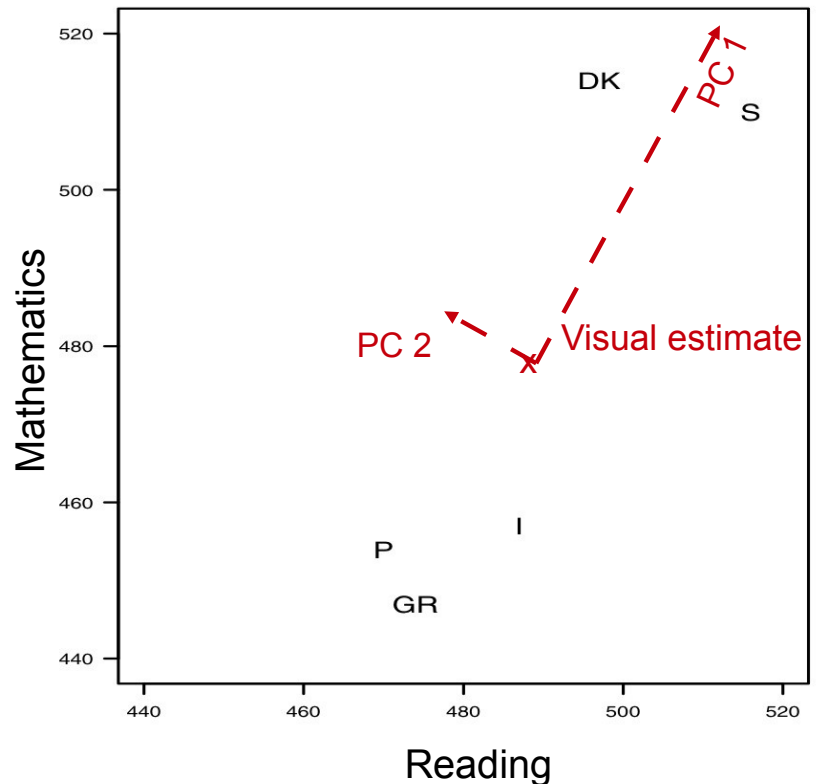
First axis explains 99% of variation!

Mathematical background of PCA

Country	Reading	Mathematics
DK	497	514
GR	474	447
I	487	457
P	470	454
S	516	510

Centering leads to

$$\tilde{X} = \begin{pmatrix} 8.2 & 37.6 \\ -14.8 & -29.4 \\ -1.8 & -19.4 \\ -18.8 & -22.4 \\ 27.2 & 33.6 \end{pmatrix}$$



Search for first axis with maximum variation!

The \sim on top of X indicates that it is an estimate for a transformed variable (centering represents a transformation). Thus, it is equivalent to the $^{\wedge}$ symbol for a non-transformed variable.

Handl, A. 2010: Multivariate Analysemethoden: Theorie und Praxis multivariater Verfahren unter besonderer Berücksichtigung von S-PLUS. Springer, Berlin.

Mathematical background of PCA

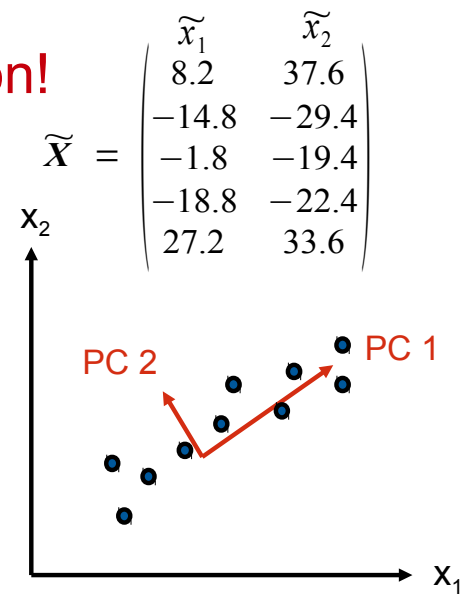
Search for first axis with maximum variation!

