

Molecular data in R

Phylogeny, evolution & R

Vojtěch Zeisek

Department of Botany, Faculty of Science, Charles University in Prague
Institute of Botany, Czech Academy of Sciences, Průhonice

<http://trapa.cz/>, zeisek@natur.cuni.cz

January 14 to 16, 2015



Outline I

① Introduction

② R

Installation

Let's start with R

Basic operations in R

Packages for our work

③ Data

Microsatellites

AFLP

Notes about data

DNA sequences, SNP

Export

④ Basic analysis

First look at the data

Statistics

Outline II

Genetic distances

Hierarchical clustering

AMOVA

MSN

NJ (and UPGMA) tree

PCoA

⑤ DAPC

⑥ SNP

PCA and NJ

⑦ Spatial analysis

Moran's I

sPCA

Monmonier

Mantel test

Geneland

Outline III

Plotting maps

⑧ Structure

Postprocessing

⑨ Alignment

MUSCLE

MAFFT

Clustal, MUSCLE and T-Coffee from ape

⑩ Trees

Tree manipulations

Seeing trees in the forest

MP

Notes about plotting the trees

⑪ Evolution

PIC

Phylogenetic autocorrelation

Outline IV

pPCA

Ancestral state reconstruction

12 The end

Resources

The end

What we will and what we will not do...

We will go through...

- Basic introduction into R
- Various methods and data for analysing phylogeny and evolution
 - DNA sequences, SNP, SSRs, AFLP, ...
 - NJ, UPGMA, PCoA, DAPC, Bayesian clustering, ML, maximum parsimony, ...
- Plotting, maps, ...
- Basic creation of scripts

We will not dig deep into...

- All differences among operating systems (OS; not big, but there are some...)
- Theory behind used methods (with some exceptions)
- Programming in R
- Other software related to the methods used but not directly related to R (with exceptions of applications called from R)

About R

- Project for Statistical Computing
- open-source — freely available with source code — anyone can use it and contribute its development
- Development is organised by non-governmental non-profit organisation from Wien
- Thousands of packages extending its functionality are available — all fields of computing
- Provides only command line interface — full control over analysis, easy to rerun and/or modify analysis in the future, easy creation of scripts for batch analysis
- Several projects provide convenient graphical interfaces (GUI)
- More details: <http://www.r-project.org/>

Graphical user interfaces

- Provide more comfortable interface for work with scripts (source code highlight), overview of loaded packages and variables, easier work with figures, ...
- RStudio <http://www.rstudio.com/> — probably the most common, multiplatform, very powerful
- RKWard <https://rkward.kde.org/> — feature very rich, developed mainly for Linux, although available also for another operating systems
- R commander (Rcmdr) <http://www.rcommander.com/> — multiplatform, not so rich as previous
- Tinn-R <http://nbcgib.uesc.br/lec/software/editores/tinn-r/en> — Windows only
- Pick one you like and install it...

RKWard

Soubory Otevřít Vložit Použití Profilovat Provedení Data Analýza Plots Distribuce Okna Nastavení Napovídka

Projekty

Vložit Ne-funkce Eukleje

Uzálohovat všechna prostředí Uzálohovat skryté předměty

Název Znázka Typ Údaje

Glossary.rnw Unknown

- Number integer data frame
- atp data frame
- sars data frame
- sars_genind data frame
- sars_genind2 data frame
- sars_genind3 data frame
- sars_genind4 data frame
- sars_genind5 data frame
- sars_genind6 data frame
- sars_genind7 data frame
- sars_genind8 data frame
- sars_genind9 data frame
- sars_genind10 data frame
- sars_genind11 data frame
- sars_genind12 data frame
- sars_genind13 data frame
- sars_genind14 data frame
- sars_genind15 data frame
- sars_genind16 data frame
- sars_genind17 data frame
- sars_genind18 data frame
- sars_genind19 data frame
- sars_genind20 data frame
- sars_genind21 data frame
- sars_genind22 data frame
- sars_genind23 data frame
- sars_genind24 data frame
- sars_genind25 data frame
- sars_genind26 data frame
- sars_genind27 data frame
- sars_genind28 data frame
- sars_genind29 data frame
- sars_genind30 data frame
- sars_genind31 data frame
- sars_genind32 data frame
- sars_genind33 data frame
- sars_genind34 data frame
- sars_genind35 data frame
- sars_genind36 data frame
- sars_genind37 data frame
- sars_genind38 data frame
- sars_genind39 data frame
- sars_genind40 data frame
- sars_genind41 data frame
- sars_genind42 data frame
- sars_genind43 data frame
- sars_genind44 data frame
- sars_genind45 data frame
- sars_genind46 data frame
- sars_genind47 data frame
- sars_genind48 data frame
- sars_genind49 data frame
- sars_genind50 data frame
- sars_genind51 data frame
- sars_genind52 data frame
- sars_genind53 data frame
- atp_kfind Unknown
- atp_kfind2 Unknown
- atp_kfind3 Unknown
- atp_kfind4 Unknown
- atp_kfind5 Unknown
- atp_kfind6 Unknown
- atp_kfind7 Unknown
- atp_kfind8 Unknown
- atp_kfind9 Unknown
- atp_kfind10 Unknown
- atp_kfind11 Unknown
- atp_kfind12 Unknown
- atp_kfind13 Unknown
- atp_kfind14 Unknown
- atp_kfind15 Unknown
- atp_kfind16 Unknown
- atp_kfind17 Unknown
- atp_kfind18 Unknown
- atp_kfind19 Unknown
- atp_kfind20 Unknown
- atp_kfind21 Unknown
- atp_kfind22 Unknown
- atp_kfind23 Unknown
- atp_kfind24 Unknown
- atp_kfind25 Unknown
- atp_kfind26 Unknown
- atp_kfind27 Unknown
- atp_kfind28 Unknown
- atp_kfind29 Unknown
- atp_kfind30 Unknown
- atp_kfind31 Unknown
- atp_kfind32 Unknown
- atp_kfind33 Unknown
- atp_kfind34 Unknown
- atp_kfind35 Unknown
- atp_kfind36 Unknown
- atp_kfind37 Unknown
- atp_kfind38 Unknown
- atp_kfind39 Unknown
- atp_kfind40 Unknown
- atp_kfind41 Unknown
- atp_kfind42 Unknown
- atp_kfind43 Unknown
- atp_kfind44 Unknown
- atp_kfind45 Unknown
- atp_kfind46 Unknown
- atp_kfind47 Unknown
- atp_kfind48 Unknown
- atp_kfind49 Unknown
- atp_kfind50 Unknown
- atp_kfind51 Unknown
- atp_kfind52 Unknown
- atp_kfind53 Unknown
- atp_kfind54 Unknown
- atp_kfind55 Unknown
- atp_kfind56 Unknown
- atp_kfind57 Unknown
- atp_kfind58 Unknown
- atp_kfind59 Unknown
- atp_kfind60 Unknown
- atp_kfind61 Unknown
- atp_kfind62 Unknown
- atp_kfind63 Unknown
- atp_kfind64 Unknown
- atp_kfind65 Unknown
- atp_kfind66 Unknown
- atp_kfind67 Unknown
- atp_kfind68 Unknown
- atp_kfind69 Unknown
- atp_kfind70 Unknown
- atp_kfind71 Unknown
- atp_kfind72 Unknown
- atp_kfind73 Unknown
- atp_kfind74 Unknown
- atp_kfind75 Unknown
- atp_kfind76 Unknown
- atp_kfind77 Unknown
- atp_kfind78 Unknown
- atp_kfind79 Unknown
- atp_kfind80 Unknown
- atp_kfind81 Unknown
- atp_kfind82 Unknown
- atp_kfind83 Unknown
- atp_kfind84 Unknown
- atp_kfind85 Unknown
- atp_kfind86 Unknown
- atp_kfind87 Unknown
- atp_kfind88 Unknown
- atp_kfind89 Unknown
- atp_kfind90 Unknown
- atp_kfind91 Unknown
- atp_kfind92 Unknown
- atp_kfind93 Unknown
- atp_kfind94 Unknown
- atp_kfind95 Unknown

cardamine_dapc.r

```

library(dape)
#> library(dape)
#> library(DAPC)
#> table_value(table(pop(atp_genind)), atp_kfind$group), col.lab=paste("Inferred cluster", 1:length(atp_kfind$size)),
#> 
```

Záložky Otevřít Použití Profilovat Data Analýza Plots Distribuce Okna Nastavení Napovídka

Cardamine_dapc.r - R kód

Dane získané zhozku (DAPC) - R kód

Inferred cluster 1

Inferred cluster 2

POP7568
POP7688
POP5449
POP6250
POP5478
POP5847
POP5555
POP6244
POP6344
POP6250
POP6350
POP6450
POP7251
POP7448
POP6746
POP7148
POP7147
POP6845
POP7046
POP6245
POP6345
POP6249
POP6349
POP6449
POP6045

RStudio

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Left Panel:** A tree view labeled "workspace" showing various objects like "Source on Base", "ssrs3n.locl.missingno", and "ssrs3n.locl.missingno".
- Middle Left Panel:** An "R Script" editor containing R code related to population genetics analysis.
- Middle Right Panel:** An "Environment" pane showing a list of loaded packages and their versions.
- Bottom Left Panel:** A "Console" window displaying R version information and a warning about the "Pumpkin Helmet".
- Bottom Center:** A "Update Packages" dialog box showing a list of packages with their installed and available versions.
- Bottom Right:** A navigation bar with icons for back, forward, search, and other functions.

MS Windows, Apple Mac OS X

- Got to <http://cran.at.r-project.org/>
- Download appropriate version and install as usual
- Download and install selected GUI
- Most of packages are available as pre-compiled and can be immediately installed from R
- Convenient, but usually not tuned for particular computer architecture (type of CPU)
- Usually there are some problems every time new version of OS is released — it takes time to modify and recompile packages for new version of OS
- You have to check for new version of R manually

Linux — general and Debian/Ubuntu

- R, and usually also GUI, is available in repositories — use standard package management according to distribution
- Linux repositories provide automatic updates
- Packages are also partially available in repositories and can be installed and updated as usual application or from R
- Packages commonly have to be compiled – R will do it automatically, just install basic Linux packages for building (in Debian and Ubuntu install package **build-essential**)
- Debian (and derivatives): follow instructions at
<http://cran.at.r-project.org/bin/linux/debian/>
- Ubuntu (and derivatives): follow instructions at
<http://cran.at.r-project.org/bin/linux/ubuntu/>
- As <**my.favorite.cran.mirror**> select **cran.at.r-project.org**

Linux — openSUSE

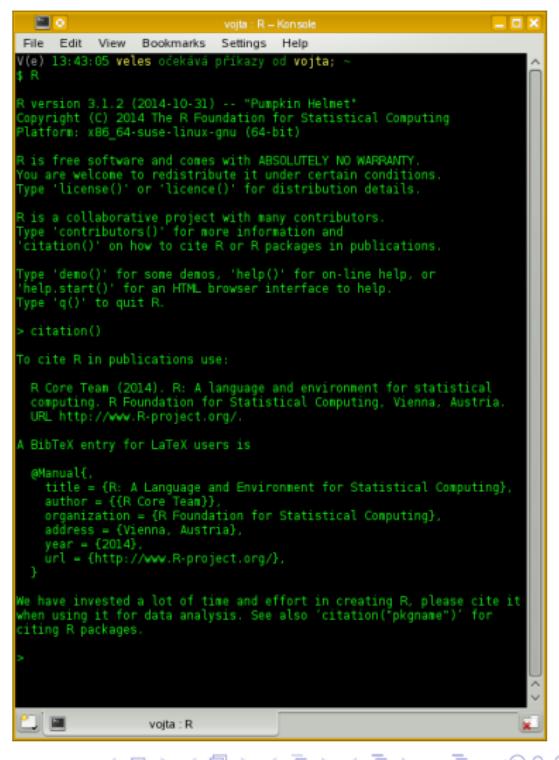
- See instructions at <http://cran.r-project.org/bin/linux/suse/>
- Add repositories for your distribution
 - <http://download.opensuse.org/repositories/devel:/languages:/R:/base/>
 - <http://download.opensuse.org/repositories/devel:/languages:/R:/patched/>
 - <http://download.opensuse.org/repositories/devel:/languages:/R:/released/>
- Install packages **R-patched** (the R) and **rstudio** and/or **rkward**
- Install package **patterns-openSUSE-devel_basis** for compilation of R packages when installing them from R (only some R packages are available in repositories)
- Compilation takes longer time and there are sometimes issues, but it can then provide higher performance...

Sources of R packages

- R CRAN <http://cran.r-project.org/> — main and largest source of R packages (over 6000 packages)
- Bioconductor <http://bioconductor.org/> — mainly bioinformatic packages, genomic data (~1000 packages)
- R-Forge <http://r-forge.r-project.org/> (over 1800 packages)
- RForge <http://www.rforge.net/> (much smaller)
- Omega Project for Statistical Computing <http://www.omegahat.org/>
- And more...
- Some packages are available from more resources — same name for function can be used in different packages — to distinguish them you can call functions like this: **muscle::read.fasta()** vs. **seqinr::read.fasta()** — call function **read.fasta()** from package **muscle** or **seqinr** (and their parameters can be different...)

First steps in R

- Linux: open any terminal, type **R** and hit Enter, Windows and Mac: find it as normal application in menu
- R is ruled by typing commands... :-)
- Use arrows up and down to navigate in history
- Ctrl+R works as reverse search — searches typed string in history of commands
- We will mostly use GUI



```
V(e) 13:43:05 veles očekává příkazy od vojta; ~
$ R
R version 3.1.2 (2014-10-31) -- "Pumpkin Helmet"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-suse-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'Citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> citation()
To cite R in publications use:

  R Core Team (2014). R: A language and environment for statistical
  computing. R Foundation for Statistical Computing, Vienna, Austria.
  URL http://www.R-project.org/.

A BibTeX entry for LaTeX users is

@Manual{,
  title = {R: A Language and Environment for Statistical Computing},
  author = {{R Core Team}},
  organization = {R Foundation for Statistical Computing},
  address = {Vienna, Austria},
  year = {2014},
  url = {http://www.R-project.org/},
}

We have invested a lot of time and effort in creating R, please cite it
when using it for data analysis. See also 'citation("pkgname")' for
citing R packages.

>
```

How it works

- General look of R commands:
- `function(argument1="SomeName", argument2=SomeVariable, argument3=8)`
- `ModifiedObject <- SomeFunction(argument1=MyData, argument2=TRUE)`
- Object is on the left: “`<-`” says to insert result of the function on the right into the object
- Functions have various parameters/arguments (in brackets, separated by commas): `argument=ItsValue`
- Arguments are named — if you keep order, no need to name them:
- `SomeFunction(MyData, TRUE, 123, "SomeName")`
- When only some of the arguments are in use, use the names (order doesn't matter any more)
- `SomeFunction(argument2=TRUE, argument3=123, argument1=MyData)`

Get help in R

```
1 # "#" marks comments - notes within code which are not executed
2 help(function) # Help for particular function (package must be loaded)
3 ?function # Help for particular function (package must be loaded)
4 ??SearchedTerm # Search for the term within all installed packages
5 help.search("searched phrase") # Search for the phrase within all
6   # installed packages - return list of hits sorted according to
7   # type and package (i.e. package::function)
8 install.package("sos") # More comprehensive search from package sos
9 library(sos)
10 findFN("function")# Search for function name
```

? shows help for questioned function:

- Name of the package (top left)
- Function name (headline)
- Description
- Usage
- Comments on arguments
- Details
- About author(s)
- References to cite
- Example code

Where we are?

- In Linux, R starts in current directory (use **cd** to change it before launching R)
- Set and check working directory in R:

```
1 setwd("/some/path/") # Or "~/...". In Windows "C:\..."  
2 getwd() # Verifies where we are  
3 dir() # Lists files and folders on the disk  
4 ls() # Lists currently available R objects
```

- In Windows (File | working directory) or in RStudio (Session | Set working directory) set it in menu or by above command
- R saves history of commands into **.Rhistory** file within working directory (by default hidden in Linux/Mac)
- When closing R by **q()** you can save all R data in **.RData** file and it can be loaded next time
- RStudio and RKWard help with this very much

Materials to help you...

- Follow all code we will use at
https://trapa.cz/sites/default/rcourse/course_commands.html
- Download the script from
https://trapa.cz/sites/default/rcourse/course_commands.r, use it and write your comments and notes to it during the course
- **Note:** Open the R script in some **good text editor** — showing syntax highlight, line numbers, etc. (**NO** Windows notepad); the file is in UTF-8 encoding and with Linux end of lines (so that too silly programs like Windows notepad won't be able to open it correctly)
- Information about the course:
<https://trapa.cz/en/course-molecular-data-r-2015> or
<https://trapa.cz/cs/kurz-molekularni-data-r>

Types of objects

- Vectors — numbers or characters
- Matrices — columns are of same type (numeric, character, etc.) and the same length
- Arrays — like matrices, but with possibly more dimensions
- Data frames — more general — columns can be of different type
- Lists — ordered collections of objects (vectors, matrices, ...) — not necessarily of the same type
- Factors — a vector of levels, e.g. populations, colours, etc.
- More “advanced” objects to store plots, genetic data, ...
 - Commonly called “S3” and “S4” objects
 - Technically commonly just lists putting together various information
 - We will meet many of them...
- Functions require particular object types — take care about it

Popular object classes (we are going to use) I

- `dist` — distance matrices
- `genind` — stores various genetic information for individuals
- `genpop` — like `genind`, but on population level
- `genlight` — variant of `genind` to store large multiple genomes
- `SNPbin` — stores large SNP data for single genome
- `DNAbin` — stores DNA sequences (aligned or not)
- `haplotype` — unique sequences from `DNAbin`
- `alignment` — aligned sequences (`seqinr`)
- `phyDat` — “preparation” of data for some phylogenetic analysis
- `loci` – extension of DF, stores information about loci
- `hclust` – output of hierarchical clustering, can be converted to `phylo`
- `treeshape` — derived from `hclust`

Popular object classes (we are going to use) II

- phylo — phylogenetic information, typically trees
- phylo4 — derived from phylo, S4 instead of S3
- matching — binary phylogenetic trees
- haplonet — networks without reticulation
- spca — results of sPCA
- pco; dudi — results of PCA, PCoA, ...
- dapc — results of DAPC
- and more...common task is converting among formats...
- ...not all formats are (easily) convertible among each other...

To get information about content of each data type see

```
getClassDef("genind") # Or any other class name
```

There are information about slots within that classes you can access.

Conversions among data types I

From	To	Command	Package
phylo	phylo4	as(x, "phylo4")	phylobase
phylo	matching	as.matching(x)	ape
phylo	treeshape	as.treeshape(x)	apTreeshape
phylo	hclust	as.hclust(x)	ape
phylo	prop.part	prop.part(x)	ape
phylo	splits	as.splits(x)	phangorn
phylo	evonet	evonet(x, from, to)	ape
phylo	network	as.network(x)	ape
phylo	igraph	as.igraph(x)	ape
phylo4	phylo	as(x, "phylo")	phylobase
matching	phylo	as.phylo(x)	ape
treeshape	phylo	as.phylo(x)	apTreeshape
splits	phylo	as.phylo(x)	phangorn

Conversions among data types II

From	To	Command	Package
splits	networx	as.networx(x)	phangorn
evonet	phylo	as.phylo(x)	ape
evonet	networx	as.networx(x)	ape
evonet	network	as.network(x)	ape
evonet	igraph	as.igraph(x)	ape
haploNet	network	as.network(x)	pegas
haploNet	igraph	as.igraph(x)	pegas
hclust	phylo	as.phylo(x)	ape
hclust	dendrogram	as.dendrogram(x)	stats
DNABin	character	as.character(x)	ape
DNABin	alignment	as.alignment(x)	ape
DNABin	phyDat	as.phyDat(x)	phangorn
DNABin	genind	DNABin2genind(x)	adegenet
character	DNABin	as.DNABin(x)	ape

Conversions among data types III

From	To	Command	Package
character	loci	as.loci(x)	pegas
alignment	DNAbin	as.DNAbin(x)	ape
alignment	phyDat	as.phyDat(x)	phangorn
alignment	character	as.matrix(x)	seqinr
alignment	genind	alignment2genind(x)	adegenet
phyDat	DNAbin	as.DNAbin(x)	phangorn
phyDat	character	as.character(x)	phangorn
loci	genind	loci2genind(x)	pegas
loci	data frame	class(x) <- "data.frame"	
genind	loci	genind2loci(x)	pegas
data frame	phyDat	as.phyDat(x)	phangorn
data frame	loci	as.loci(x)	pegas
data frame	genind	df2genind(x)	adegenet
matrix	phyDat	as.phyDat(x)	phangorn

Basic operations with data I

```
1 x <- c(1, 2, 3, 4, 5) # Creates vector (see also ?rep)
2 c() # Is generic function to concatenate objects into new one
3 length(x) # Length of the object - for matrices and DF use dim()
4 str(x) # Information about structure of the object, try also ls.str()
5 mode(x) # Gets type of storage mode of the object
6 class(x) # Shows class of the object
7 x[2] # Shows second element of the object
8 x <- x[-5] # Removes fifth element
9 y <- matrix(data=5:20, nrow=4, ncol=4) # Creates a matrix
10 is.matrix(y) # Is it matrix? Try is.<TAB><TAB>
11 # TAB key shows available functions and objects starting by typed text
12 y # Prints the matrix
13 y[,2] # Prints second column
14 y[3,] # Prints third row
15 y[4,3] # Prints element from fourth row and third column
16 x <- y[2,] # Replaces "x" by second row of "y" (no warning)
17 # R doesn't ask neither notifies when overwriting objects! Be careful!
```

Basic operations with data II

```
1 rm(x) # Deletes x
2 y[,1:3] # Prints first through third column of the matrix
3 y[3,] <- rep(x=20, each=4) # Replaces third line by value of 20
4 y[y==20] <- 10 # If value of y's element is 20, replace it by 10
5 summary(y) # Basic statistics - according to columns
6 colnames(y) <- c("A", "B", "C", "D") # Set column names
7 # Objects and functions are without quotation marks; files, text with
8 colnames(y) # Prints column names, use rownames() in same way
9 y[, "C"] # Prints column C (R is case sensitive!)
10 t(y) # Transposes the matrix
11 y <- as.data.frame(y) # Turns into DF (see other functions as.*)
12 y[y==17] <- "NA" # Removes values of 17
13 y$B # Gets variable B of data frame y ($ works similarly in S3 objects)
14 save(list=ls(), file="test.RData") # Saves all objects during the work
15 load("test.RData") # Loads saved R environment with all objects
16 # When loading saved project, you have to load again libraries and
17 # scripts (see further), data objects are restored
```

Set repositories

```
1 # Basic package installation
2 install.packages("PackageName") # Case sensitive!
3 ?install.packages # Shows all available parameters
4 # We will need extra repositories:
5 options(repos=c("http://cran.at.r-project.org",
6 "http://r-forge.r-project.org/", "http://www.rforge.net/",
7 "http://bioconductor.statistik.tu-dortmund.de/packages/3.1/bioc",
8 "http://bioconductor.statistik.tu-dortmund.de/packages/3.1/data/
9 annotation", "http://bioconductor.statistik.tu-dortmund.de/
10 packages/3.1/data/experiment",
11 "http://bioconductor.statistik.tu-dortmund.de/packages/3.1/extra"))
12getOption("repos") # Shows actual repositories
13options() # Generic function to modify various settings
14?options # Gives details
```

- All the time keep newest version of R and and newest versions of packages!
- Installation of multiple packages may sometimes fail — install then packages in smaller groups or one by one

Install needed packages

```
1 install.packages(pkgs=c("Geneland", "PBSmapping", "ParallelStructure",
2 "RandomFields", "RgoogleMaps", "TeachingDemos", "XML", "ade4",
3 "adegenet", "adephylo", "akima", "ape", "colorspace", "combinat",
4 "corrplot", "fields", "gplots", "grid", "ips", "lattice", "mapdata",
5 "mapproj", "maps", "maptools", "muscle", "pegas", "permute",
6 "phangorn", "phyloch", "phytools", "polysat", "poppr", "rworldmap",
7 "seqinr", "sp", "spam", "tcltk", "vegan"), repos=getOption("repos"),
8 dependencies=TRUE)
9 library(package) # Loads installed package (we will do it on the fly)
10 update.packages(repos=getOption("repos")) # Updates installed packages
```

Installation of packages in GUI

- RStudio: set repositories by command from slide 28 and in bottom right pane select “Packages” and click on “Install Packages”...
- RKWard: go to menu “Settings” | “Configure ‘RKWard’” and select “Packages R”. Add URLs pf repositories from slide 28. OK. Go to menu “Settings” | “Manage R packages”, click to “Install...”, select and install desired packages...

Bioconductor

- Tools for analysis of genomic data, see <http://bioconductor.org/>
- To install it, add appropriate repositories (repository use to have same version number as R — change them accordingly when upgrading R) or use Bioconductor's special function (recommended by Bioconductor):

```
1 # Load basic Bioconductor function
2 source("http://bioconductor.org/biocLite.R")
3 # Get help how to use it
4 ?biocLite
5 # Install packages
6 biocLite(c("package1", "package2", "..."))
7 # Upgrades installed Bioconductor packages
8 biocLite("BiocUpgrade")
```

- Explore available packages:

<http://bioconductor.org/packages/release/BiocViews.html>

Population genetics and phylogenetics in R

Microsatellites, AFLP, SNP & sequences

- Now we will use mainly packages `adegenet` and `poppr` for dominant (presence/absence) data
 - Other important genetic packages: `ape`, `ade4` and `pegas`
 - Dominant/co-dominant marker data of any ploidy level including SSRs, SNP, and AFLP
 - Most of methods are able for polyploids (although not all)
 - Some methods are unavailable
- Mixing of ploidy levels is tricky — it doesn't matter when data are encoded as PA, otherwise it is problematic

```
1 library(ape)
2 library(ade4)
3 library(adegenet)
4 library(pegas)
5 library(poppr)
```

Load data

Source data:

```
pop msta93 msta101 msta102 msta103 msta105 msta131...
H01 He 269/269 198/198 221/223 419/419 197/197 196/196...
H02 He 275/283 198/198 221/223 419/419 193/193 168/190...
... ... ... ... ... ... ... ...
```

```
1 # Load training data (Taraxacum haussknechtii from Macedonia)
2 hauss.loci <- read.loci(file="http://trapa.cz/sites/default/rcourse/
3   haussknechtii_ssrs.txt", header=TRUE, loci.sep="\t", allele.sep="/",
4   col.pop=2, col.loci=3:14, row.names=1) # \t means TAB key
5 # Data control
6 hauss.loci
7 print(hauss.loci, details=TRUE)
```

First line starts with empty cell (if **header** is presented), there can be any extra column, just take care about **col.loci**. **row.names** are individual names (first column). Take care about **loci.sep** (here TAB) and **allele.sep** according to data formatting.

Prepare genind object for analysis and load coordinates

```

1 # Conversion of loci to genind - used for many analysis
2 hauss.genind <- loci2genind(hauss.loci)
3 # See population names
4 pop(hauss.genind)
5 hauss.genind$pop.names # "$" separates extra slots within object

```

Coordinates can be in any projection or scale — according to aim

Source data:

Ind	lon	lat
H01	21.3333	41.1
...

```

1 # Read coordinates
2 hauss.coord <- read.csv("http://trapa.cz/sites/default/rcourse/
3   haussknechtii_coordinates.csv", header=TRUE, sep="\t", quote="",
4   dec=". ", row.names=1)
5 hauss.coord

```

Take care about parameters of `read.csv()`! See `?read.csv` for details.

Add coordinates to genind and create genpop object

```
1 # Add coordinates - note identification of slots within object
2 hauss.genind$other$xy <- hauss.coord
3 # See result
4 hauss.genind$other$xy
5 hauss.genind
6 # Removes missing data - see ?missingno for types of dealing them
7 hauss.genind <- missingno(pop=hauss.genind, type="mean", cutoff=0.1)
8 # Convert corrected genind to loci
9 hauss.loci.cor <- genind2loci(hauss.genind)
10 # Writes loci file to the disk
11 write.loci(hauss.loci.cor, file="hauss.loci.cor.txt",
12   loci.sep="\t", allele.sep="/")
13 # Conversion to genpop - for population-level analysis
14 hauss.genpop <- genind2genpop(hauss.genind, missing="NA",
15   process.other=TRUE)
16 # See result
17 hauss.genpop
```

Import existing data set from popular software

```
1 read.genalex() # poppr - reads *.csv file  
2 read.fstat() # adegenet - reads *.dat files, only haploid/diploid data  
3 read.genetix() # adegenet - reads *.gtx files, only haploid/diploid data  
4 read.genepop() # adegenet - reads *.gen files, only haploid/diploid data  
5 read.structure() # adegenet - reads *.str files, only haploids/diploids  
6 import2genind() # adegenet - more automated version of above functions
```

One function rules them all...

All those functions (including read-loci and read.csv) are only modifications of **read.table()**. You can use it directly to import any data. Look at **?read.table** and play with it. Take care about parameters. Does the table use quotes to mark cell (quote=...)? How are columns separated (sep=...)? Is there a header with names of populations/loci/whatever (header=T/F)? What is decimal separator (dec=...)? Are there row names (used typically as names of individuals; row.names=...)? Check data after import!

Import of polyploid microsatellites

- adegent, poppr and related packages can for most of functions handle any ploidy level (but no mixing of various ploidy levels for microsatellites)
- polysat can handle mixed ploidy levels for microsatellites, but range of methods is limited

As for AFLP, we need two files: the data matrix and individual's populations:

	msat58	msat31	msat78	msat61	...
ala1	124/124/124	237/237/237	164/164/172	136/136/138	...
ala2	124/124/124	237/237/237	164/164/172	136/136/138	...
ala4-1	124/124/124	237/237/237	164/164/172	136/136/138	...
...

How to import polyploid microsatellites

```
1 # Import if table is as usual. Last column contains populations
2 tarax3n.table <- read.table("http://trapa.cz/sites/default/rcourse/
3   tarax3n.txt", header=TRUE, sep="\t", quote="", row.names=1)
4 # Check the data
5 tarax3n.table
6 class(tarax3n.table)
7 dim(tarax3n.table)
8 # See parameter "X" - we don't import whole tarax3n.table as last column
9 # contains populations - this column we use for "pop" parameter (note
10 # different style of calling the column - just to show the possibility).
11 # Check "ploidy" and "ncode" (how many digits code one allele - must be
12 # same everywhere). See ?df2genind for more details.
13 tarax3n.genind <- df2genind(X=tarax3n.table[,1:6], sep="/", ncode=3,
14   pop=tarax3n.table[["pop"]], missing=NA, ploidy=3, type="codom")
15 # See resulting genind object
16 tarax3n.genind
17 summary(tarax3n.genind)
```

Import of AFLP data — background

Source data (two files - AFLP data with individual names and populations)

AFLP – presence/absence data:

L1	L2	L3	L4	L5	L6	L7	L8	L9	...	
Ind1	0	0	1	1	1	0	0	0	1	...
IndG	0	0	1	1	0	0	0	0	0	...
...	

Individual's populations:

POP
pop1
popZ
...

- Use any names, just keep one word (no spaces) and don't use special characters
- Keep names of loci as simple as possible, there are some issues when they contain dots
- As soon as one line of data = one individual, ploidies and their mixing doesn't matter
- Not all methods are available/meaningful for PA

Import of AFLP data — the code

```
1 amara.aflp <- read.table(file="http://trapa.cz/sites/default/rcourse/
2   amara_aflp.txt", header=TRUE, sep="\t", quote="")
3 amara.aflp
4 dim(amara.aflp)
5 class(amara.aflp) # Must be matrix or data frame
6 # Populations - just one column with population names for all inds
7 amara.pop <- read.table(file="http://trapa.cz/sites/default/rcourse/
8   amara_pop.txt", header=TRUE, sep="\t", quote="")
9 amara.pop
10 # You can use just one file, where populations are in last column and
11 # in df2genind() use for example X=aflp[,1:XXX] and pop=aflp[,YYY]
12 # Create genind object - ind.names and loc.names are taken from X
13 aflp.genind <- df2genind(X=amara.aflp, sep=NULL, ind.names=NULL,
14   loc.names=NULL, pop=amara.pop[,1], missing=NA, type="PA")
15 aflp.genind
16 summary(amara.aflp)
17 # You can add any other variables into genind$other$XXX
```

Another data manipulation

SNPs can be imported into genind in same way as AFLP

```
1 genind2df() # adegenet - export into data frame
2 genind2genalex() # poppr - export for genalex
3 splitcombine() # poppr - edits population hierarchy
4 popsub() # poppr - extracts only selected population(s)
5 clonecorrect() # poppr - corrects for clones
6 informloci() # poppr - removes uninformative loci
7 seppop() # adegenet - separates populations from genind or genlight
8 seploc() # adegenet - splits genind, genpop or genlight by markers
9 # seppop and seploc return lists of objects - for further analysis
10 # read manual pages (...) of those functions before usage
```

Notes about getting data into R

- When importing fragmentation data, we somehow use function **read.table()** — it is important to understand it
- I recommend to use TAB (TSV — tab separated values; encoded as “\t” in R) to separate columns (no quotation marks, no commas)
- When importing microsatellites, all alleles **must** have same number of digits. Separate alleles by “/”, “|” or something similar and correctly specify it in **read.loci()** or **df2genind()**
- Do not use underscores (“_”) to name objects in R — only numbers, letters, dots and minuses
- **read.loci()** sometimes doesn't work correctly on AFLP or polyploid microsatellites — try **write.table()** instead...
- Theoretically, genind object could be currently able to store mixing ploidy data, but most of analysis will probably fail...

Import of DNA data

```
1 # Reading FASTA (reads also another formats, see ?read.dna)
2 # Sequences of flu viruses from various years
3 usflu.dna <- read.dna(file="http://adegenet.r-forge.r-project.org/
4   files/usflu.fasta", format="fasta")
5 # Another possibility (only for FASTA alignments, same result):
6 usflu.dna <- fasta2DNAbin(file="http://adegenet.r-forge.r-project.org/
7   files/usflu.fasta") # Normally keeps only SNP -- see ?fasta2DNAbin
8 # Check the object
9 usflu.dna
10 class(usflu.dna)
11 as.character(usflu.dna)[1:5,1:10]
12 # Read annotations
13 usflu.annot <- read.csv("http://adegenet.r-forge.r-project.org/files/
14   usflu.annot.csv", header=TRUE, row.names=1)
15 head(usflu.annot) # See result
16 # Convert DNAbin to genind - only polymorphic loci (SNPs) are retained
17 # When converting DNAbin to genind, the sequences must be aligned!
18 usflu.genind <- DNAbin2genind(x=usflu.dna, pop=usflu.annot[["year"]])
```

Import sequences from GeneBank

```
1 # read.fasta() from seqinr package reads DNA or AA in FASTA format
2 # returns a list (DNAbin is for us now better choice)
3 usflu.dna <- read.fasta(file="http://adegenet.r-forge.r-project.org/
4   files/usflu.fasta", seqtype="DNA")
5 # Importing sequences according to sequence ID
6 # data from http://www.ncbi.nlm.nih.gov/popset/608602125 (Meles meles)
7 meles.dna <- read.GenBank(c("KJ161355.1", "KJ161354.1", "KJ161353.1",
8   "KJ161352.1", "KJ161351.1", "KJ161350.1", "KJ161349.1", "KJ161348.1",
9   "KJ161347.1", "KJ161346.1", "KJ161345.1", "KJ161344.1", "KJ161343.1",
10  "KJ161342.1", "KJ161341.1", "KJ161340.1", "KJ161339.1", "KJ161338.1",
11  "KJ161337.1", "KJ161336.1", "KJ161335.1", "KJ161334.1", "KJ161333.1",
12  "KJ161332.1", "KJ161331.1", "KJ161330.1", "KJ161329.1", "KJ161328.1"))
13 meles.dna
14 class(meles.dna)
15 # Converts DNAbin to genind - extracts SNP - for large datasets can be
16 # computationally very intensive
17 meles.genind <- DNAbin2genind(meles.dna)
```

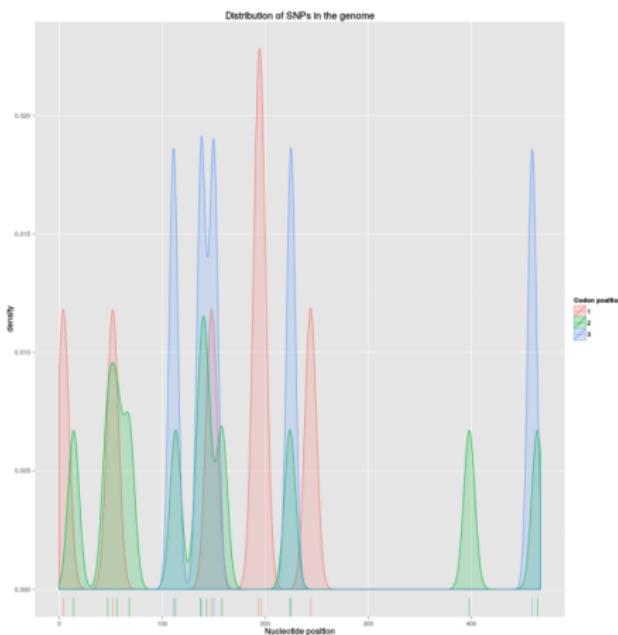
To query on-line database as through web interface we use [seqinr](#).

Query on-line sequence databases

```
1 library(seqinr)
2 choosebank() # Genetic banks available for seqinr
3 choosebank("embl") # Choose some bank
4 ?query # See how to construct the query
5 # Query selected database - there are much possibilities
6 nothofagus <- query(listname="nothofagus",
7   query="SP=Nothofagus AND K=rbcl", verbose=TRUE)
8 # See the sequences information
9 nothofagus$req
10 # Get the sequences as a list
11 nothofagus.sequences <- getSequence(nothofagus$req)
12 # See sequences
13 nothofagus.sequences
14 # Get annotations
15 nothofagus.annot <- getAnnot(nothofagus[["req"]])
16 # Close the bank when work is over
17 closebank()
18 # Convert sequences from a list to DNAbin (functions as.DNAbin*)
19 nothofagus.dna <- as.DNAbin.list(nothofagus.sequences)
```

Checking SNPs

```
1 # Position of polymorphism within
2 # alignment - snpposi.plot requires
3 # input data in form of matrix
4 snpposi.plot(x=as.matrix(
5   meles.dna), codon=FALSE)
6 # Position of polymorphism within
7 # alignment - differentiating codons
8 snpposi.plot(as.matrix(meles.dna))
9 # When converting to genind object,
10 # only polymorphic loci are kept -
11 # threshold for polymorphism can
12 # be arbitrary
13 meles.genind <- DNAbin2genind(x=
14   meles.dna, polyThres=0.01)
15 # Test of random distribution of SNP
16 snpposi.test(as.matrix(meles.dna))
```



Importing SNP

Import from **PLINK** requires saving of data with option “**-recodeA**”:

```
read.PLINK(file="PLINKfile", ...) # See ?readPLINK
```

Extracting SNP from alignments reads FASTA alignments and keep only SNPs. The method is relatively efficient even for large data sets with several genomes:

```
1 usflu.genlight <- fasta2genlight(file=
2   "http://adegenet.r-forge.r-project.org/files/usflu.fasta",
3   quiet=FALSE, saveNbAlleles=TRUE)
4 ?fasta2genlight # Function has several options to speed up reading
5 # If it crashes (on Windows), try add parameter "parallel=FALSE"
```

- For small datasets it doesn't matter if resulting dataset is genind or genlight, for large genlight is more efficient

```
1 # Convert genlight to genind
2 usflu.genind <- as.genind(usflu.genlight)
3 # Resulting data matrix might be slightly different than when using
4 # DNAbin2genind - depends on used settings
```

Notes about using genlight (vs. genind)

- Genlight is “just” version of more common genind object to store large data sets with (nearly) complete multiple genomes
- “Large” is tricky — there is no easy criterion — try genind and when work fails because of not enough computer resources, go on with genlight
- Can be quickly converted to genind (not back)
- Use is basically same as when working with genind — but not all functions are able to deal with it (on the other hand, others are optimised to work well on large data sets)
- SNPbin is version of genind/genlight to store one large genome — serves basically as storage, no need to deal with it
- genlight allows varying ploidy level — genind has issues with it

Export data

```
1 # Convert genind into DF using genind2genotype() or genind2df()
2 hauss.df <- genind2genotype(x=hauss.genind, pop=hauss.genind@pop,
3   res.type="matrix")
4 hauss.df <- genind2df(x=hauss.genind, pop=NULL, sep="/",
5   usepop=TRUE, oneColPerAll=FALSE)
6 write.table(x=hauss.df, file="haussdata.txt", quote=FALSE,
7   sep="\t", na="NA", dec=". ", row.names=TRUE, col.names=TRUE)
8 # Export of DNA sequences into FASTA format
9 write.dna(x=usflu.dna, file="usflu.fasta", format="fasta",
10   append=FALSE, nbcol=6)
11 write.fasta(sequences=meles.dna, names=names(meles.dna),
12   file.out="meles.fasta", open="w")
13 # Export trees (objects of class phylo)
14 # Writes tree(s) in NEWICK format
15 write.tree(phy=hauss.nj.bruvo, file="haussknechtii.nwk")
16 # Writes tree(s) in NEXUS format
17 write.nexus(hauss.nj.bruvo, file="haussknechtii.nexus")
```