Scores on new axis given as

$$PC_1 = a_1 \tilde{x}_1 + a_2 \tilde{x}_2$$

and maximum variation means:

$$\max_{a_1, a_2} \text{Var}(a_1 \tilde{x}_1 + a_2 \tilde{x}_2)$$



Generalise problem: define a_1, a_2 as elements of vector a and likewise \tilde{x}_1, \tilde{x}_2 of matrix \tilde{X} : $\max_a \text{Var}(a \tilde{X})$

Trivial solution: choose high values for a_1, a_2, \dots, a_n

→ introduce condition: $a_1^2 + a_2^2 + \dots + a_n^2 = 1$

26

The score of objects on the first PC axis are obtained through a linear combination of the scores from the initial axes (after centering). The general equation for q dimensions for the first PC axis is:

$$PC_1 = a_1 x_1 + a_2 x_2 + \dots + a_q x_q$$

As we could artificially inflate the variance, we introduce a condition for the coefficients of the linear combination.

Regarding the Pisa data, for $a_1 = 1$ und $a_2 = 0$ the variance is given by the variance for \tilde{x}_1 . If you set $a_1 = 0.6$ and $a_2 = 0.8$ (meeting the condition: $0.6^2 + 0.8^2 = 1$), this results in a variance of 1317, which is higher than the individual variances of 1071 and 346. We need to find the values for a_1 and a_2 (under condition outlined above) that maximise variation.

Mathematical background of PCA

Solve:

$$\max_a \text{Var}(a\widetilde{X}) \text{ with } a^T a = 1$$

This can be expressed as (see Handl 2010: 79)

$$\max_a (a^T \Sigma a) \text{ with } a^T a = 1$$

Covariance matrix

Using the Lagrange function yields:

$$L(a, \lambda) = a^T \Sigma a - \lambda (a^T a - 1)$$

$$\frac{\partial L(a, \lambda)}{\partial a} = 2 \Sigma a - 2 \lambda a \longrightarrow \text{Eigenvalue problem}$$

$\Sigma a = \lambda a$

$$\frac{\partial L(a, \lambda)}{\partial \lambda} = 1 - a^T a$$

27

The description here uses the known variance-covariance matrix Σ , in other words the variance-covariance matrix of the statistical population. In practice, we usually have to estimate this matrix from the sample data and use this sample variance-covariance matrix \mathbf{S} . If the data are standardized before analysis (divided by the variance) then this matrix equals the correlation matrix \mathbf{R} .

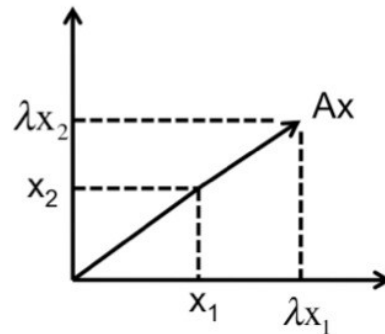
Mathematical basics II: Eigenvalues

Idea: Conversion of a matrix into a matrix with linear independent variables

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \dots & a_{22} & \dots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \xrightarrow{\text{Conversion}} \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_n \end{pmatrix}$$

Eigenvalue problem: $Ax = \lambda x$ Eigenvector
Eigenvalue

Eigenvectors form canonical basis and are only stretched or shrunk by λ when multiplied with A .



28

Eigenvalues and eigenvectors play an important role in the mathematics behind most ordination methods. The eigenvectors and eigenvalues of matrix A can be obtained by solving the (special) eigenvalue problem. It is important to understand that the eigenvalue problem can also be rewritten as linear function $f(x) = \lambda x$ and that this vector x is then only stretched by the factor λ for all $f(x)$. The matrix with the eigenvalues is also termed “canonical” form, a label encountered in many methods (e.g. canonical correspondence analysis, canonical correlation etc.).

Geometrically, the multiplication of a matrix with a vector, where the vector defines a point in a source coordinate system, rotates and stretches the vector to a new position. Consider the matrix $\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$ that defines the rotation/stretching of a point $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

Examples for related matrices:

$$\begin{pmatrix} 5 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{Stretching } x_1 \qquad \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \quad 180 \text{ degree rotation}$$

Mathematical basics II: Eigenvalues

$$Ax = \lambda x \Leftrightarrow Ax - \lambda x = 0$$

$$\Leftrightarrow (A - \lambda E)x = 0$$

$$\Leftrightarrow \begin{pmatrix} a_{11} - \lambda_1 & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{m1} & \dots & a_{mn} - \lambda_n \end{pmatrix} \begin{pmatrix} x_1 \\ \dots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ \dots \\ 0 \end{pmatrix} \quad \text{Homogeneous linear equation system (HLS)}$$

$$\text{Ignore trivial solution: } x = 0 \quad \Rightarrow \quad A - \lambda E = 0$$

The given HLS has only a non-trivial solution if the columns of $A - \lambda E$ are linearly dependent, which is the case if the determinant = 0.

$$\Rightarrow \det(A - \lambda E) = 0 \Leftrightarrow |A - \lambda E| = 0$$

29

E is the identity matrix.