Descriptive statistics I

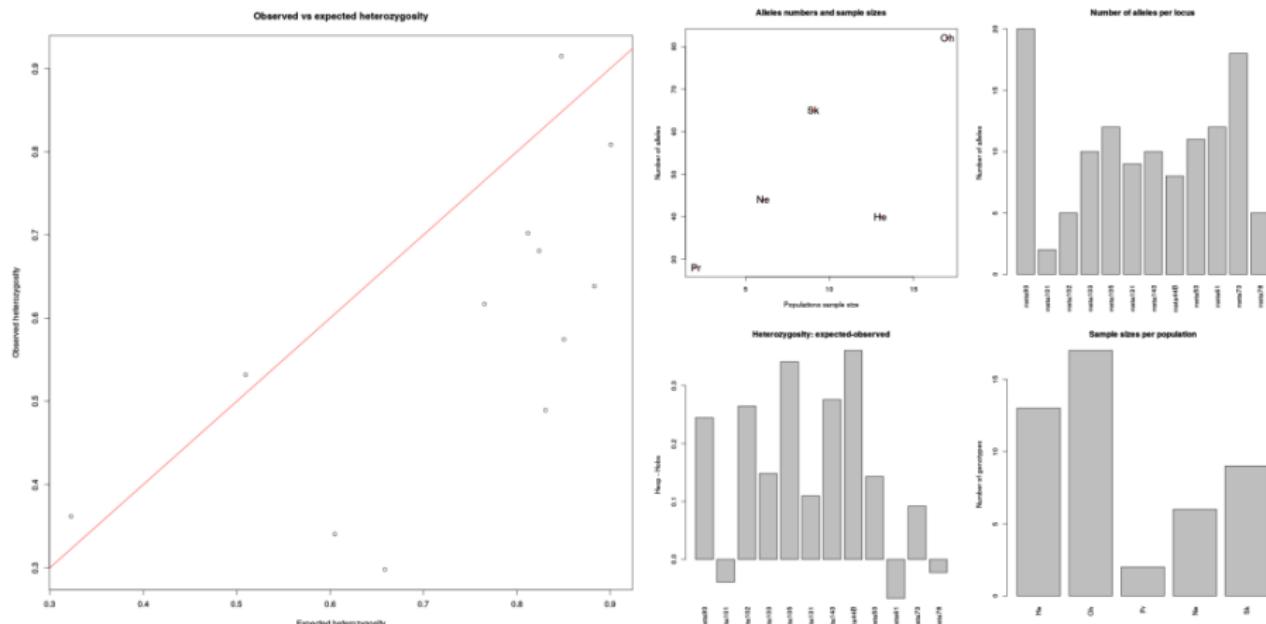
We will now work mainly with diploid SSRs for *Taraxacum haussknechtii*, you can try other data examples by yourselves

```
1 # Get summary - names and sizes of populations,
2 # heterozygosity, some info about loci
3 hauss.summ <- summary(hauss.genind)
4 # Plot expected vs. observed heterozygosity
5 # it looks like big difference
6 plot(x=hauss.summ$Hexp, y=hauss.summ$Hobs,
7       main="Observed vs expected heterozygosity",
8       xlab="Expected heterozygosity", ylab="Observed heterozygosity")
9 abline(0, 1, col="red")
10 # T-test of difference between
11 # observed and expected heterozygosity - strongly significant
12 t.test(x=hauss.summ$Hexp, y=hauss.summ$Hobs, paired=TRUE, var.equal=T)
13 # Bartlett's K-squared of difference
14 # between observed and expected heterozygosity - not significant
15 bartlett.test(list(hauss.summ$Hexp, hauss.summ$Hobs))
```

Descriptive statistics II

```
1 # Create pane with some information
2 par(mfrow=c(2,2)) # Divide graphical devices into 4 smaller spaces
3 # Plot alleles number vs. population sizes
4 plot(x=hauss.summ$pop.eff, y=hauss.summ$pop.nall, xlab="Populations
5     sample size", ylab="Number of alleles", main="Alleles numbers and
6     sample sizes", col="red", pch=20)
7 # Add text description to the point
8 text(x=hauss.summ$pop.eff, y=hauss.summ$pop.nall,
9     lab=names(hauss.summ$pop.eff), cex=1.5)
10 # Barplots of various data
11 barplot(height=hauss.summ$loc.nall, ylab="Number of alleles",
12     main="Number of alleles per locus",
13     names.arg=hauss.genind@loc.names, las=3)
14 barplot(height=hauss.summ$Hexp-hauss.summ$Hobs, main="Heterozygosity:
15     expected-observed", ylab="Hexp - Hobs",
16     names.arg=hauss.genind@loc.names, las=3)
17 barplot(height=hauss.summ[["pop.eff"]], main="Sample sizes per
18     population", ylab="Number of genotypes", las=3)
```

Graphs from previous two slides



```
1 dev.off() # Closes graphical device - otherwise following  
2 # graphs would still be divided into 4 parts
```

Population statistics by poppr

It returns various population statistics — **png()** will produce set of figures for each population. You can similarly use **svg()**, **pdf()**, ...

```
1 png(filename="poppr_%03d.png", width=720, height=720, bg="white")
2 # Output figures will be saved to the disk
3 poppr(dat=hauss.genind, total=TRUE, sample=1000, method=4,
4 missing="geno", cutoff=0.15, quiet=FALSE, clonecorrect=FALSE,
5 hist=TRUE, minsamp=1, legend=TRUE)
6 dev.off() # Closes graphical device - needed after use of functions
7 # like png(), svg(), pdf(), ... to write the file(s) to the disk
```

- Pop — Population name (Note that "Total" also means "Pooled")
- N — Number of individuals observed
- MLG — Number of multilocus genotypes (MLG) observed
- eMLG — The number of expected MLG at the smallest sample size 10 based on rarefaction
- SE — Standard error based on eMLG
- H — Shannon-Wiener Index of MLG diversity
- G — Stoddart and Taylor's Index of MLG diversity
- Hexp — Nei's 1978 genotypic diversity (corrected for sample size), or Expected Heterozygosity
- E.5 — Evenness, E5
- Ia — The index of association, IA
- rbarD — The standardized index of association
- See [poppr's manual](#) for details

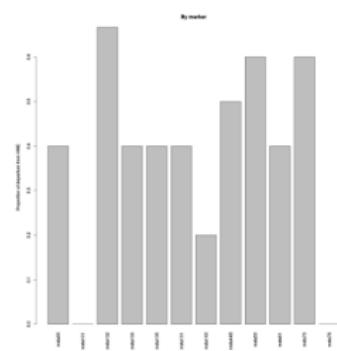
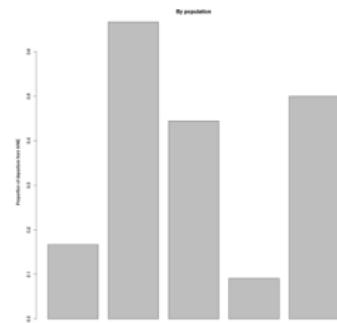
Departure from Hardy-Weinberg equilibrium

```
1 # Departure from HWE
2 hauss.hwe <- HWE.test.genind(x=hauss.genind, permut=TRUE,
3   nsim=1000, res.type="matrix")
4 # According to populations
5 barplot(apply(X=hauss.hwe<0.05, 1, FUN=mean, na.rm=TRUE),
6   las=3, ylab="Proportion of departure from HWE", main="By
7   population", names.arg=hauss.genind@pop.names)
8 # According to loci
9 barplot(apply(X=hauss.hwe<0.05, 2, FUN=mean, na.rm=TRUE),
10  las=3, ylab="Proportion of departure from HWE", main="By marker")
11 hauss.hwe.test <- hw.test(x=hauss.loci, B=1000)
12 hauss.hwe.test
13 # List of lists for all combinations of loci and populations
14 hauss.hwe.full <- HWE.test.genind(x=hauss.genind, permut=TRUE,
15   nsim=1000, res.type="full")
16 hauss.hwe.full # Very looong list...
```

Departure from HWE — results

```
> hauss.hwe.test
```

	chi^2	df	Pr(chi^2 >)	Pr.exact
msta93	383.5519728	190	3.552714e-15	0.000
msta101	0.6927242	1	4.052393e-01	0.657
msta102	83.0741964	10	1.250111e-13	0.000
msta103	77.1819098	45	1.998865e-03	0.000
msta105	282.3441356	66	0.000000e+00	0.000
msta131	175.8647805	36	0.000000e+00	0.000
msta143	155.4567196	45	4.474199e-14	0.000
msta44B	112.6283892	28	4.149125e-12	0.000
msta53	110.2581985	55	1.436375e-05	0.002
msta61	70.7719907	66	3.215148e-01	0.070
msta73	310.9811293	153	7.842615e-13	0.000
msta78	4.4688908	10	9.237260e-01	0.569



F-statistics and population theta

Functions return tables of Fst/theta values for populations/loci

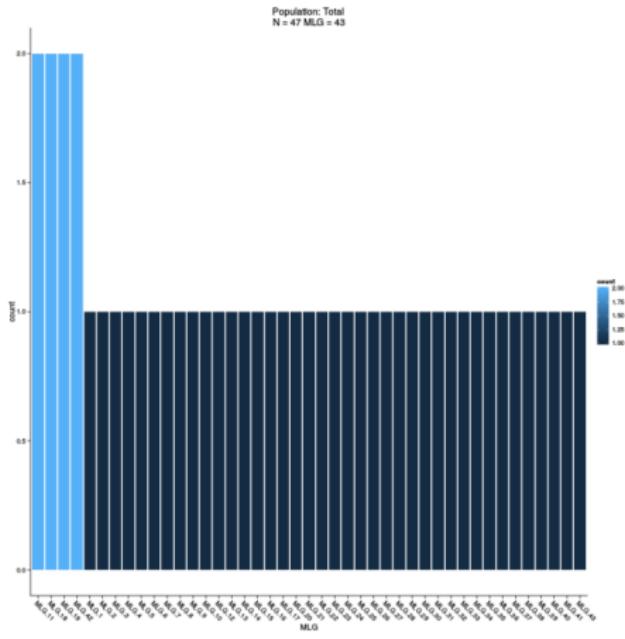
```
1 # Fit, Fst and Fis for each locus
2 Fst(x=hauss.loci.cor, pop=1)
3 # multilocus estimators of variance components and F-statistics,
4 # alternative to Fst
5 fstat(x=hauss.genind, pop=NULL, fstonly=FALSE)
6 # Nei's pairwise Fst between all pairs of populations.
7 # Difference in res.type="dist"/"matrix" is only in format of output
8 pairwise.fst(x=hauss.genind, pop=NULL, res.type="matrix", truenames=T)
9 # Estimates of population theta according to Kimmel et al. 1998
10 theta.msat(hauss.loci.cor)

Fst(x=hauss.loci.cor, pop=1)
```

	Fit	Fst	Fis
msta93	0.31835291	0.17867087	0.17006829
msta101	-0.09968472	0.04064928	-0.14628018
...

Multi locus genotypes

```
1 # Total number of MLGs
2 # (simple value)
3 mlg(pop=hauss.genind, quiet=FALSE)
4 # MLGs shared among populations
5 mlg.crosspop(pop=hauss.genind,
6   df=TRUE, quiet=FALSE)
7 # Detailed view on distribution
8 # of MLGs into populations
9 # (table and/or plot)
10 mlg.table(pop=hauss.genind,
11   bar=TRUE, total=TRUE,
12   quiet=FALSE)
13 # See manual pages for details
```



Functions from poppr package — the best for microsatellites, although available also from another data types

Basic distances

```
1 # Euclidean distance for individuals
2 hauss.dist <- dist(x=hauss.genind, method="euclidean", diag=T, upper=T)
3 hauss.dist
4 # Nei's distance (not Euclidean) for populations
5 # (other methods are available, see ?dist.genpop)
6 hauss.dist.pop <- dist.genpop(x=hauss.genpop, method=1, diag=T, upper=T)
7 # Test if it is Euclidean
8 is.euclid(hauss.dist.pop, plot=TRUE, print=TRUE, tol=1e-10) # FALSE = No
9 # Turn to be Euclidean
10 hauss.dist.pop <- cailliez(distmat=hauss.dist.pop, print=FALSE,
11 tol=1e-07, cor.zero=TRUE)
12 # Test if it is Euclidean
13 is.euclid(hauss.dist.pop, plot=TRUE, print=TRUE, tol=1e-10) # TRUE = OK
14 # Show it
15 hauss.dist.pop
```

Most of analysis based on distances more or less require Euclidean (non-negative) distances. If the distance matrix contains non-Euclidean distances, the result can be weird...

Distances reflecting microsatellite repeats

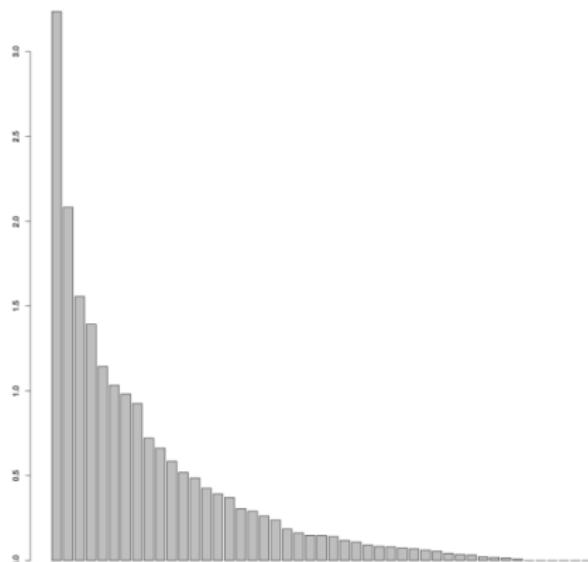
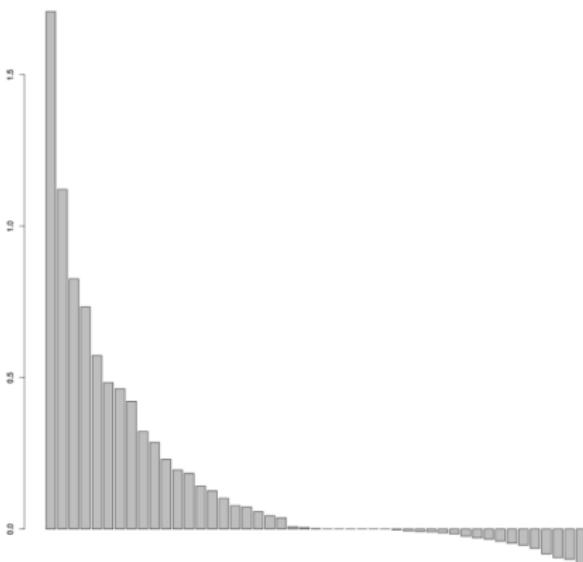
```
1 # Bruvo's distances weighting SSRs repeats - take care about replen
2 # parameter - requires repetition length for every SSRs locus
3 hauss.dist.bruvo <- bruvo.dist(pop=hauss.genind, replen=rep(2, 12),
4 loss=TRUE)
5 # Test if it is Euclidean
6 is.euclid(hauss.dist.bruvo, plot=TRUE, print=TRUE, tol=1e-10)
7 # Turn to be Euclidean
8 hauss.dist.bruvo <- cailliez(distmat=hauss.dist.bruvo, print=FALSE,
9 tol=1e-07, cor.zero=TRUE)
10 # Test if it is Euclidean
11 is.euclid(hauss.dist.bruvo, plot=TRUE, print=TRUE, tol=1e-10)
12 # Show it
13 hauss.dist.bruvo
```

See poppr's manual and manual pages of the functions for details and different possibilities of settings.

Turning distance matrix into Euclidean is controversial...

How to deal with zero distances in original matrix? There is no really good solution...

Histograms of Bruvo distance before and after transformation:



More distances...

```

1 # Nei's distance (not Euclidean) for individuals
2 # (other methods are available, see ?nei.dist from poppr package)
3 hauss.dist.nei <- nei.dist(x=hauss.genind, warning=TRUE)
4 hauss.dist.nei
5 # Dissimilarity matrix returns a distance reflecting the number of
6 # allelic differences between two individuals
7 hauss.dist.diss <- diss.dist(x=hauss.genind, percent=FALSE, mat=TRUE)
8 hauss.dist.diss

```

Import own distance matrix from another software:

	Fe	He	Oh	...
Fe	0.00000	132.019	109.159	...
He	132.0191	0.00000	9.89111	...
Oh	109.1590	9.89111	0.00000	...
Pr	139.5669	8.55312	4.40562	...
Ne	156.7619	9.96143	16.6927	...
...

```

1 MyDistance <- read.csv("distances.
2   txt", header=TRUE, sep="\t",
3   dec=". ", row.names=1)
4 MyDistance <- as.dist(MyDistance)
5 class(MyDistance)
6 dim(MyDistance)
7 MyDistance

```

Distances among DNA sequences

The sequences must be aligned before calculating distances among them!

```
1 ?dist.dna # There are various models available
2 usflu.dist <- dist.dna(x=usflu.dna, model="TN93")
3 usflu.dist
4 class(usflu.dist)
5 dim(as.matrix(usflu.dist))
6 meles.dist <- dist.dna(x=meles.dna, model="F81")
7 meles.dist
8 class(meles.dist)
9 dim(as.matrix(meles.dist))
```

Distances and genlight object

Pairwise genetic distances for each data block (genlight objects with whole genome data) — sensitive to missing data (not useful in every case):

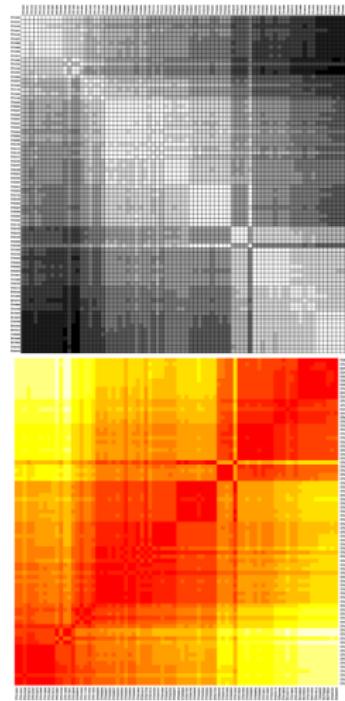
```
1 usflu.dists <- lapply(X=usflu.genlight, FUN=function(DDD)
2   dist(as.matrix(DDD)))
3 class(usflu.dists)
4 names(usflu.dists)
5 class(usflu.dists[[1]])
6 usflu.distr <- Reduce(f="+", x=usflu.dists)
7 class(usflu.distr)
8 usflu.distr
9 # It is possible to use just basic dist function on whole
10 # genlight object
11 usflu.distg <- dist(as.matrix(usflu.genlight))
```

Rationale of this approach is to save resources when dividing whole data set into smaller blocks — useful for huge data, nor for most of the cases

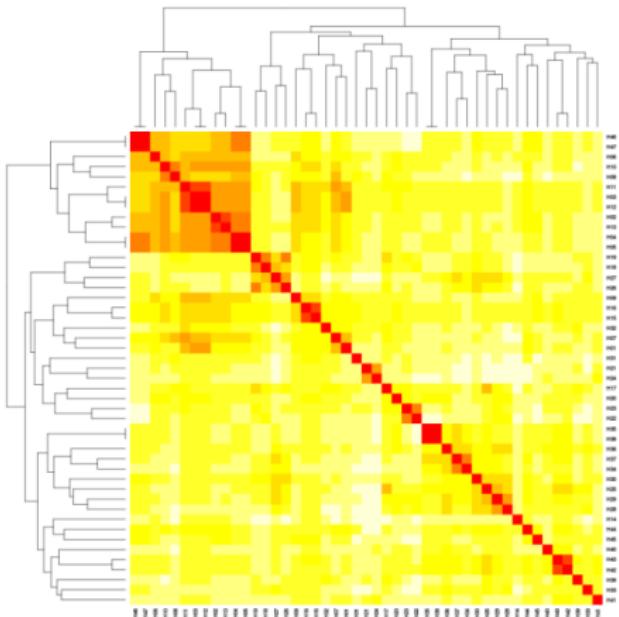
Visualize pairwise genetic similarities

```
1 # table.paint() requires data
2 # frame, dist can't be directly
3 # converted to DF
4 table.paint(df=as.data.frame(
5   as.matrix(usflu.dist)), cleg=0,
6   clabel.row=0.5, clabel.col=0.5)
7 # Same visualization, coloured
8 # heatmap() reorders values
9 # because by default it plots
10 # also dendrograms on the edges
11 heatmap(x=as.matrix(usflu.dist),
12   Rowv=NA, Colv=NA, symm=TRUE)
```

Another possibility is to use function **corrplot()** from package **corrplot** — it plots correlation plots.



Heatmaps



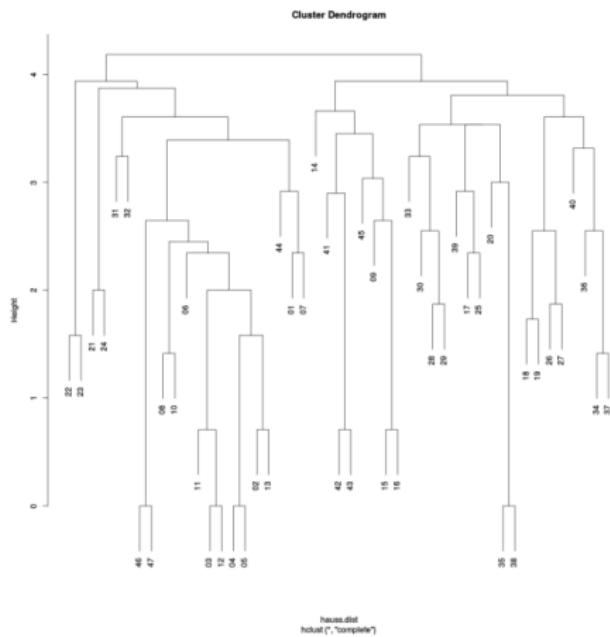
```
1 # Based on various distances
2 heatmap(as.matrix(hauss.dist),
3         symm=TRUE, labRow=rownames(
4             as.matrix(hauss.dist.bruvo)),
5             labCol=colnames(as.matrix(
6                 hauss.dist.bruvo)))
7 # hauss.dist doesn't contain
8 # names of individuals - add here
9 heatmap(as.matrix(hauss.dist.pop),
10        symm=TRUE)
11 heatmap(as.matrix(hauss.dist.
12            bruvo), symm=TRUE)
13 heatmap(as.matrix(hauss.dist.
14            diss), symm=TRUE)
```

There are various settings — colours, dendrogram, ... See **?heatmap**.

Hierarchical clustering — UPGMA and others

```
1 # According to distance used
2 # see ?hclust for methods
3 plot(hclust(d=hauss.dist,
4   method="complete"))
5 plot(hclust(d=hauss.dist.pop,
6   method="complete"))
7 plot(hclust(d=hauss.dist.
8   bruvo, method="complete"))
```

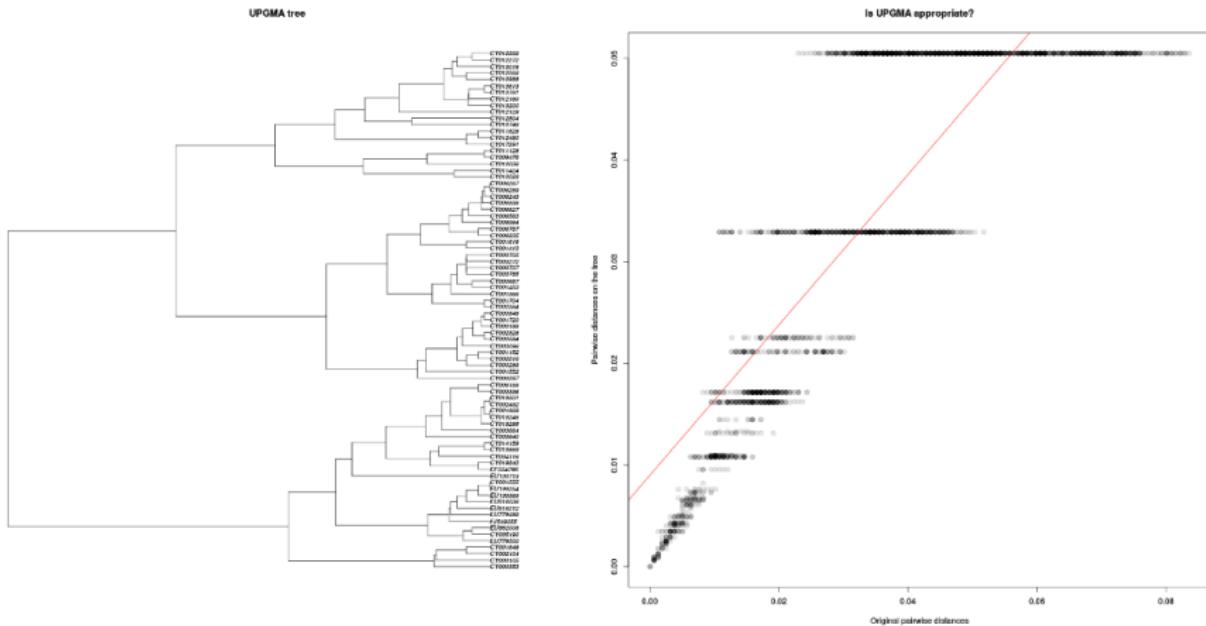
This is very basic function to make dendrogram. There are better possibilities (NJ etc).



UPGMA and its test

```
1 # Calculate it
2 # Saving as phylo object (and not hclust) gives more
3 # possibilities for further plotting and manipulations
4 usflu.upgma <- as.phylo(hclust(d=usflu.dist, method="average"))
5 plot.phylo(x=usflu.upgma, cex=0.75)
6 title("UPGMA tree")
7 # Test quality - tests correlation of original distance in the matrix
8 # and reconstructed distance from hclust object
9 plot(x=as.vector(usflu.dist), y=as.vector(as.dist(
10   cophenetic(usflu.upgma))), xlab="Original pairwise distances",
11   ylab="Pairwise distances on the tree", main="Is UPGMA
12   appropriate?", pch=20, col=transp(col="black",
13   alpha=0.1), cex=2)
14 abline(lm(as.vector(as.dist(cophenetic(usflu.upgma))) ~
15   as.vector(usflu.dist)), col="red")
```

UPGMA is not the best choice here...



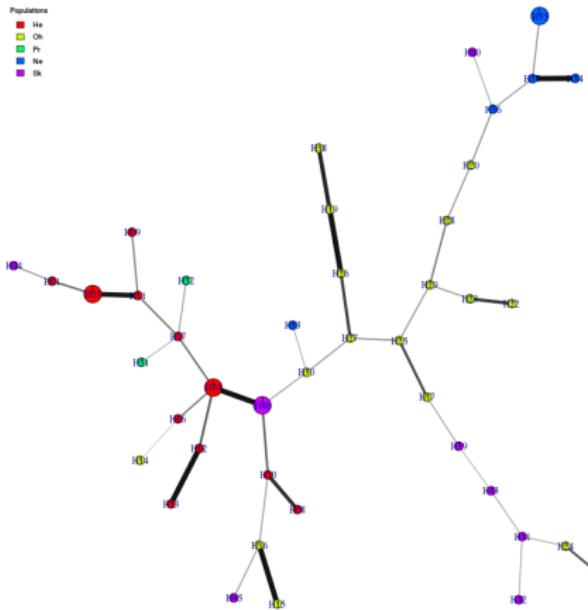
AMOVA

```
1 # From package pegas (doesn't directly show percentage of variance)
2 hauss.pop <- hauss.genind$pop
3 hauss.amova <- amova(hauss.dist~hauss.pop, data=NULL,
4   nperm=1000, is.squared=TRUE)
5 hauss.amova
6 # Output:
7 $tab
8           SSD      MSD df
9 hauss.pop 15.11916 3.779790 4
10 Error     59.62501 1.419643 42
11 Total     74.74417 1.624873 46
12 # Another possibility is poppr.amova - for more complicated hierarchy
13 # see ?poppr.amova
```

Minimum Spanning Network

Populations

- He
- Ob
- Pr
- Ne
- Sk

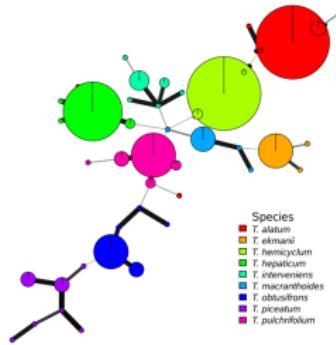


```

1 bruvo.msn(pop=hauss.genind,
2 replen=rep(2, 12), loss=TRUE,
3 palette=rainbow, vertex.label=
4 "inds", gscale=TRUE,
5 wscale=TRUE, showplot=TRUE)

```

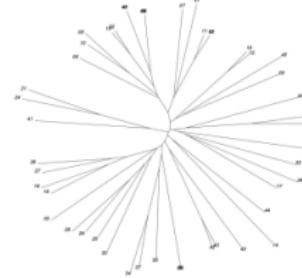
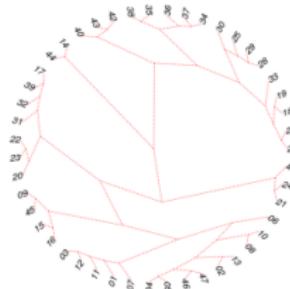
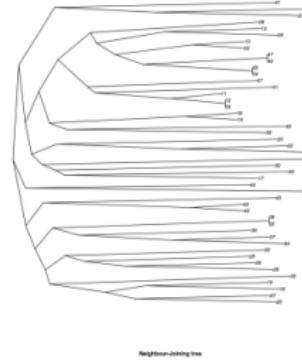
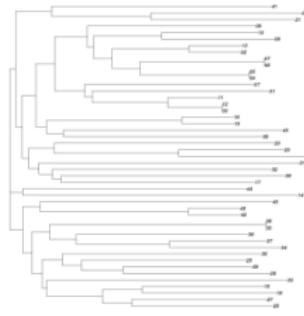
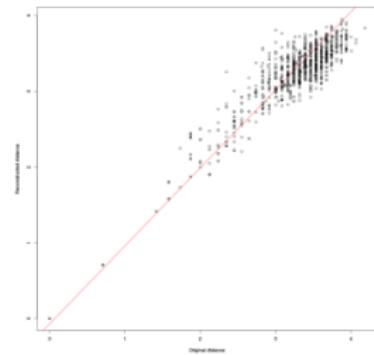
Package poppr, based on Bruvo's distance (for SSRs).



Calculate and test NJ tree

```
1 # Calculates the tree (try with various distances)
2 hauss.nj <- nj(hauss.dist)
3 # Test tree quality - plot original vs. reconstructed distance
4 plot(as.vector(hauss.dist), as.vector(as.dist(cophenetic(hauss.nj))),
5      xlab="Original distance", ylab="Reconstructed distance")
6 abline(lm(as.vector(hauss.dist) ~
7          as.vector(as.dist(cophenetic(hauss.nj)))), col="red")
8 summary(lm(as.vector(hauss.dist) ~
9          as.vector(as.dist(cophenetic(hauss.nj))))) # Prints summary text
10 # Bootstrap
11 hauss.boot <- boot.phylo(phy=hauss.nj, x=hauss.loci.cor,
12   FUN=function(XXX) nj(dist(loci2genind(XXX))), B=1000)
13 # Plot a basic tree - see ?plot.phylo for details
14 plot.phylo(x=hauss.nj, type="phylogram")
15 plot.phylo(x=hauss.nj, type="cladogram", edge.width=2)
16 plot.phylo(x=hauss.nj, type="fan", edge.width=2, edge.lty=2)
17 plot.phylo(x=hauss.nj, type="radial", edge.color="red",
18   edge.width=2, edge.lty=3, cex=2)
```

Choose your tree...

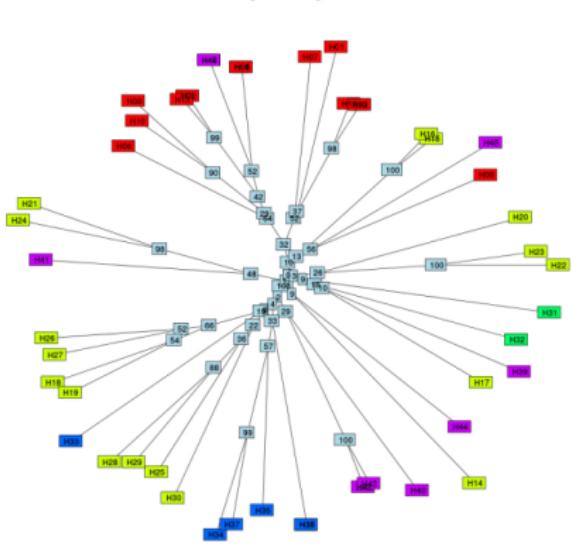


Nicer trees

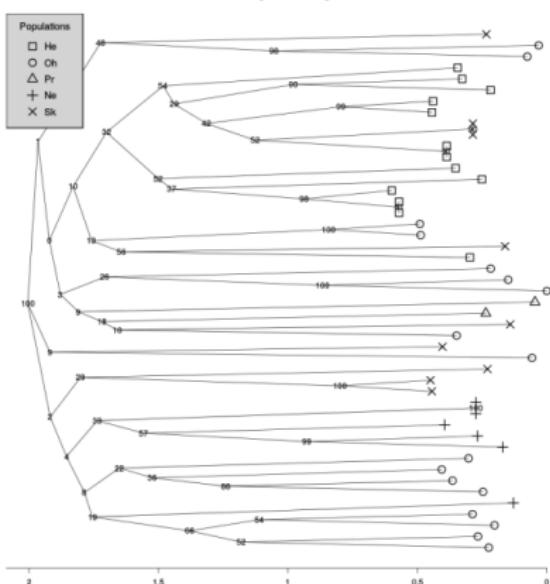
```
1 ## Plot a nice tree with colored tips
2 plot.phylo(x=hauss.nj, type="unrooted", show.tip=F, lwd=3, main="NJ")
3 # Labels for nodes - bootstrap - see ?nodelabels for graphical settings
4 nodelabels(text=round(hauss.boot/10))
5 # Colored labels - creates vector of colors according to populations
6 nj.rainbow<-colorRampPalette(rainbow(length(levels(pop(hauss.genind)))))
7 tiplabels(text=hauss.genind$ind.names, bg=fac2col(x=hauss.genind$pop,
8   col.pal=nj.rainbow)) # Coloured tips
9 ## Plot BW tree with tip symbols and legend
10 plot.phylo(x=hauss.nj, type="cladogram", show.tip=F, lwd=3, main="NJ")
11 axisPhylo() # Add axis with distances
12 # From node labels let's remove unneeded frame
13 nodelabels(text=round(hauss.boot/10), frame="none", bg="white")
14 # As tip label we use only symbols - see ?points for graphical details
15 tiplabels(frame="none", pch=rep(0:4,times=c(13,17,2,6,9)), lwd=2, cex=2)
16 # Plot a legend explainindg symbols
17 legend(x="topleft", legend=c("He", "Oh", "Pr", "Ne", "Sk"),
18   border="black", pch=0:4, pt.lwd=2, pt.cex=2, bty="o", bg="lightgrey",
19   box.lwd=2, cex=1.2, title="Populations")
```

Choose your tree...

Neighbour-Joining tree



Neighbour-Joining tree



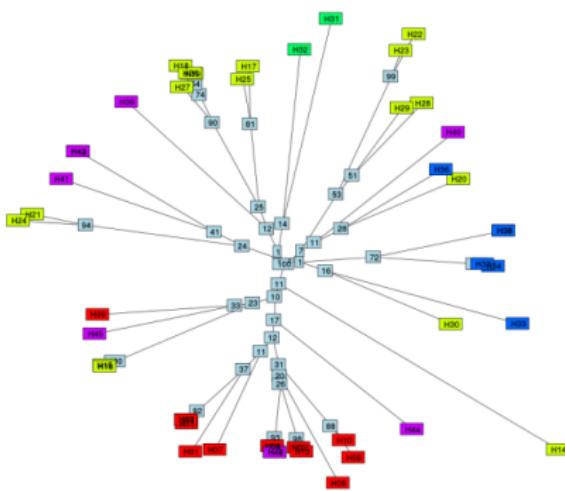
Trees based on Bruvo's distance

Package poppr (bootstrap is incorporated within the function)

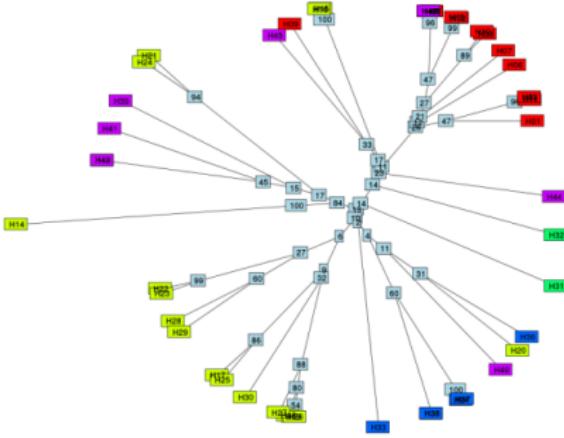
```
1 # NJ
2 hauss.nj.bruvo <- bruvo.boot(pop=hauss.genind, replen=rep(2, 12),
3   sample=1000, tree="nj", showtree=TRUE, cutoff=1, quiet=FALSE)
4 plot.phylo(x=hauss.nj.bruvo, type="unrooted", show.tip=FALSE,
5   lwd=3, main="Neighbor-Joining tree.")
6 # you can call node labels as phylo$node.labels or phylo[["node.labels"]]
7 nodelabels(hauss.nj.bruvo[["node.labels"]])
8 tiplabels(hauss.nj.bruvo[["tip.label"]], bg=fac2col(x=hauss.genind$pop,
9   col.pal=nj.rainbow))
10 # UPGMA
11 hauss.upgma <- bruvo.boot(pop=hauss.genind, replen=rep(2, 12),
12   sample=1000, tree="upgma", showtree=TRUE, cutoff=1, quiet=FALSE)
13 plot.phylo(hauss.upgma, type="unrooted", show.tip=FALSE, lwd=3,
14   main="UPGMA tree")
15 nodelabels(hauss.upgma[["node.labels"]])
16 tiplabels(hauss.upgma[["tip.label"]], bg=fac2col(x=hauss.genind@pop,
17   col.pal=nj.rainbow))
```

Choose your tree...

Neighbor-Joining tree.



UPGMA tree



NJ tree of populations

```
1 # Populations
2 hauss.nj.pop <- nj(hauss.dist.pop)
3 # Bootstrap - source() loads external scripts
4 # boot.phylo doesn't work for population trees
5 source("http://trapa.cz/sites/default/rcourse/boot_phylo_nj_pop.r")
6 hauss.boot.pop <- boot.phylo.nj.pop(hauss.nj.pop, hauss.genind, 1000)
7 # Plot a tree
8 plot(hauss.nj.pop, type="radial", cex=1.2, lwd=3,
9     main="Neighbor-Joining tree of populations")
10 # Labels - bootstrap
11 nodelabels(round(hauss.boot.pop/10), frame="none")
12 # Print information about phylo object
13 print.phylo(hauss.nj.pop)
```

NJ tree based on DNA sequences

```
1 usflu.tree <- nj(X=usflu.dist)
2 # Plot it
3 plot.phylo(x=usflu.tree, type="unrooted", show.tip=FALSE)
4 title("Unrooted NJ tree")
5 # Coloured tips
6 usflu.pal <- colorRampPalette(topo.colors(length(levels(as.factor(
7 usflu.annot[["year"]])))))
8 # Tip labels
9 tiplabels(text=usflu.annot$year, bg=num2col(usflu.annot$year,
10 col.pal=usflu.pal), cex=0.75)
11 # Legend - describing years - pretty() automatically shows best
12 # values from given range, num2col() selects colours from colour scale
13 legend(x="bottomright", fill=num2col(x=pretty(x=1993:2008, n=8),
14 col.pal=usflu.pal), leg=pretty(x=1993:2008, n=8), ncol=1)
```

Root the tree

```
1 # Root the tree - "outgroup" is name of accession (in quotation
2 # marks) or number (position within phy object)
3 usflu.tree.rooted <- root(phy=usflu.tree, outgroup=1)
4 # Plot it
5 plot.phylo(x=usflu.tree.rooted, show.tip=FALSE, edge.width=2)
6 title("Rooted NJ tree")
7 # Labeling of tips
8 tiplabels(text=annot$year, bg=transp(num2col(x=annot$year,
9   col.pal=usflu.pal), alpha=0.7), cex=0.75, fg="transparent")
10 # Add axis with phylogenetic distance
11 axisPhylo()
12 # Legend - describing years - pretty() automatically shows best
13 # values from given range, num2col() selects colours from colour scale
14 legend(x="topright", fill=num2col(x=pretty(x=1993:2008, n=8),
15   col.pal=usflu.pal), leg=pretty(x=1993:2008, n=8), ncol=1)
```