A homogenous linear equation system $Ax = 0$ has the unique solution $x_1 = x_2 = x_3 = \dots = x_n = 0$ for a regular matrix (symmetric matrix for which the inverse exists). If the matrix is singular, at least one row is a linear combination of the other rows and consequently at least one $x \neq 0$. However, there is no unique solution in this case. The determinant (a special function, which assigns a unique number to every $n \times n$ matrix) for singular matrices is 0. This means that for the non-trivial case, the determinant is 0, otherwise all rows would be linearly independent (which also means that they do not share any variance and a PCA would not produce meaningful results).

In PCA, the variance-covariance or correlation matrix is used, which are symmetric. However, the eigenvalues of a non-symmetric matrix can also be computed by singular value decomposition (special case of the Schur decomposition). This is implemented in R with `svd()`.

Example I: Calculation of Eigenvalues

Sample Variance-Covariance matrix from Pisa example:

$$\mathbf{S} = \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix}$$

Following $|\mathbf{A} - \lambda \mathbf{E}| = 0$ we obtain:

$$\left| \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right| = 0 \Leftrightarrow$$

$$\left| \begin{pmatrix} 345.7 & 528.35 \\ 528.35 & 1071.30 \end{pmatrix} - \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \right| = 0 \Leftrightarrow$$

$$\begin{vmatrix} 345.7 - \lambda & 528.35 \\ 528.35 & 1071.30 - \lambda \end{vmatrix} = 0 \Leftrightarrow$$

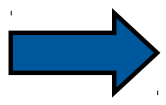
$$(345.7 - \lambda)(1071.30 - \lambda) - 528.35^2 = 0$$

30

For small, symmetric matrices, the determinant can be calculated by hand. For example The determinant of a matrix $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$ is given as $a d - b c$.