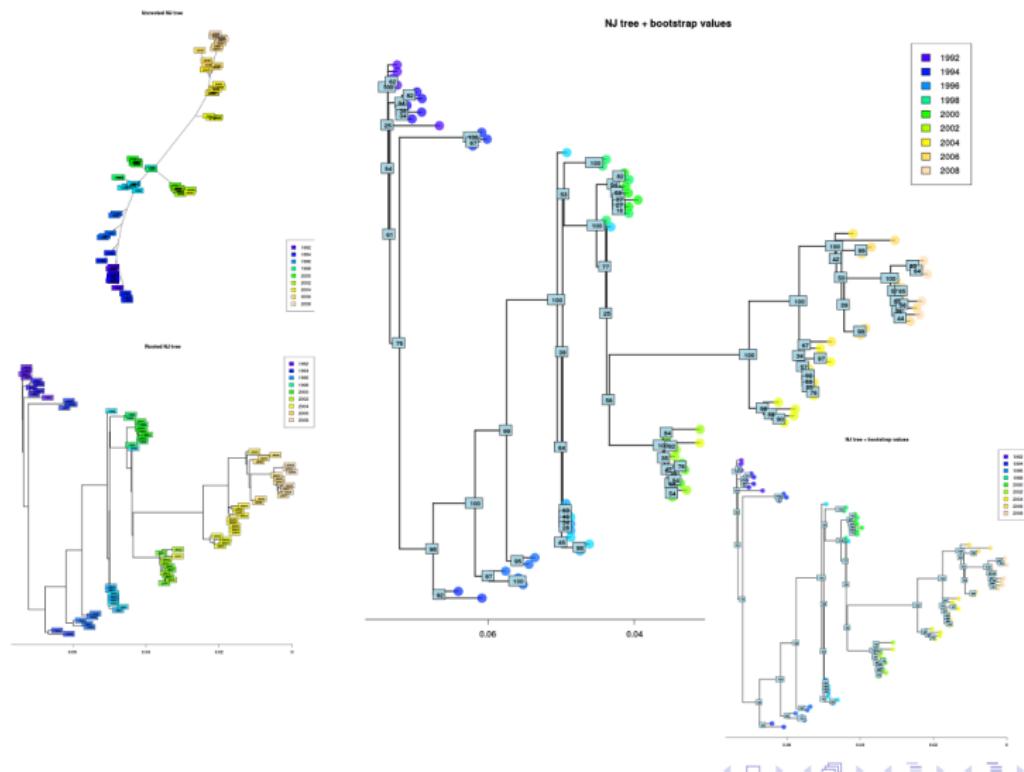
Bootstrap rooted tree

```
1 # Calculate it
2 usflu.boot <- boot.phylo(phy=usflu.tree.rooted, x=usflu.dna,
3   FUN=function(EEE) root(nj(dist.dna(EEE, model="TN93")),
4   outgroup=1), B=1000)
5 # Plot the tree
6 plot.phylo(x=usflu.tree.rooted, show.tip=FALSE, edge.width=2)
7 title("NJ tree + bootstrap values")
8 tiplabels(frame="none", pch=20, col=transp(num2col(x=annot[["year"]]),
9   col.pal=usflu.pal), alpha=0.7, cex=3.5, fg="transparent")
10 axisPhylo()
11 # Legend - describing years - pretty() automatically shows best
12 # values from given range, num2col() selects colours from colour scale
13 legend(x="topright", fill=num2col(x=pretty(x=1993:2008, n=8),
14   col.pal=usflu.pal), leg=pretty(x=1993:2008, n=8), ncol=1)
15 # Plots bootstrap support - note usflu.boot contains raw numbers
16 # transform it into percent
17 nodelabels(text=round(usflu.boot/10), cex=0.75)
```

Collapse branches with low bootstrap support

```
1 usflu.tree.temp <- usflu.tree.rooted
2 # Determine branches with low support - note BS values are in raw
3 # numbers - use desired percentage with respect to number of bootstraps
4 usflu.tocollapse <- match(x=which(usflu.boot < 700) +
5   length(usflu.tree.rooted$tip.label), table=usflu.tree.temp$edge[,2])
6 # Set length of bad branches to zero
7 usflu.tree.temp$edge.length[usflu.tocollapse] <- 0
8 # Create new tree
9 usflu.tree.collapsed <- di2multi(phy=usflu.tree.temp, tol=0.00001)
10 # Plot the consensus tree
11 plot.phylo(x=usflu.tree.collapsed, show.tip=FALSE, edge.width=2)
12 title("NJ tree after collapsing weak nodes")
13 tiplabels(text=annot$year, bg=transp(num2col(x=annot[["year"]]),
14   col.pal=usflu.pal), alpha=0.7, cex=0.5, fg="transparent")
15 axisPhylo()
16 legend(x="topright", fill=num2col(x=pretty(x=1993:2008, n=8),
17   col.pal=usflu.pal), leg=pretty(x=1993:2008, n=8), ncol=1)
```

The trees



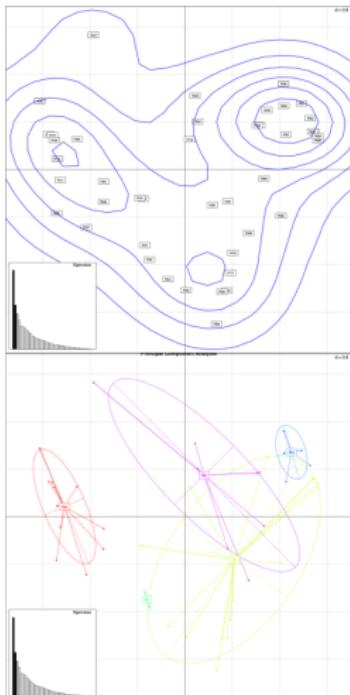
NJ is death. Long live NJ!

- “Basic” NJ has many limitations – there are several tries to overcome them
- Package **phangorn** has functions **NJ()** and unweighted version **UNJ()**
- Package **ape** has functions **njs()** and **bionjs()** which are designed to perform well on distances with (more) missing values
- Function **bionj()** from **ape** implements BIONJ algorithm
- FastME functions (package **ape**) perform the minimum evolution algorithm and aim to be replacement of NJ — read **?fastme** before use
- All those functions read distance matrix and their use is same as with “classical” **nj()** (read manual pages before using them) — it is also from package **ape**

PCoA I

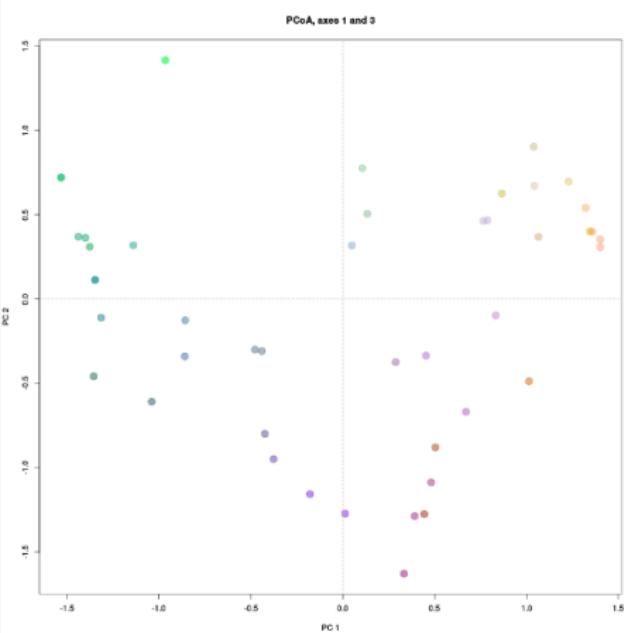
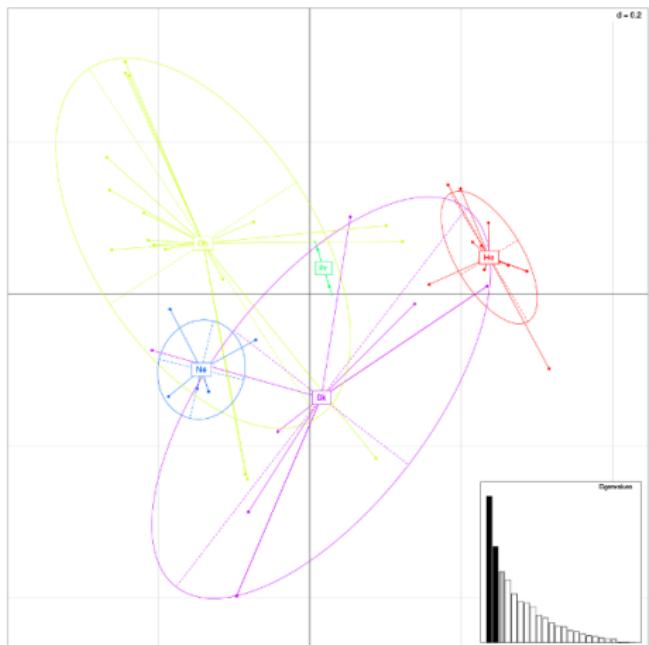
```
1 hauss.pcoa <- dudi.pco(d=dist(x=scaleGen(x=hauss.genind, center=TRUE,
2   scale=FALSE, truenames=TRUE, missing="NA"), method="euclidean"),
3   scannf=FALSE, nf=3)
4 # Basic display
5 s.label(dfxy=hauss.pcoa$li, clabel=0.75)
6 # To plot different axes use for example dfxy=hauss.pcoa$li[c(2, 3)]
7 # Add kernel density
8 s.kde2d(dfxy=hauss.pcoa$li, cpoint=0, add.plot=TRUE)
9 # Adds histogram of Eigenvalues
10 add.scatter.eig(w=hauss.pcoa$eig, nf=3, xax=1, yax=2,
11   posi="bottomleft", sub="Eigenvalues")
12 # Colored display according to populations
13 # Creates vector of colors according to populations
14 hauss.pcoa.col <- rainbow(length(levels(pop(hauss.genind))))
15 s.class(dfxy=hauss.pcoa$li, fac=pop(hauss.genind), col=hauss.pcoa.col)
16 add.scatter.eig(w=hauss.pcoa$eig, nf=3, xax=1, yax=2,
17   posi="bottomleft", sub="Eigenvalues")
18 title("Principal Coordinates Analysis") # Adds title to the graph
```

PCoA II



```
1 hauss.pcoa.bruvo <- dudi.pco(d=
2   bruvo.dist(pop=hauss.genind,
3   rep=rep(2, 12)), scannf=F,
4   nf=3)
5 s.class(dfxy=hauss.pcoa.bruvo$li,
6   fac=pop(hauss.genind),
7   col=hauss.pcoa.col)
8 add.scatter.eig(hauss.pcoa.bruvo$eig,
9   posi="bottomright", 3, 1, 2)
# Another possibility for coloured
# plot (see ?colorplot for details)
12 colorplot(xy=hauss.pcoa$li[c(1, 2)],
13   X=hauss.pcoa$li, transp=TRUE,
14   cex=3, xlab="PC 1", ylab="PC 2")
15 title(main="PCoA, axes 1 and 3")
16 abline(v=0, h=0, col="grey", lty=2)
```

PCoA — Bruvo and colorplot



DAPC

- Discriminant Analysis of Principal components
- Runs K-means Bayesian clustering on data transformed with PCA (reduces number of variables, speeds up process)
- Finally it runs discriminant analysis
- Various modes of displaying of results — “Structure-like”, “PCA-like” and more
- More information at <http://adegenet.r-forge.r-project.org/> and **adegenetTutorial("dapc")**
- If following commands would seem to complicated to you, try web interface by this command:

```
adegenetServer("DAPC")
```

K-find — Bayesian K-means clustering

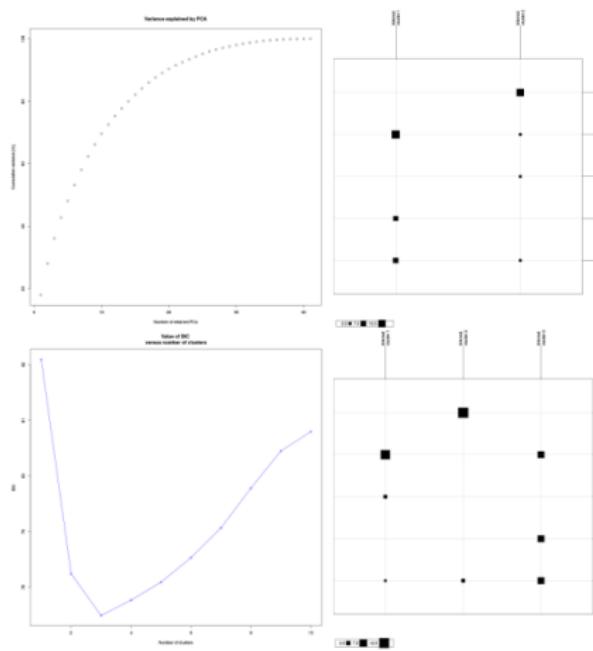
```
1 # Retain all informative PC (here about 35).
2 # According to second graph select best K (here 2 or 3).
3 hauss.kfind <- find.clusters(x=hauss.genind, stat="BIC",
4   choose.n.clust=TRUE, max.n.clust=10, n.iter=100000, n.start=100,
5   scale=FALSE, truenames=TRUE)
6 # See results as text
7 table(pop(hauss.genind), hauss.kfind$grp)
8 hauss.kfind
9 # Graph showing table of original and inferred populations and
10 # assignment of individuals
11 table.value(df=table(pop(hauss.genind), hauss.kfind$grp),
12   col.lab=paste("Inferred\ncluster", 1:length(hauss.kfind$size)), grid=T)
13 # For K=3 - note parameters n.pca and n.clust - we just rerun the
14 # analysis and when results are stable, no problem here
15 hauss.kfind3 <- find.clusters(x=hauss.genind, n.pca=35, n.clust=3,
16   stat="BIC", choose.n.clust=FALSE, max.n.clust=10, n.iter=100000,
17   n.start=100, scale=FALSE, truenames=TRUE)
```

K-find outputs

```

1 # See results as text
2 table(pop(hauss.genind),
3     hauss.kfind3$grp)
4 hauss.kfind3
5 # Graph showing table of original
6 # and inferred populations and
7 # assignment of individuals
8 table.value(
9     df=table(pop(hauss.
10    genind), hauss.kfind3$grp),
11    col.lab= paste("Inferred\ncluster",
12    1:length(hauss.kfind3$size)),
13    grid=TRUE)

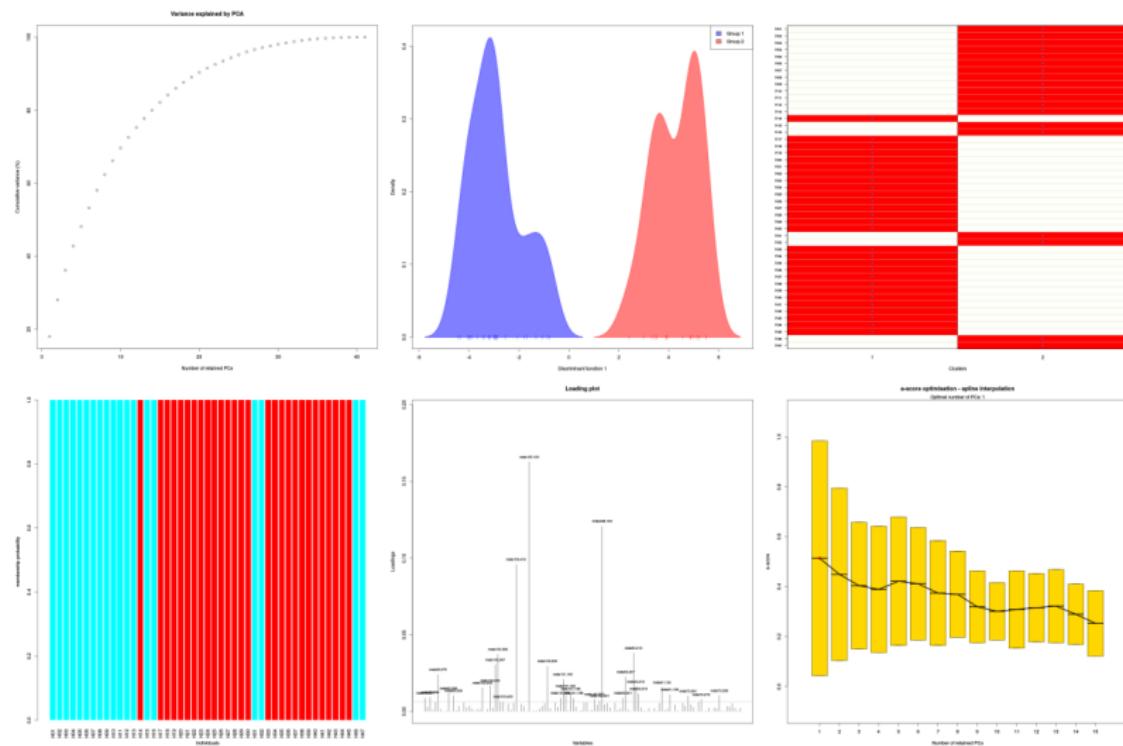
```



DAPC code I

```
1 ## K=2
2 # Creates DAPC
3 # Number of informative PC (Here 15, adegenet recommends < N/3). Select
4 # number of informative DA (here only one is available - no PCA graph)
5 hauss.dapc <- dapc(x=hauss.genind, pop=hauss.kfind$grp, center=TRUE,
6 scale=FALSE, var.contrib=TRUE, pca.info=TRUE, truenames=TRUE)
7 # Information
8 hauss.dapc
9 # Density functions
10 scatter(x=hauss.dapc, xax=1, yax=1, main="DAPC", bg="white", solid=0.5,
11 leg=TRUE, txt.leg=c("Group 1", "Group 2"), posi.leg="topright")
12 # Assignment of individuals to clusters
13 assignplot(x=hauss.dapc)
14 # Structure-like plot
15 compoplot(x=hauss.dapc, xlab="Individuals", leg=FALSE)
16 # Loadingplot - alleles the most adding to separation of individuals
17 loadingplot(x=hauss.dapc$var.contr)
```

DAPC for K=2

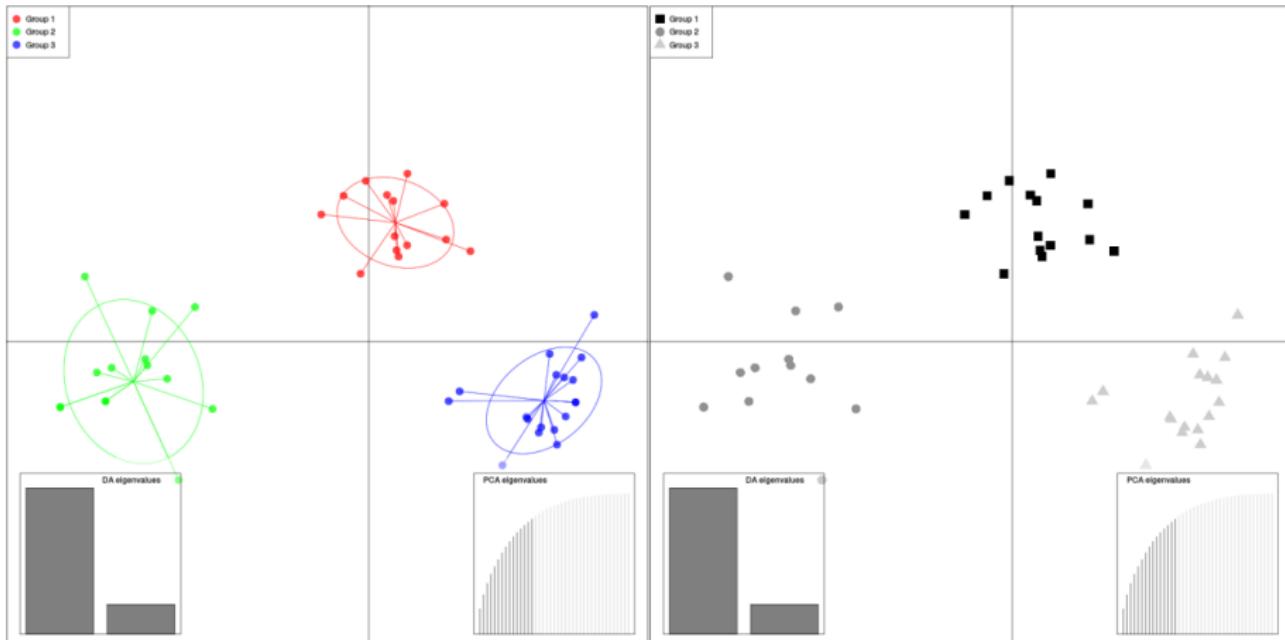


DAPC code II

```
1 ## alfa-score - according to number of PC axis
2 optim.a.score(x=hauss.dapc)
3 # K=3
4 # Creates DAPC
5 # Number of informative PC (Here 15, adegenet recommends < N/3).
6 # Select number of informative DA (here 2).
7 hauss.dapc3 <- dapc(x=hauss.genind, pop=hauss.kfind3$grp, center=TRUE,
8 scale=FALSE, var.contrib=TRUE, pca.info=TRUE, truenames=TRUE)
9 # Information
10 hauss.dapc
11 # A la PCA graph
12 scatter(x=hauss.dapc3, main="DAPC, Cardamine", bg="white", cex=3,
13 clab=0, col=rainbow(3), posi.da="bottomleft", scree.pca=TRUE,
14 posi.pca="bottomright", leg=TRUE, txt.leg=c("Group 1", "Group 2",
15 "Group 3"), posi.leg="topleft")
```

Especially graphical parameters have huge possiilities...

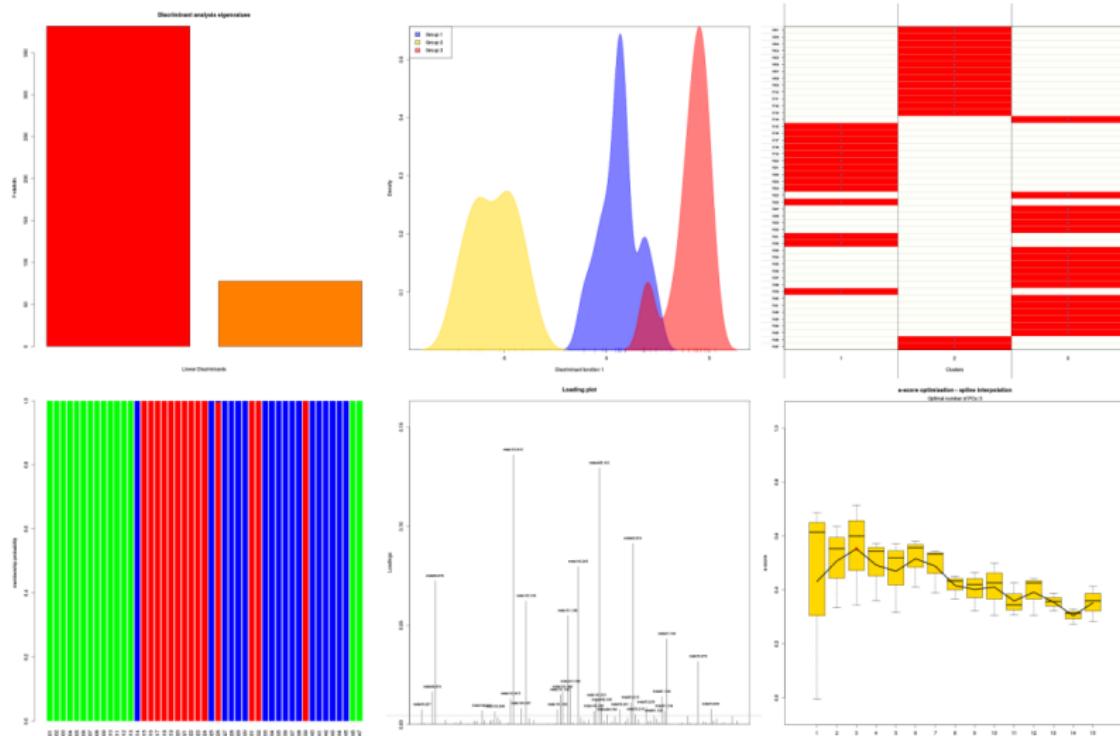
DAPC for K=3



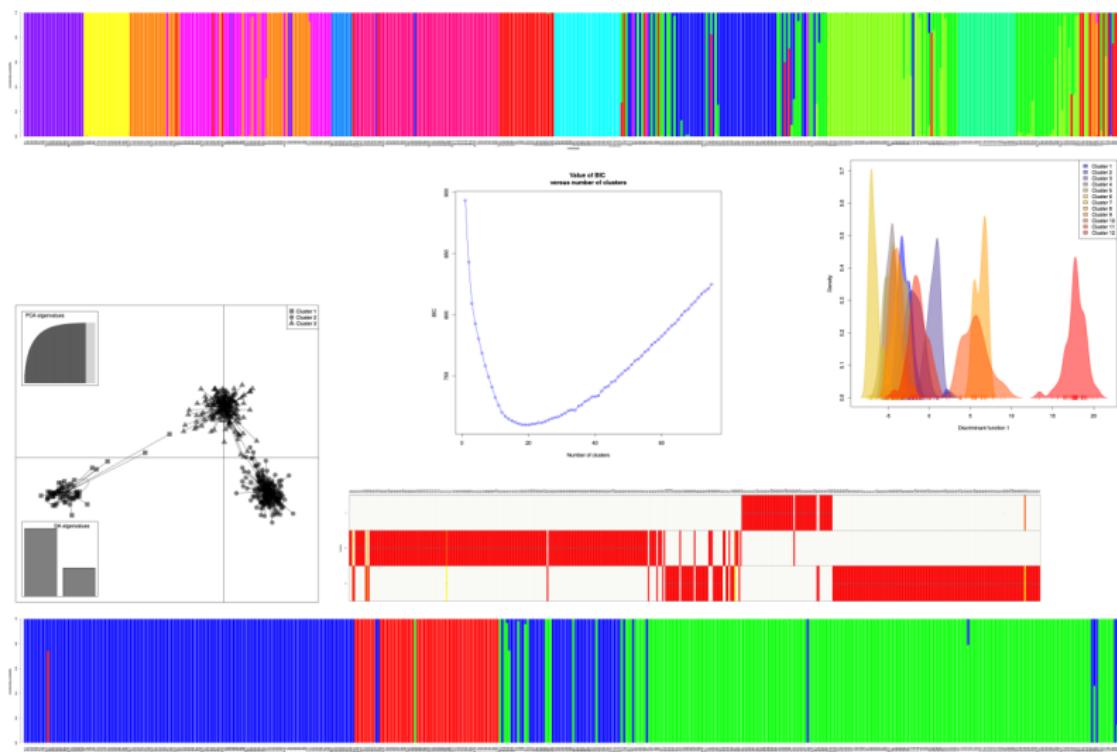
DAPC code III

```
1 # Same in BW
2 scatter(x=hauss.dapc3, main="DAPC, Cardamine", bg="white", pch=c(15:17),
3   cell=0, cstar=0, solid=1, cex=2.5, clab=0, col=grey.colors(3, start=0,
4   end=0.8, gamma=2, alpha=0), posi.da="bottomleft", scree.pca=TRUE,
5   posi.pca="bottomright", leg=TRUE, txt.leg=c("Group 1", "Group 2",
6   "Group 3"), posi.leg="topleft")
7 # Density functions
8 scatter(x=hauss.dapc3, xax=1, yax=1, main="DAPC", bg="white", solid=0.5,
9   leg=T, txt.leg=c("Group 1", "Group 2", "Group 3"), posi.leg="topleft")
10 # Assignment of individuals to clusters
11 assignplot(hauss.dapc3)
12 # Structure-like plot
13 compoplot(hauss.dapc3, xlab="Individuals", leg=FALSE)
14 # Loadingplot - alleles the most adding to separation of individuals
15 loadingplot(hauss.dapc3$var.contr)
16 # alfa-score - according to number of PC axis
17 optim.a.score(hauss.dapc3)
```

DAPC for K=3, extra information



Another DAPC example



Some functions to work with huge SNP data sets

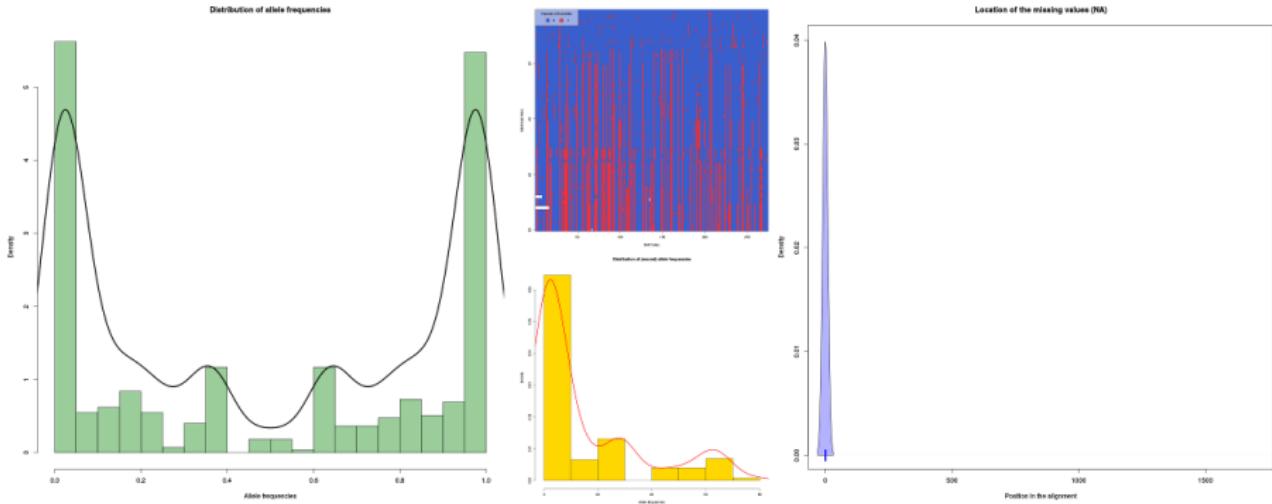
```
1 # Plot of missing data (white) and number of 2nd alleles
2 glPlot(x=usflu.genlight, legend=TRUE, posi="topleft")
3 # Sum of the number of second allele in each SNP
4 usflu.freq <- glSum(usflu.genlight)
5 # Plot distribution of (second) allele frequencies
6 hist(x=usflu.freq, proba=TRUE, col="gold", xlab="Allele
7   frequencies", main="Distribution of (second) allele frequencies")
8 lines(x=density(usflu.freq)$x, y=density(usflu.freq)$y*1.5,
9   col="red", lwd=3 )
10 # Mean number of second allele in each SNP
11 usflu.mean <- glMean(usflu.genlight)
12 usflu.mean <- c(usflu.mean, 1-usflu.mean)
13 # Plot distribution of allele frequencies
14 hist(x=usflu.mean, proba=TRUE, col="darkseagreen3",
15   xlab="Allele frequencies", main="Distribution of allele
16   frequencies", nclass=20)
17 lines(x=density(usflu.mean, bw=0.05)$x, y=density(usflu.mean,
18   bw=0.05)$y*2, lwd=3)
```

Number of missing values in each locus

```
1 # Play with bw parameter to get optimal image
2 usflu.na.density <- density(glNA(usflu.genlight), bw=10)
3 # Set range of xlim parameter from 0 to the length
4 # of original alignment
5 plot(x=usflu.na.density, type="n", xlab="Position in the alignment",
6     main="Location of the missing values (NA)", xlim=c(0, 1701))
7 polygon(c(usflu.na.density$x, rev(usflu.na.density$x)),
8         c(usflu.na.density$y, rep(0, length(usflu.na.density$x))),
9         col=transp("blue", alpha=0.3))
10 points(glNA(usflu.genlight), rep(0, nLoc(usflu.genlight)),
11        pch="|", cex=2, col="blue")
```

Those tools are designed mainly for situation when having multiple (nearly) complete genomes. It is not needed for smaller (normal) datasets. Lets keep hoping in fast development of computers...

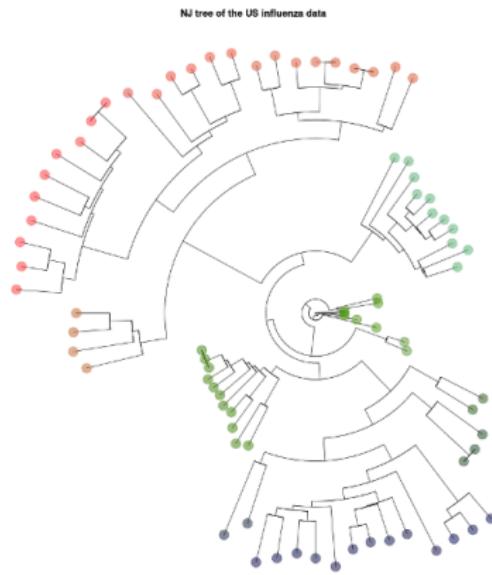
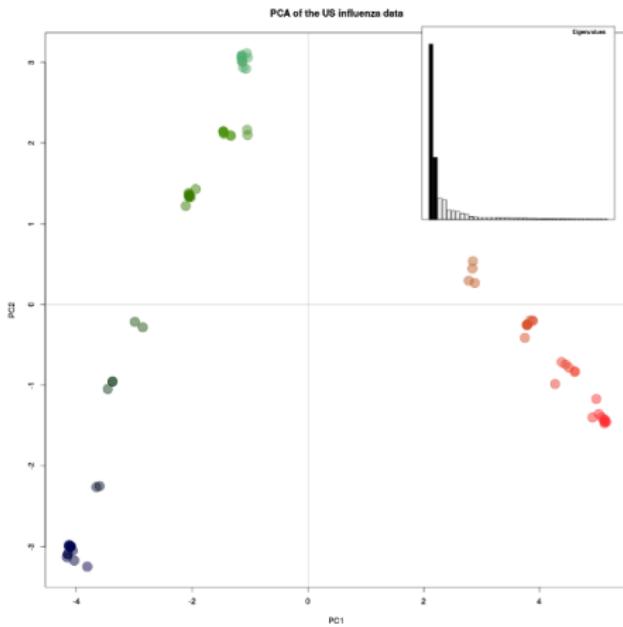
Basic information about SNP: distribution of 2nd allele frequencies, missing data and number of 2nd allele, distribution of allele frequencies, and number of missing values in each locus



PCA, NJ and genlight objects

```
1 usflu.pca <- glPca(x=usflu.genlight, center=TRUE, scale=FALSE,
2   loadings=TRUE) # Select number of retained PC axes, about 10 here
3 # Plot PCA
4 scatter.glPca(x=usflu.pca, posi="bottomright")
5 title("PCA of the US influenza data")
6 # Coloured plot
7 colorplot(usflu.pca$scores, usflu.pca$scores, transp=TRUE, cex=4)
8 title("PCA of the US influenza data")
9 abline(h=0, v=0, col="grey")
10 add.scatter.eig(usflu.pca[["eig"]][1:40], 2, 1, 2, posi="topright",
11   inset=0.05, ratio=0.3)
12 # Calculate phylogenetic tree
13 usflu.tree.genlight <- nj(dist(as.matrix(usflu.genlight)))
14 # Plot coloured phylogenetic tree
15 plot.phylo(x=usflu.tree.genlight, typ="fan", show.tip=FALSE)
16 tiplabels(pch=20, col=usflu.col, cex=4)
17 title("NJ tree of the US influenza data")
```

PCA, NJ and genlight objects



Short overview of spatial genetics (in R)

- Basic approaches

- Moran's I — several implementations (here as basic autocorrelation index, in sPCA and in Monmonier's algorithm), generally it is autocorrelation coefficient with broader use
- Mantel test — several implementations, popular, although recently criticised as biologically irrelevant, generally correlation of two matrices (here genetic and geographical)
- Bayesian clustering using geographical information as a proxy and showing results in geographical context (here as implemented in Geneland)
- There are unlimited possibilities with connections with GIS software — check specialised courses and literature

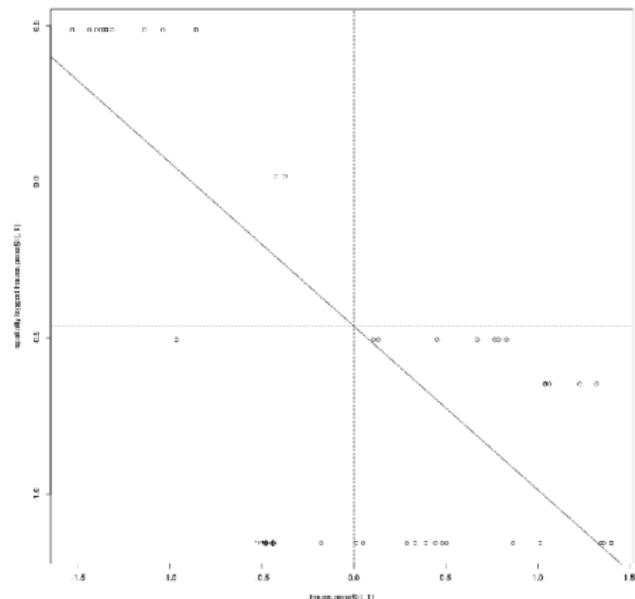
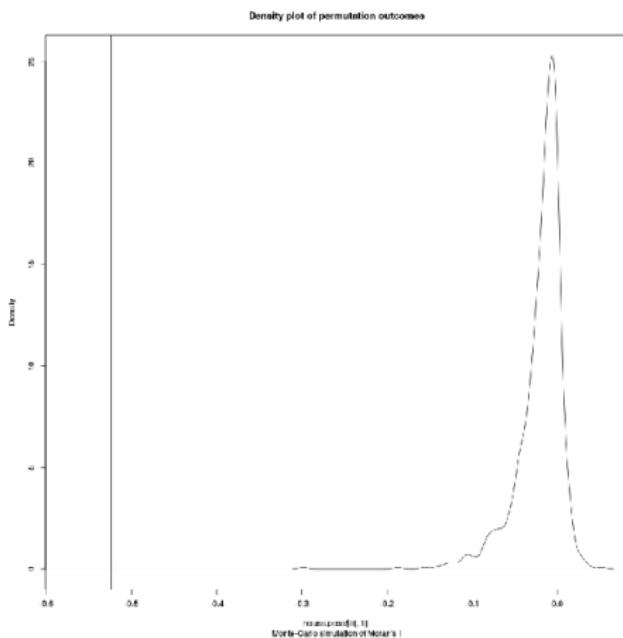
Calculation of Moran's I

```
1 # Creates connection network
2 hauss.connectivity <- chooseCN(xy=hauss.genind$other$xy, type=5,
3   d1=0, d2=1, plot.nb=TRUE, result.type="listw", edit.nb=FALSE)
4 hauss.connectivity
5 # Output:
6 Characteristics of weights list object:
7 Neighbour list object:
8 Number of regions: 47
9 Number of nonzero links: 1120
10 Percentage nonzero weights: 50.70167
11 Average number of links: 23.82979
12 Weights style: W
13 Weights constants summary:
14     n    nn S0      S1      S2
15 W 47 2209 47 4.033607 207.037
16 ###
```

Calculation of Moran's I / II

```
1 # Test of Moran's I for 1st PCoA axis
2 # Results can be checked against permuted values of moran.mc()
3 moran.test(x=hauss.pcoa$li[,1], listw=hauss.connectivity,
4   alternative="greater", randomisation=TRUE)
5 # Output:
6 Moran's I test under randomisation
7 data: hauss.pcoa$li[, 1]
8 weights: hauss.connectivity
9 Moran I statistic standard deviate = -18.514, p-value = 1
10 alternative hypothesis: greater
11 sample estimates:
12 Moran I statistic      Expectation      Variance
13 -0.5232003724      -0.0217391304      0.0007336276
14 hauss.pcoa1.mctest <- moran.mc(x=hauss.pcoa$li[,1],
15   listw=hauss.connectivity, alternative="greater", nsim=1000)
16 hauss.pcoa1.mctest
17 plot(hauss.pcoa1.mctest) # Plot of density of permutations
18 moran.plot(x=hauss.pcoa$li[,1], listw=hauss.connectivity) # PC plot
```

Moran's I for our 1st axis isn't significant

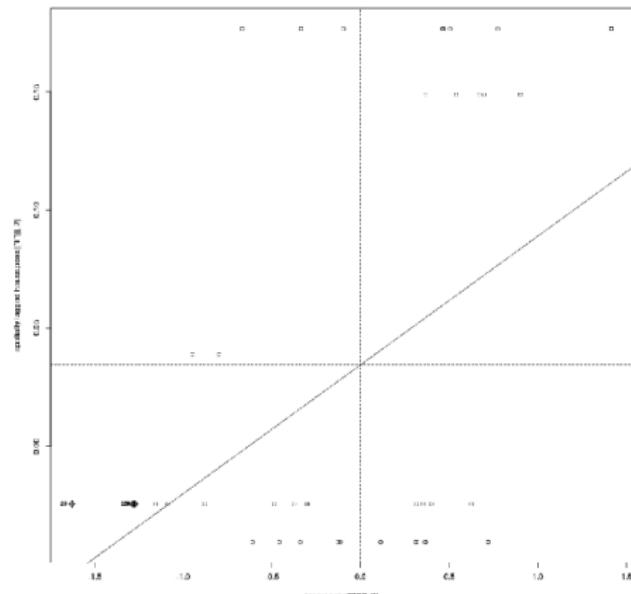
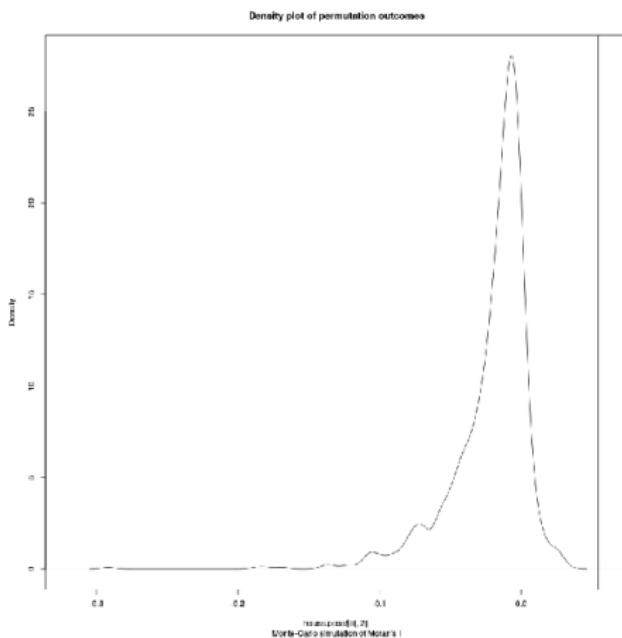


We tested hypothesis “greater” — there is **no** significant positive autocorrelation

Calculation of Moran's I (2nd axis)

```
1 # Test of Moran's I for 2nd PCoA axis
2 moran.test(x=hauss.pcoa$li[,2], listw=hauss.connectivity,
3   alternative="greater", randomisation=TRUE)
4 hauss.pcoa2.mctest <- moran.mc(x=hauss.pcoa$li[,2],
5   listw=hauss.connectivity, alternative="greater", nsim=1000)
6 hauss.pcoa2.mctest
7 # Output
8 Monte-Carlo simulation of Moran's I
9 data: hauss.pcoa$li[, 2]
10 weights: hauss.connectivity
11 number of simulations + 1: 1001
12 statistic = 0.0545, observed rank = 1001, p-value = 0.000999
13 alternative hypothesis: greater
14 # Plots
15 plot(hauss.pcoa2.mctest)
16 moran.plot(x=hauss.pcoa[,2], listw=hauss.connectivity)
```