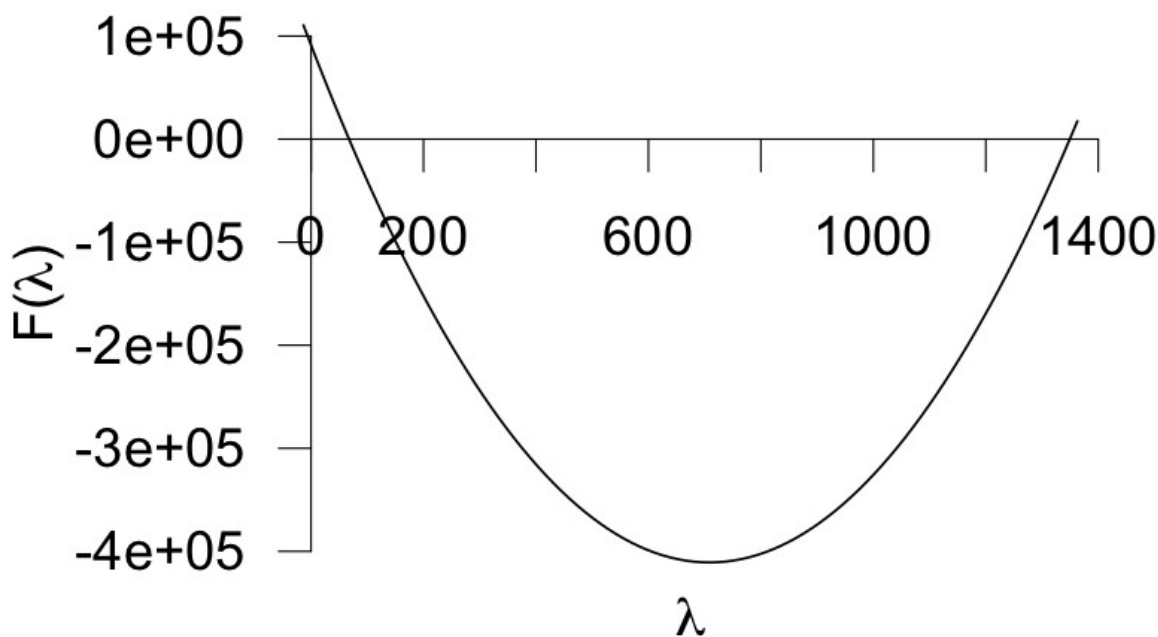
Example I: Calculation of Eigenvalues

$$(345.7 - \lambda)(1071.30 - \lambda) - 528.35^2 = 0$$



Characteristical polynom

$$\lambda^2 - 1417\lambda + 91194.69 = 0$$



31

Handl 2010: 123

In the quadratic case, the eigenvalues (λ) can easily be found using the pq formula:

$$\lambda_{1,2} = -\frac{p}{2} \pm \sqrt{\left(\frac{p}{2}\right)^2 - q}$$

Evaluating the equation for our numbers yields to:

$$\lambda_{1,2} = -\frac{-1417}{2} \pm \sqrt{\left(\frac{-1417}{2}\right)^2 - 91194.69}$$

Thus, λ_1 yields to 1349.42 and λ_2 yields to 67.6.

Handl, A. 2010: Multivariate Analysemethoden: Theorie und Praxis multivariater Verfahren unter besonderer Berücksichtigung von S-PLUS. Springer, Berlin.

Example II: Calculation of Eigenvalues and -vectors

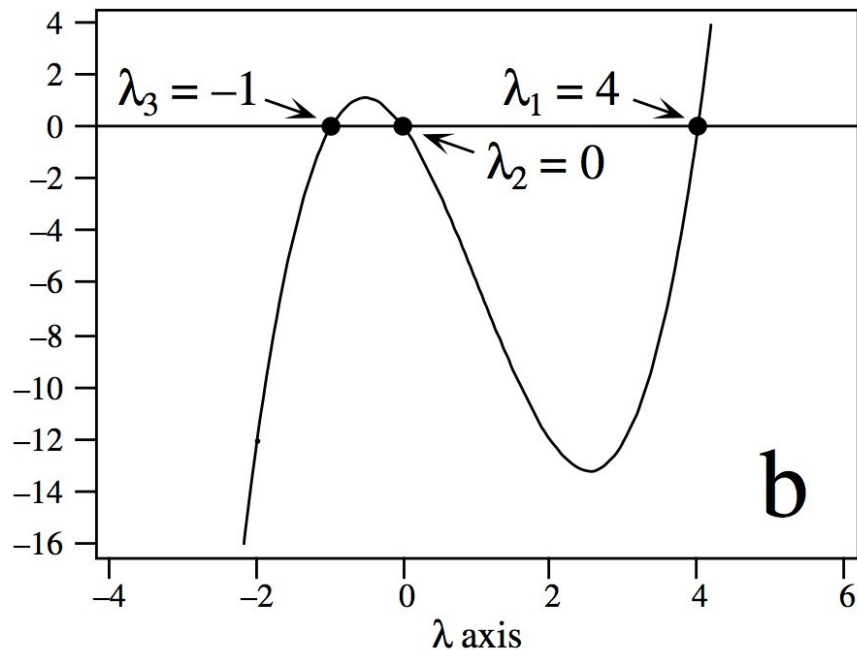
$$\begin{pmatrix} 1 & 3 & -1 \\ 0 & 1 & 2 \\ 1 & 4 & 1 \end{pmatrix}$$

EV calculation
following Sarrus

Characteristical polynom

$$\lambda^3 - 3\lambda^2 - 4\lambda = 0$$

**Eigenvalues λ : 4,
0 and -1**



Note that this example is only to demonstrate the calculation of eigenvalues and vectors, but is not based on real data. The eigenvalues of a PCA would usually all be positive.

Legendre P. & Legendre L. (2012) Numerical ecology, 3rd English ed. Elsevier, Amsterdam; Boston.

Example II: Calculation of Eigenvalues and -vectors

Calculation of eigenvector for $\lambda = 4$

$$(A - \lambda E)x = 0$$

$$\left(\begin{pmatrix} 1 & 3 & -1 \\ 0 & 1 & 2 \\ 1 & 4 & 1 \end{pmatrix} - \lambda_1 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0 \Leftrightarrow \begin{pmatrix} 1-4 & 3 & -1 \\ 0 & 1-4 & 2 \\ 1 & 4 & 1-4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$\begin{array}{cccc} -3x_1 & 3x_2 & -x_3 & 0 \\ 0 & -3x_2 & 2x_3 & = 0 \\ x_1 & 4x_2 & -3x_3 & 0 \end{array} \Leftrightarrow \begin{array}{cccc} -3x_1 & 0 & x_3 & 0 \\ 0 & 1.5x_2 & 0 & = x_3 \\ x_1 & 4x_2 & -3x_3 & 0 \end{array}$$

Matrix is singular (no unique solution)

→ fix value of one variable e.g. $x_1 = 1$.

Example II: Calculation of Eigenvalues and -vectors

$$x_1=1 \Rightarrow \begin{pmatrix} -3 & 0 & 0 \\ 0 & 1.5x_2 & 0 \\ 1 & 4x_2 & -3x_3 \end{pmatrix} = \begin{pmatrix} -x_3 \\ x_3 \\ 0 \end{pmatrix} \Rightarrow x_1=1; x_2=2; x_3=3$$

Calculation of eigenvectors for all eigenvalues yields the following matrix of eigenvectors (or multiples of columns):

$$\begin{pmatrix} 1 & 7 & 2 \\ 2 & -2 & -1 \\ 3 & 1 & 1 \end{pmatrix}$$

Eigenvalues: 4; 0 and -1

34

As an exercise, compute the eigenvectors for the eigenvalues 0 and -1.

Note that the condition $a^T a = 1$ will typically lead to values in the matrix of eigenvectors < 1 , contrary to the given example. Moreover, the eigenvalues are typically positive.

Introduction to multivariate analysis, ordination and PCA

Contents

1. Learning targets and introduction to multivariate analysis
2. Overview ordination
3. Introduction to PCA and mathematical background
- 4. PCA results and interpretation**
5. PCA diagnosis and extensions, brief tutorial

Results of PCA II

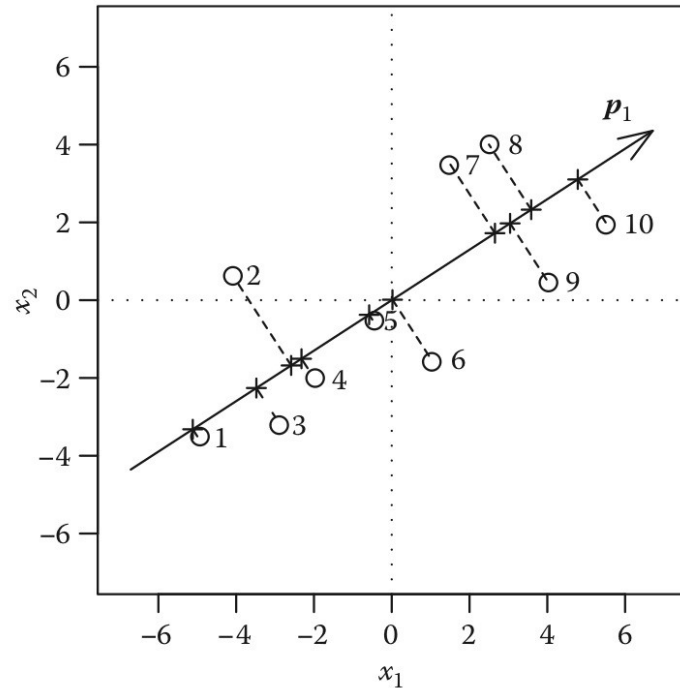
PC Scores

Demo Example for PCA with 10 Objects and Two Mean-Centered Variables x_1 and x_2

i	x_1	x_2	t_1	t_2
1	-5.0	-3.5	-6.10	-0.21
2	-4.0	0.5	-3.08	2.60
3	-3.0	-3.0	-4.15	-0.88
4	-2.0	-2.0	-2.77	-0.59
5	-0.5	-0.5	-0.69	-0.15
6	1.0	-1.5	0.02	-1.80
7	1.5	3.5	3.16	2.12
8	2.5	4.0	4.27	1.99
9	4.0	0.5	3.63	-1.76
10	5.5	2.0	5.70	-1.32
\bar{x}	0.00	0.00	0.00	0.00
v	12.22	6.72	16.22	2.72
$v\%$	64.52	35.48	85.64	14.36

Note: i , Object number; t_1 and t_2 are the PCA scores of PC1 and PC2, respectively; \bar{x} , mean; v , variance; $v\%$, variance in percent of total variance.