Second axis is surprisingly significant



We tested hypothesis "greater" — there **is** significant positive autocorrelation

Spatial Analysis of Principal Components (sPCA)

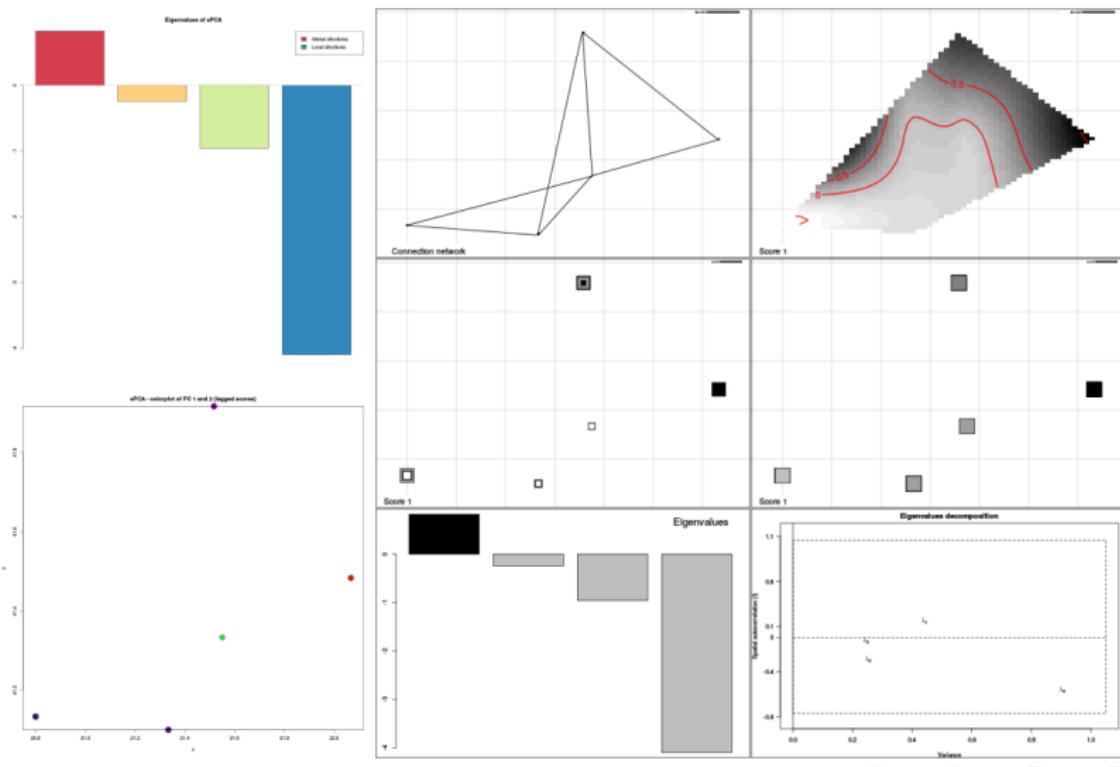
- Implemented in adegenet, see **adegenetTutorial("spca")**
- Analyses matrix of relative allele frequencies of genotypes/populations and spatial weighting matrix
- The geographical matrix is usually (as for Moran's *I*) created by **chooseCN()** — creates connectivity network among entities (genotypes/populations) — spatial coordinates are not directly used
- When using **chooseCN()**, look at the documentation and try various methods with changing settings to see differences

```
1 data(rupica) # Loads adegenet's training dataset
2 # Try various settings vor chooseCN (type=X) - type 1-4 as there
3 # are identical coordinates (multiple sampling from same locality)
4 chooseCN(xy=rupica$other$xy, ask=T, type=5/6/7, plot.nb=T, edit.nb=F)
5 # See ?chooseCN for more details
```

Calculations of sPCA

```
1 hauss.spca <- spca(obj=hauss.genind, cn=hauss.connectivity,
2   scale=TRUE, scannf=TRUE)
3 # Plot eigenvalues of sPCA - global vs. local structure
4 barplot(height=hauss.spca$eig, main="Eigenvalues of sPCA",
5   col=spectral(length(hauss.spca$eig)))
6 legend("topright", fill=spectral(2), leg=c("Global structures",
7   "Local structures")) # Add legend
8 abline(h=0, col="grey") # Add line showing zero
9 print.spca(hauss.spca) # Information about sPCA
10 summary.spca(hauss.spca) # Summary of sPCA results
11 # Shows connectivity network, 3 different scores
12 # barplot of eigenvalues and eigenvalues decomposition
13 plot.spca(hauss.spca)
14 colorplot.spca(hauss.spca, cex=3) # Display of scores in color canals
15 title("sPCA - colorplot of PC 1 and 2 (lagged scores)", line=1, cex=1.5)
16 # Spatial and variance components of the eigenvalues
17 screeplot.spca(x=hauss.spca, main=NULL)
```

sPCA outputs

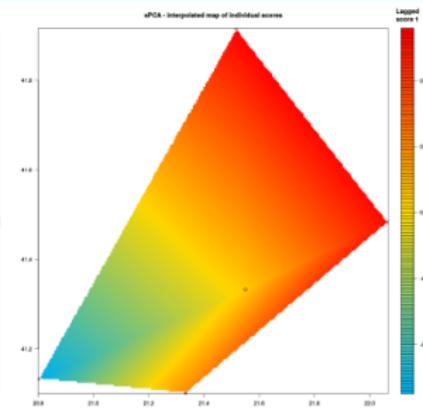
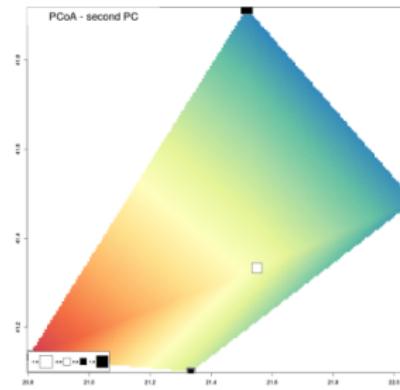
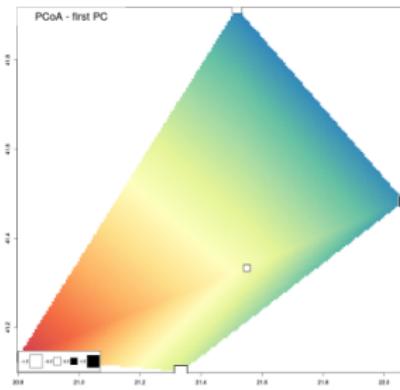


Map of genetic clines

```
1 library(akima) # This library is needed for manipulation with coordinates
2 hauss.spca.temp <- interp(other(hauss.genind)$xy[,1], other(
3   hauss.genind)$xy[,2], hauss.spca$ls[,1], xo=seq(min(other(
4     hauss.genind)$xy[,1]), max(other(hauss.genind)$xy[,1]), le=200),
5   yo=seq(min(other(hauss.genind)$xy[,2]), max(other(hauss.genind)$xy[,2]),
6   le=200), duplicate="median")
7 # For 1st axis
8 image(x=hauss.spca.temp, col=spectral(100))
9 s.value(dfxy=hauss.genind$other$xy, z=hauss.pcoa$li[,1], add.p=TRUE,
10   csizen=0.5, sub="PCoA - first PC", csub=2, possub="topleft")
11 # For 2nd axis
12 image(x=hauss.spca.temp, col=spectral(100))
13 s.value(dfxy=hauss.genind$other$xy, z=hauss.pcoa[["li"]][,2], add.p=TRUE,
14   csizen=0.5, sub="PCoA - second PC", csub=2, possub="topleft")
```

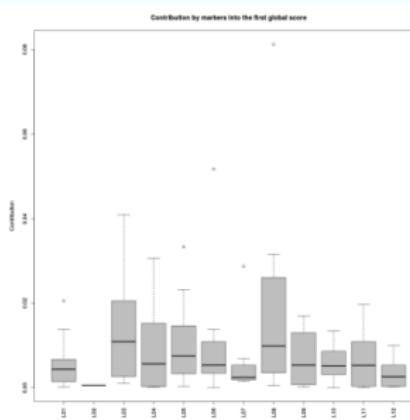
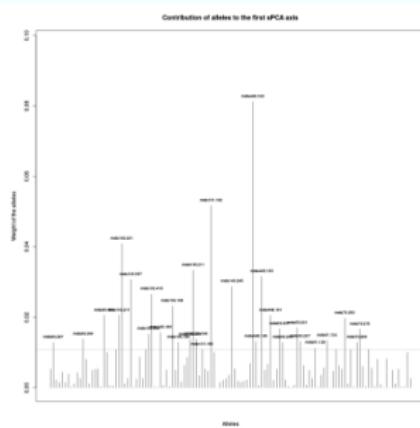
sPCA outputs

```
1 # Interpolated lagged score on a map
2 hauss.spca.annot <- function() {
3   title("sPCA - interpolated map of individual scores")
4   points(other(hauss.genind)$xy[,1], other(hauss.genind)$xy[,2])
5 }
6 filled.contour(hauss.spca.temp, color.pal=colorRampPalette(
7   lightseasun(100)), pch=20, nlevels=100, key.title=title("Lagged\n
8   score 1"), plot.title=hauss.spca.annot())
```



Loading plots - which alleles contribute the most?

```
1 hauss.spca.loadings <- hauss.spca[["c1"]][,1]^2
2 names(hauss.spca.loadings) <- rownames(hauss.spca$c1)
3 loadingplot(x=hauss.spca.loadings, xlab="Alleles", ylab="Weight of the
4     alleles", main="Contribution of alleles to the first sPCA axis")
5 boxplot(formula=hauss.spca.loadings~hauss.genind$loc.fac, las=3,
6     ylab="Contribution", xlab="Marker", main="Contribution by markers
7     into the first global score", col="grey")
```



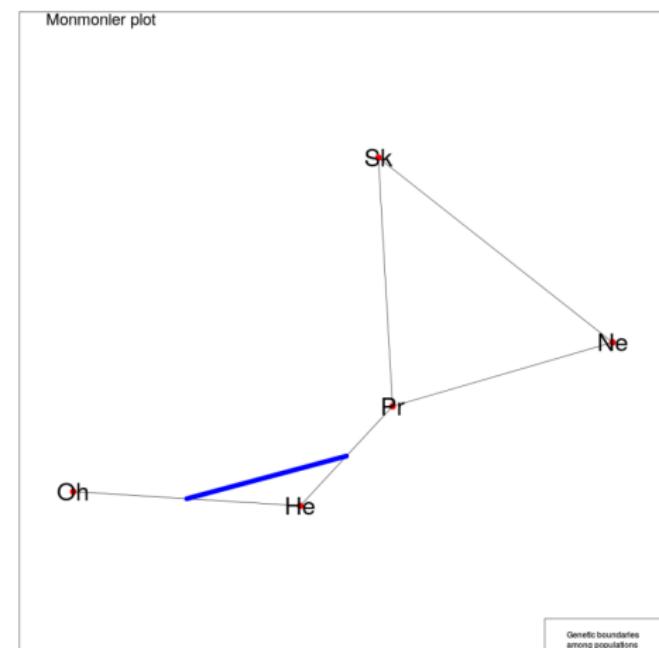
Monmonier's algorithm — genetic boundaries

- Finds boundaries of maximum differences between contiguous polygons of a tessellation
- Detects genetic boundaries among georeferenced genotypes (or populations)
- For more information see **adegenetTutorial("basics")**
- Requires every point to have unique coordinates (use **jitter()** when having same coordinates from one locality)
- Example here is on population level, which is not ideal for all cases

```
1 # Calculates Monmonier's function (for threshold use 'd')
2 hauss.monmonier <- monmonier(xy=hauss.genpop$other$xy, dist=
3   dist(hauss.genpop@tab), cn=chooseCN(hauss.genpop$other$xy,
4   ask=FALSE, type=2, plot.nb=FALSE, edit.nb=FALSE), nrun=1)
5 coords.monmonier(hauss.monmonier) # See result as text
```

Plot genetic boundaries

```
1 plot.monmonier(hauss.monmonier,
2   add.arrows=FALSE, bwd=10,
3   sub="Monmonier plot", csub=2)
4 points(hauss.genpop$other$xy,
5   cex=2.5, pch=20, col="red")
6 text(x=hauss.genpop$other$xy$lon,
7   y=hauss.genpop$other$xy$lat,
8   labels=hauss.genpop@pop.names,
9   cex=3)
10 legend("bottomright",
11   leg="Genetic boundaries\namong populations")
12 # See ?points, ?text and ?legend
```



Monmonier notes

- Sometimes it is needed to get rid of random noise in data. To do so use as parameter **dist** of **monmonier()** table from PCA (**pcaObject\$li**):

```
1 monmonier(..., dist=dudi.pco(d=dist(x=GenindObject$tab),  
2 scannf=FALSE, nf=1)$li, ...)
```

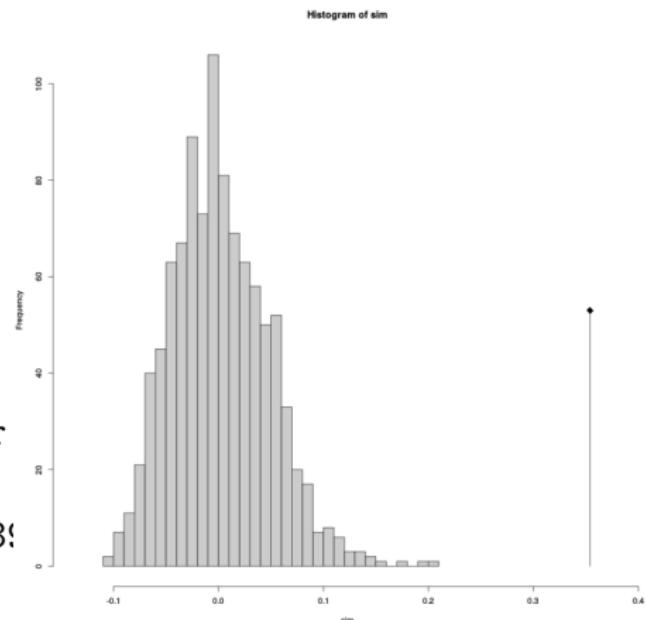
- Generally (when dataset is bigged and more diverse) it is recommended to run it several times (parameter **nrun**) — there will be several iterations
- See **?plot.monmonier** for various graphical parameters to customize the plot
- Use **points()** to add for example coloured symbols of samples and/or **text()** to add text labels

Mantel test — isolation by distance

```
1 # Geographical distance
2 hauss.gdist <- dist(x=hauss.genind$other$xy, method="euclidean",
3   diag=TRUE, upper=TRUE)
4 # Mantel test
5 hauss.mantel <- mantel.randtest(m1=hauss.dist, m2=hauss.gdist,
6   nrepet=1000)
7 hauss.mantel # See text output
8 plot(hauss.mantel, nclass=30)
9 # Libraries required by mantel.correlog:
10 library(permute)
11 library(lattice)
12 library(vegan)
13 # Different implementation of Mantel test testing distance classes
14 hauss.mantel.cor <- mantel.correlog(D.eco=hauss.dist, D.geo=hauss.gdist,
15   XY=NULL, n.class=0, break.pts=NULL, cutoff=FALSE, r.type="pearson",
16   nperm=1000, mult="holm", progressive=TRUE)
17 # See results for respective classes
18 hauss.mantel.cor
19 summary(hauss.mantel.cor)
```

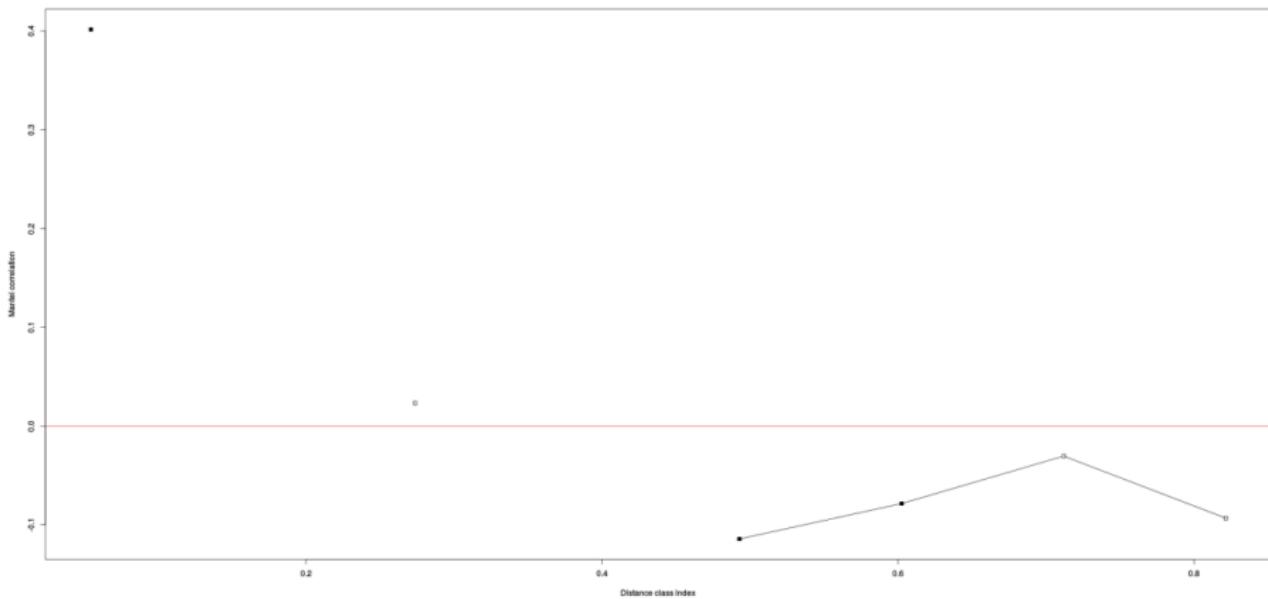
Mantel test outputs — strongly significant

```
> hauss.mantel
Monte-Carlo test
Call: mantel.randtest(m1 =
  hauss.dist, m2 =
  hauss.gdist, nrepet = 1000)
Observation: 0.35409
Based on 1000 replicates
Simulated p-value: 0.000999001
Alternative hypothesis: greater
  Std.Obs  Expectation  Variance
7.61967545 0.001687140 0.0021389
```



```
1 # Plot correlogram (next slide)
2 plot(hauss.mantel.cor)
```

Plot of Mantel Correlogram Analysis



About Geneland

- Works with haploid and diploid codominant markers (microsatellites or SNPs)
- Spatially explicit Bayesian clustering
- Produces maps of distribution of inferred genetic clusters
- Relative complicated tool with various modelling options

```
1 # Load needed libraries
2 library(PBSmapping)
3 library(RandomFields)
4 library(fields)
5 library(spam)
6 library(grid)
7 library(maps)
8 library(tcltk)
9 library(Geneland)
10 # Graphical interface is available
11 # we will use only command line
12 Geneland.GUI()
```

For more information see <http://www2.imm.dtu.dk/~gigu/Geneland/>

Loading and conversions of coordinates

```
1 # Geneland requires specific coordinate space
2 # hauss.coord is DF, we need just plain matrix
3 hauss.geneland.coord <- as.matrix(hauss.coord)
4 colnames(hauss.geneland.coord) <- c("X", "Y")
5 attr(hauss.geneland.coord, "projection") <- "LL"
6 attr(hauss.geneland.coord, "zone") <- NA
7 hauss.geneland.coord.utm <- convUL(hauss.geneland.coord)
8 dim(hauss.geneland.coord)
9 hauss.geneland.coord
10 dim(hauss.geneland.coord.utm)
11 hauss.geneland.coord.utm # Final coordinates
12 # Load data (only haploid or diploid data are supported)
13 # only plain table with alleles
14 hauss.geneland.data <- read.table(file=
15   "http://trapa.cz/sites/default/rcourse/haussknechtii_geneland.txt",
16   na.string="-999", header=FALSE, sep="\t")
17 dim(hauss.geneland.data)
18 hauss.geneland.data
```

Settings and running MCMC

```
1 hauss.geneland.nrun <- 5 # Set number of independent runs
2 hauss.geneland.burnin <- 100 # Set length of burnin chain
3 hauss.geneland.maxpop <- 10 # Set maximal K (number of populations)
4 # FOR loop will run several independent runs and produce output maps
5 # of genetic clusters - outputs are written into subdirectory within
6 # geneland directory (this has to exist prior launching analysis)
7 for (hauss.geneland.irun in 1:hauss.geneland.nrun) {
8   hauss.geneland.path.mcmc <- paste("geneland/", hauss.geneland.irun,
9     "/", sep="")
10  # On Windows, remove following line and create subdirectories from
11  # 1 to max K manually (creating subdirs in Windows is complicated)
12  system(paste("mkdir ", hauss.geneland.path.mcmc)) # Creates subdirs
13  # Inferrence - MCMC chain - see ?MCMC for details
14  MCMC(coordinates=hauss.geneland.coord.utm, geno.dip.codom=
15    hauss.geneland.data, path.mcmc=hauss.geneland.path.mcmc,
16    delta.coord=0.001, varnpop=TRUE, npopmin=1, npopmax=
17    hauss.geneland.maxpop, nit=10000, thinning=10,
18    freq.model="Uncorrelated", spatial=TRUE)
19  # For loop continues on next slide
```

Running MCMC

```
1 # Start of FOR loop is on previous page
2 # Post-process chains
3 PostProcessChain(coordinates=hauss.geneland.coord.utm,
4     path.mcmc=hauss.geneland.path.mcmc, nxdom=500, nydom=500,
5     burnin=hauss.geneland.burnin)
6 # Output
7 # Simulated number of populations
8 Plotnpop(path.mcmc=hauss.geneland.path.mcmc, printit=TRUE,
9     file=paste(hauss.geneland.path.mcmc, "/geneland-number_of_clusters
10    .pdf", sep=""), format="pdf", burnin=hauss.geneland.burnin)
11 dev.off() # We must close graphical device manually
12 # Map of estimated population membership
13 PosteriorMode(coordinates=hauss.geneland.coord.utm,
14     path.mcmc=hauss.geneland.path.mcmc, printit=TRUE, format="pdf",
15     file=paste(hauss.geneland.path.mcmc, "/geneland-map.pdf", sep=""))
16 dev.off() # We must close graphical device manually
17 } # End of FOR loop from previous slide
```

Estimate F_{ST}

```
1 # Prepare list to record values of Fst for all runs
2 hauss.geneland.fstat <- list()
3 # Estimate Fst
4 for (hauss.geneland.irun in 1:hauss.geneland.nrun) {
5   hauss.geneland.path.mcmc <- paste("geneland/", hauss.geneland.irun,
6                                     "/", sep="")
7   # F-statistics - Fis and Fst
8   hauss.geneland.fstat[[hauss.geneland.irun]] <- Fstat.output(
9     coordinates=hauss.geneland.coord.utm, genotypes=hauss.geneland.data,
10    burnin=hauss.geneland.burnin, ploidy=2,
11    path.mcmc=hauss.geneland.path.mcmc)
12 }
13 # Print Fst output
14 hauss.geneland.fstat
```

MCMC inference under the admixture model

```
1 for (hauss.geneland.irun in 1:hauss.geneland.nrun) {  
2   hauss.geneland.path.mcmc <- paste("geneland/",  
3   hauss.geneland.irun, "/", sep="")  
4   hauss.geneland.path.mcmc.adm <- paste(hauss.geneland.path.mcmc,  
5   "admixture", "/", sep="")  
6   # On Windows, remove following line of code and create in each  
7   # result directory (from 1 to max K) new subdirectory "admixture"  
8   # (creating subdirs in Windows is complicated)  
9   system(paste("mkdir ", hauss.geneland.path.mcmc.adm))  
10  HZ(coordinates=hauss.geneland.coord.utm, geno.dip.codom=  
11  hauss.geneland.data, path.mcmc.noadm=hauss.geneland.path.mcmc,  
12  nit=10000, thinning=10, path.mcmc.adm=hauss.geneland.path.mcmc.adm)  
13 }
```

Currently, there is no much use for admixture results, at least not without extra work...