PC scores result from multiplication of scores from initial axes with eigenvectors

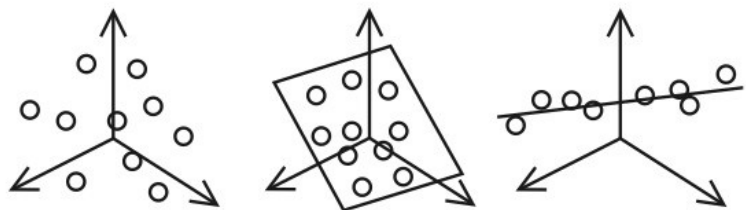


36 Varmuza & Filzmoser 2009: Chapter 3

The figure displays the scores t_1 for the points i on the first PC.

Number of principal components

- Number of descriptors/explanatory variables determines number of eigenvalues and thus principal components
- Largest eigenvalue (and the corresponding eigenvector) explains highest share of total variance
- Aim is to represent the major variation with a few principal components → How many components are needed?



Number of variables

3

3

3

Number of relevant components
= intrinsic dimensionality

3

2

1

How many principal components needed?

Some criteria to evaluate the optimal number of axes:

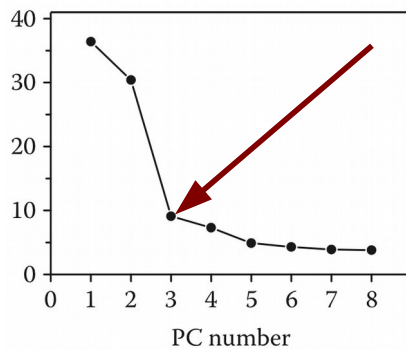
1. Sum criterion

$$\frac{\sum_{i=1}^r \lambda_i}{\sum_{j=1}^p \lambda_j} \geq \alpha$$

2. Broken-Stick criterion

$$\lambda_i > \frac{1}{p} \sum_{i=1}^p \frac{1}{i}$$

3. Scree plot



4. Cross-validation

$$\min_S \text{MSEP}(S) = \frac{1}{I K} \sum_{i=1}^I \sum_{k=1}^K \left(x_{i,k} - (\hat{x}_{i,k})^{(S)} \right)^2$$

For the matrix $X_{I \times K}$, we search the number of PC S minimizing the mean square error of prediction (MSEP).

38

Varmuza & Filzmoser 2009: Chapter 3, Jolliffe 2002: Chapter 6, Josse & Husson 2012

The sum criterion and the scree plot are relatively simple rules without a thorough statistical foundation. In case of the sum criterion, the first $i = 1, \dots, r$ eigenvalues that are larger than α (proportion of cumulative variance) are selected. In other words, the minimal number of eigenvalues that capture at least % of total variance. Typically, a value between 0.7 and 0.9 is chosen for α . Generally, the optimal α will decrease with an increase in the number of variables p in the data set.

The Broken-Stick model is based on the idea that we break a stick of unit length into p pieces, where p is given by the number of eigenvalues/principal components. If the resulting pieces are sorted in decreasing order with respect to their length, the i -th length of the stick is given by the formula. The i -th eigenvalue should be larger than the i -th length from the broken-stick model, because otherwise it would not capture more variance than a random process. A modelling study (Jackson 1993) found the broken stick model among the most reliable methods for the determination of the number of principal components, though novel methods have been developed in the last two decades. Example for broken stick model: Imagine we have 2 variables, then the largest eigenvalue should be higher than $1/2 (1+1/2) = 0.75$. The second eigenvalue should then be higher than $1/2 * 1/2 = 0.25$.

Cross-validation is computationally costly and may not be applicable for large data sets (for example see Saccenti & Camacho (2015)). Hence, approximations of cross-validation have been developed (Josse & Husson 2012), though they may be less reliable for some data sets. The performance of methods varies with the properties of the data sets, for an overview see Camacho & Ferrer (2014) and Saccenti & Camacho (2015). Note that statistical tests are in many cases less reliable than cross-validation and are not discussed here, but within the given references.

Camacho J. & Ferrer A. (2014) Cross-validation in PCA models with the element-wise k-fold (ekf) algorithm: Practical aspects. *Chemometrics and Intelligent Laboratory Systems* 131, 37–50.

Jackson, D.A. (1993) Stopping rules in principal components-analysis - a comparison of heuristic and statistical approaches. *Ecology*, 74, 2204–2214.

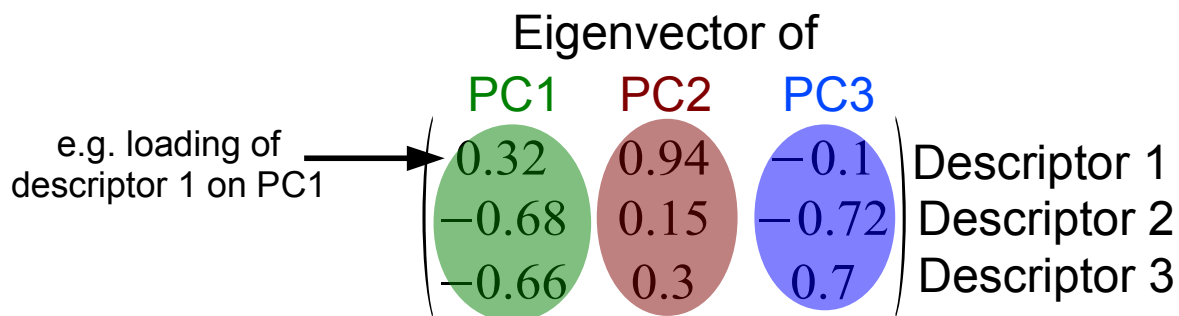
Jolliffe I.T. (2002) Principal component analysis, 2nd edn. Springer, New York.

Josse J. & Husson F. (2012) Selecting the number of components in principal component analysis using cross-validation approximations. *Computational Statistics & Data Analysis* 56, 1869–1879.

Saccenti E. & Camacho J. (2015) Determining the number of components in principal components analysis: A comparison of statistical, crossvalidation and approximated methods. *Chemometrics and Intelligent Laboratory Systems* 149, Part A, 99–116.

Importance of descriptor for PC axes

- Elements of eigenvector matrix = 'loadings', indicate weight of original descriptor on PCs



- Easier to interpret: Correlation loadings $r_i = a_i \sqrt{\lambda_i}$
- Interpretation of descriptor importance complicated if many variables load on PC
- Sparse PCA – introduces penalty term (cf. LASSO):

$$\max_a \text{Var}(a\tilde{\mathbf{X}}) - \lambda \|a\| \quad \text{with} \quad a^T a = 1$$

39

a

The correlation loadings are equivalent to correlating the original descriptors with the PC scores.

Note that the λ used in the context of sparse PCA relates to the penalty term and not to eigenvalues. This must not be confused (i.e. the λ in the context of correlation loadings are the eigenvalues).

$\|a\|$ is the norm of a (function that assigns a value to a vector)

For further details refer to the course material on the LASSO and the following articles:

Croux C., Filzmoser P. & Fritz H. (2011) Robust sparse principal component analysis. TU Vienna, Vienna.

<http://www.statistik.tuwien.ac.at/forschung/SM/SM-2011-2complete.pdf> (more or less identical to: Croux C., Filzmoser P. & Fritz H. (2013) Robust Sparse Principal Component Analysis. *Technometrics* 55, 202–214.)

Zou H., Hastie T. & Tibshirani R. (2006) Sparse Principal Component Analysis. *Journal of Computational and Graphical Statistics* 15, 265–286.

Introduction to multivariate analysis, ordination and PCA

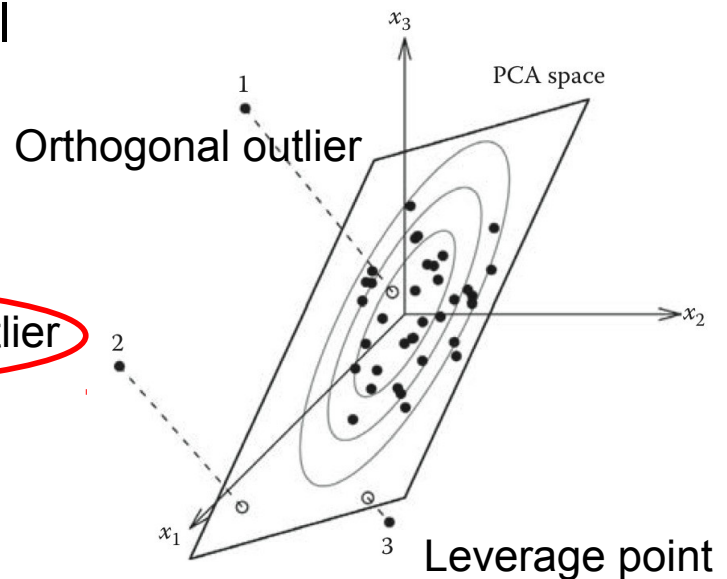
Contents

1. Learning targets and introduction to multivariate analysis
2. Overview ordination
3. Introduction to PCA and mathematical background
4. PCA results and interpretation
- 5. PCA diagnosis and extensions, brief tutorial**

PCA assumptions and diagnosis

- Independence of observations (temporal and spatial independence)
- multivariate normality (depending on research question)
- no serious outliers, otherwise use robust PCA (see Varmuza & Filzmoser 2009: chapter 3.5)
- Computation of orthogonal distance (OD) and score distance (SD)

Leverage point and orthogonal outlier
Problematic



41 Varmuza & Filzmoser 2009, ch.3

In the case of multivariate normal distribution of the data, the score distances can be approximated with a χ^2 distribution to detect leverage points. Orthogonal distances can be approximated with a normal distribution. See Varmuza & Filzmoser 2009, chapter 3.7.3 for details.

Multivariate normality is not required if the main aim of the PCA is visual exploration. In this case, also discrete variables can be reliably displayed, as long as the data values are numerical and the distances between them are meaningful. However, if the PCs are used in subsequent regression analysis (Principal component regression), multivariate normality is required. Similarly, hypothesis tests for the number of PCs require multivariate normality.

PCA assumptions and limitations

- Linear gradient of descriptors (rarely the case for species data, but often for environmental data)
- Euclidean distances inappropriate for species data with many double zeros (joint absences)
- Adding noise variables to data increases fraction of variance on first axis, but has no meaning
- Useful technique in exploratory analysis and for analysing environmental data (or other data with linear gradients)
- Best results for large n and high $n:p$ (p = descriptors)

42

The most important assumption is that of a linear gradient of the descriptors over the object space. If this assumption is not met, as for example in the case of species data, this can lead to serious problems (see the example in the related R course script. It nicely displays how unimodal gradients influence the PCA). Non-linearity can be checked using scatter plots of the descriptors before analysis. Strongly skewed data and non-linear relationships may cause problems. For such data, different ordination methods are available or the data can be transformed. See Legendre & Legendre 2012: 450ff and Borcard et al. 2011: 130f. for details.

Borcard D., Gillet F. & Legendre P. (2011) Numerical ecology with R, 1. Springer, New York, NY.

Legendre P. & Legendre L. (2012) Numerical ecology, 3rd English ed. Elsevier, Amsterdam; Boston.

Brief tutorial for PCA

1. Check if conditions for descriptors are met (quantitative, multivariate normality, linear)
2. Conduct PCA (or sparse PCA) on scaled descriptors unless they exhibit a similar variation and have been measured on similar scale
3. Check for outliers
4. Determine number of principal components
5. How informative are first two PCs?
6. Which descriptors contribute most to PCs?
7. Visualise and interpret

43

A detailed key for PCA is given in Legendre & Legendre 2012: 453

Legendre P. & Legendre L. (2012) Numerical ecology, 3rd English ed. Elsevier, Amsterdam; Boston.

Principal component regression

- Extract (unscaled) PC scores to use PCs as descriptors in multiple regression analysis
- PCs are orthogonal → fix for multicollinearity in regression analysis
- In low $n:p$ situations, the few last PCs are often removed to reduce number of predictors in regression → Can be problematic because low variance of PC does not imply low explanatory power for response variable
→ not necessarily a fix for low $n:p$ ratios

See Jolliffe (2002): 173 for examples where PCs that capture minor amounts of variance exhibit high explanatory power. This is explained when considering that the last few PCs represent the constant elements (low variance) across the descriptors.