Produce maps of respective inferred clusters

```
1 for (hauss.geneland.irun in 1:hauss.geneland.nrun) {  
2   hauss.geneland.path.mcmc <- paste("geneland/",  
3     hauss.geneland.irun, "/", sep="")  
4   # Maps - tessellations  
5   PlotTessellation(coordinates=hauss.geneland.coord.utm,  
6     path.mcmc=hauss.geneland.path.mcmc, printit=TRUE,  
7     path=hauss.geneland.path.mcmc)  
8   for (hauss.geneland.irun.img in 1:hauss.geneland.maxpop) {  
9     dev.off() } # We must close graphical device manually  
10 }
```

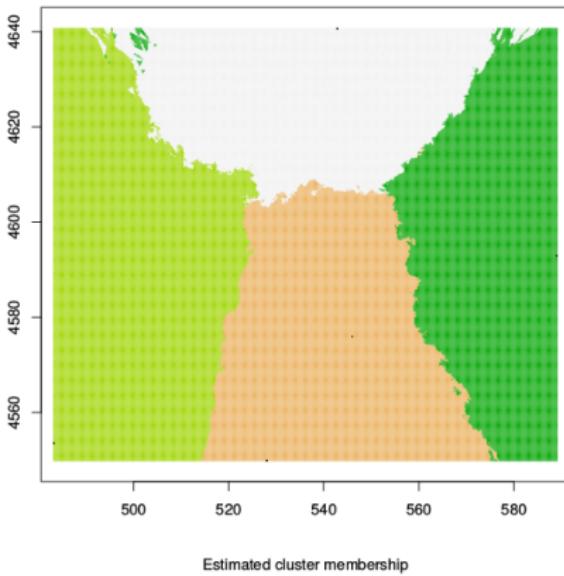
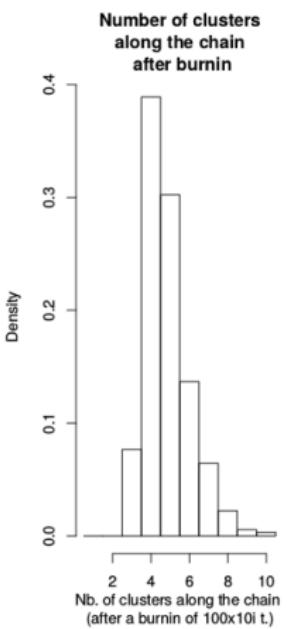
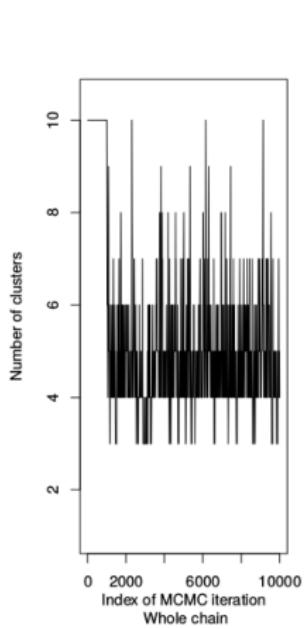
Maps are produced as PS (PostScript) files. Not every graphical software can handle them (try for example [GIMP](#)). There are as many plots as was maximal K, but only those up to inferred number of clusters have some content.

Determine which cluster is the best

```
1 # Calculate average posterior probability
2 hauss.geneland.lpd <- rep(NA, hauss.geneland.nrun)
3 for (hauss.geneland.irun in 1:hauss.geneland.nrun) {
4   hauss.geneland.path.mcmc <- paste("geneland/",
5     hauss.geneland.irun, "/", sep="")
6   hauss.geneland.path.lpd <- paste(hauss.geneland.path.mcmc,
7     "log.posterior.density.txt", sep="")
8   hauss.geneland.lpd[hauss.geneland.irun] <-
9     mean(scan(hauss.geneland.path.lpd)[-1:hauss.geneland.burnin]))
10 }
11 # Sorts runs according to decreasing posterior probability
12 # the first one is the best
13 order(hauss.geneland.lpd, decreasing=TRUE)
14 hauss.geneland.lpd
```

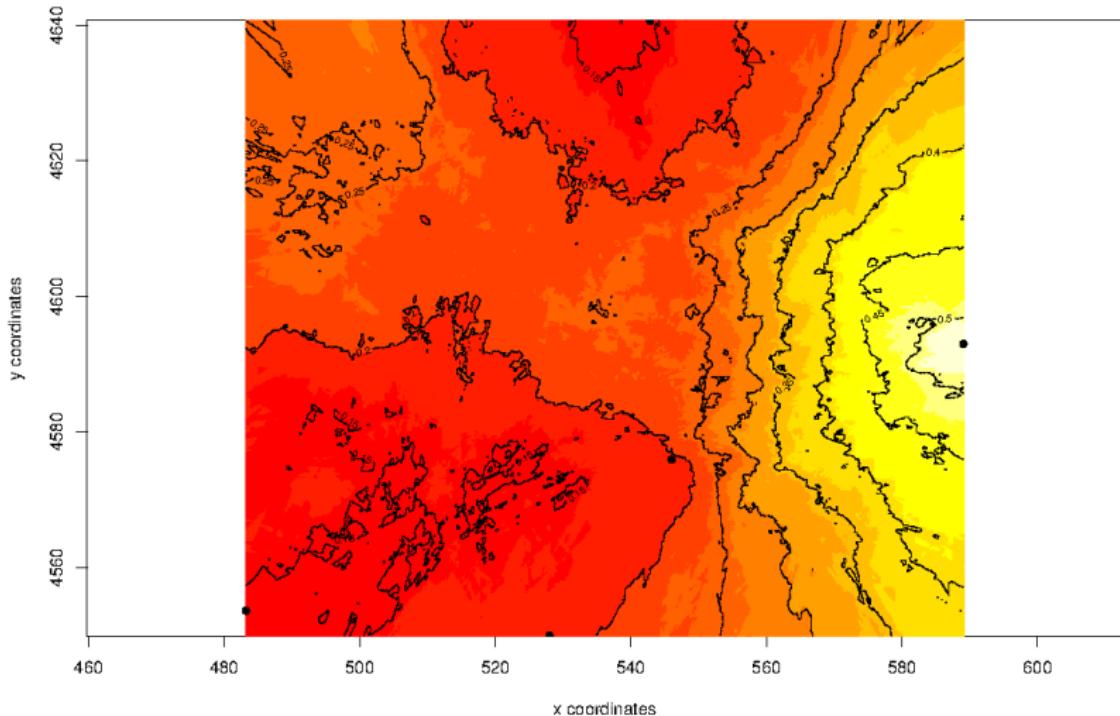
We will use figures and F_{ST} outputs only from the best run.

MCMC chain, number of clusters and their map



Map of posterior probability of belonging into cluster 1

Map of posterior probability to belong to cluster 1



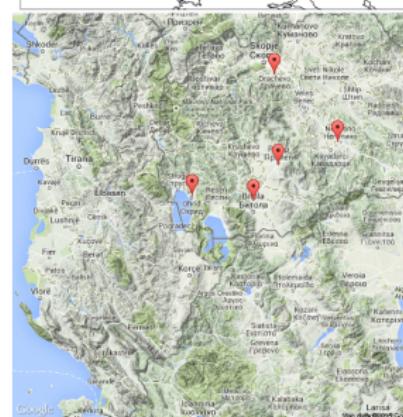
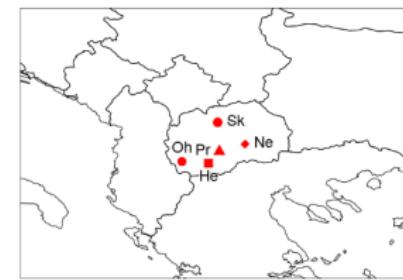
Very basic mapping in R

```
1 # Load libraries
2 library(sp)
3 library(rworldmap) # Basic world maps
4 library(TeachingDemos) # To be able to move text little bit
5 library(RgoogleMaps) # Google and OpenStreetMaps
6 # Plot basic map with state boundaries within selected range
7 plot(x=getMap(resolution="high"), xlim=c(19, 24), ylim=c(39, 44),
8      asp=1, lwd=1.5)
9 box() # Add frame around the map
10 # Plot location points
11 points(x=hauss.genpop@other$xy[["lon"]], y=hauss.genpop@other$xy
12   [ ["lat"]], pch=15:19, col="red", cex=4)
13 # Add text descriptions for points. Text is aside and with background
14 shadowtext(x=hauss.genpop@other$xy[["lon"]], y=hauss.genpop@other$xy
15   [ ["lat"]], labels=as.vector(hauss.genind@pop.names), col="black",
16   bg="white", theta=seq(pi/4, 2*pi, length.out=8), r=0.15,
17   pos=c(1, 3, 2, 4, 4), offset=0.75, cex=1.5)
```

Basic map and Google map

```

1 # Insert legend
2 legend(x="topright", inset=1/50,
3   legend=c("He", "Oh", "Pr", "Ne",
4   "Sk"), col="red", border="black",
5   pch=15:19, pt.cex=2, bty="o",
6   bg="lightgrey", box.lwd=1.5,
7   cex=1.5, title="Populations")
8 # Google map is produced into a
9 # file. Parametr markers contain
10 # data frame with coordinates and
11 # possibly with another information
12 hauss.gmap <- GetMap(center=c(lat=
13   41, lon=21), size=c(640, 640),
14   destfile="gmap.png", zoom=8,
15   markers=hauss.coord,
16   maptype="terrain")
```



OpenStreetMaps and datasets from mapproj

```
1 # Plot on OpenStreetMap - server is commonly overloaded and doesn't
2 # respond correctly - in fact, it rarely works...
3 GetMap.OSM(lonR=c(18, 24), latR=c(39, 44), scale=200000, destfile=
4   "osmmap.png", format="png", RETURNIMAGE=TRUE, GRAYSCALE=FALSE,
5   NEWMAP=TRUE, verbose=1)
6 # Plot on datasets from mapproj package
7 library(maps) # Various mapping tools (plotting, ...)
8 library(mapdata) # More detailed maps, but political boundaries often
9   # outdated, see http://cran.r-project.org/web/packages/mapdata/
10 library(mapproj) # Converts latitude/longitude into projected coordinates
11 # Plot a map, check parameters
12 # Check among others "projection" and ?mapproject for its details
13 map(database="worldHires", boundary=TRUE, interior=TRUE, fill=TRUE,
14   col="lightgrey", plot=TRUE, xlim=c(16, 27), ylim=c(37, 46))
15 # If you'd use projection, use - mapproject to convert also coordinates!
16 # See ?mapproject for details
17 points(x=hauss.genpop@other$xy[["lon"]], y=hauss.genpop@other$xy
18   [ ["lat"]], pch=15:19, col="red", cex=4)
```

Plotting on SHP files I

Get SHP files from <http://trapa.cz/sites/default/rcourse/macedonia.zip>

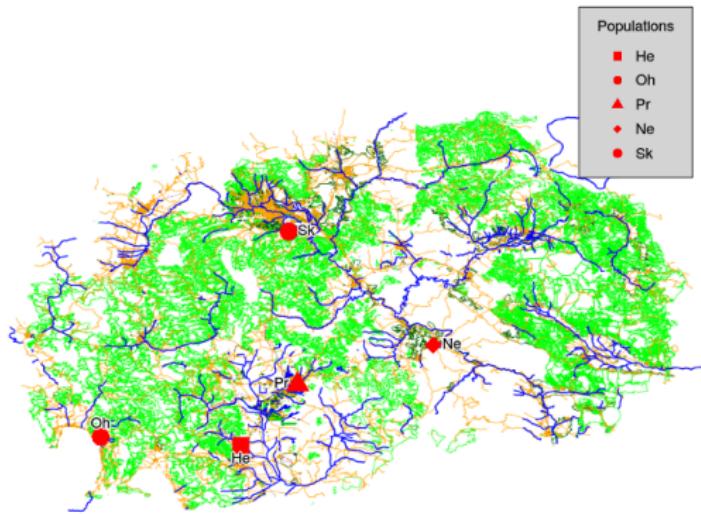
```
1 library(maptools)
2 # Load SHP file
3 # Data from http://download.geofabrik.de/europe/macedonia.html
4 # Directory has to contain also respective DBF and SHX files
5 # (same name, only different extension)
6 # There are several functions readShape* - select appropriate
7 # according to data stored in respective SHP file
8 macedonia_building <- readShapeLines(fn="macedonia_buildings.shp")
9 plot(macedonia_building)
10 macedonia_landuse <- readShapeLines(fn="macedonia_landuse.shp")
11 plot(macedonia_landuse)
12 macedonia_natural <- readShapeLines(fn="macedonia_natural.shp")
13 plot(macedonia_natural)
14 macedonia_railways <- readShapeLines(fn="macedonia_railways.shp")
15 plot(macedonia_railways)
16 macedonia_roads <- readShapeLines(fn="macedonia_roads.shp")
```

Plotting on SHP files II

```
1 plot(macedonia_roads)
2 macedonia_waterways <- readShapeLines(fn="macedonia_waterways.shp")
3 plot(macedonia_waterways)
4 # Plot all layers into single image
5 plot(macedonia_building)
6 plot(macedonia_landuse, add=TRUE, col="darkgreen", fill=TRUE)
7 plot(macedonia_natural, add=TRUE, col="green", fill=TRUE)
8 plot(macedonia_railways, add=TRUE, col="brown", lty="dotted")
9 plot(macedonia_roads, add=TRUE, col="orange")
10 plot(macedonia_waterways, add=TRUE, col="blue", lwd=2)
11 # Add sampling points
12 points(x=hauss.genpop@other$xy[["lon"]], y=hauss.genpop@other$xy[["lat"]], pch=15:19, col="red", cex=4)
13 # Add description of psampling points
14 shadowtext(x=hauss.genpop@other$xy[["lon"]], y=hauss.genpop@other$xy[["lat"]], labels=as.vector(hauss.genind@pop.names), col="black", bg="white", theta=seq(pi/4, 2*pi, length.out=8), r=0.15, pos=c(1, 3, 2, 4, 4), offset=0.75, cex=1.5)
```

Plotting on SHP files III

```
1 # Add legend  
2 legend(x="topright", inset=1/50, legend=c("He", "Oh", "Pr", "Ne",  
3 "Sk"), col="red", border="black", pch=15:19, pt.cex=2, bty="o",  
4 bg="lightgrey", box.lwd=1.5, cex=1.5, title="Populations")
```



Running Structure in parallel with R

- Population genetic software for Bayesian clustering **Structure** requires several runs for each K and various Ks. It can result in hundreds runs
- We should use power of modern multi-core CPUs
- **ParallelStructure** library can optimally distribute computations of independent Structure runs among CPU cores
- When using it, cite **Besnier & Glover 2013**
- I show slightly modified way from **The Molecular Ecologist**

```
1 # Prepare special directory and set working directory
2 setwd("~/dokumenty/fakulta/vyuka/r_mol_data/priklady/structure/")
3 library(ParallelStructure) # Load the library
4 # See Structure manual and function's documentation
5 # It takes more or less same parameter as for normal Structure
6 ?parallel_structure
7 # Author's recommend to run it without GUI and not on Windows
```

Preparing for ParallelStructure

Within working “structure” directory you need

- Subdirectory for results
- Text file describing jobs (“joblist”)
 - One row for one Structure run
 - Every line contains name of run, list of populations separated by commas (e.g. 1,2,3,4,5) — you don't have to use all populations in all runs
 - K for actual run
 - Length of burnin chain
 - Number of steps
 - Columns are separated by spaces (or TABs)
- Data input file
 - Make it as simple as possible — remove any unneeded columns
 - For population names use subsequent number from 1 to number of populations

Input files

Joblist file:

S02	1,2,3,4,5	1	500	10000
S06	1,2,3,4,5	2	500	10000
S08	1,2,3,4,5	2	500	10000
...

Input file:

		msta93	msta101	msta102	msta103	...	
H01	1	0	269	198	221	419	...
H01	1	0	269	198	223	419	...
H02	1	0	275	198	221	419	...
H02	1	0	283	198	223	419	...
...

Running ParallelStructure

```
1 parallel_structure(joblist="joblist.txt", n_cpu=7, structure_path=
2   "~/bin/", infile="hauss_stru.in", outpath="results/", numinds=47,
3   numloci=12, plot_output=1, label=1, popdata=1, popflag=1,
4   phenotypes=0, markernames=1, mapdist=0, onerowperind=0, phaseinfo=0,
5   extracol=0, missing=-9, ploidy=2, usepopinfo=0, revert_convert=1,
6   printqhat=1, locdata=0, recessivealleles=0, phased=0, noadmix=0,
7   linkage=0, locprior=0, inferalpha=1)
```

- Choose **n_cpu** according to your computer
- **structure_path** points to **directory** containing Structure binary
- **outpath** should aim to empty directory
- **plot_output=1** will produce plots for all runs
- Check all other settings according to Structure manual and your needs
- Get input file https://trapa.cz/sites/default/rcourse/hauss_stru.in and joblist <https://trapa.cz/sites/default/rcourse/joblist.txt>

ParallelStructure and Windows

- Authors do not recommend to run ParallelStructure on Windows...
- **parallel_structure()** is missing parallelizing library which is not available on Windows, instead try

```
1 # Install Rmpi library required by ParallelStructure for
2 # parallelisation on Windows
3 install.packages("Rmpi")
4 library(Rmpi)
5 # Instead of parallel_structure() use MPI_structure()
6 # with same arguments
7 MPI_structure(...) # Same arguments as on previous slide
```

- If this fails, look for some UNIX machine (Linux, Mac OS X, BSD, ...)...

Postprocess Structure results — select the best K

- Using Structure-sum-2009 R script by [Dorothee Ehrich](#)

```
1 # Load the script
2 source("structure-sum-2009.r")
3 # Create new directory with result files results_job_*_f and set
4 # working directory accordingly
5 setwd("/home/vojta/dokumenty/fakulta/vyuka/r_mol_data/priklady/
6 structure/structure_sum/")
```

- When using it, cite [Ehrich 2006](#). If you don't have the script, ask (it is not my so I don't want to post it on the net)
- Prepare list_k.txt containing on each line K and name of output file
- Get list K from https://trapa.cz/sites/default/rcourse/list_k.txt
- Joblist look like this:

```
2     results_job_S010_f
3     results_job_S011_f
...
...
```

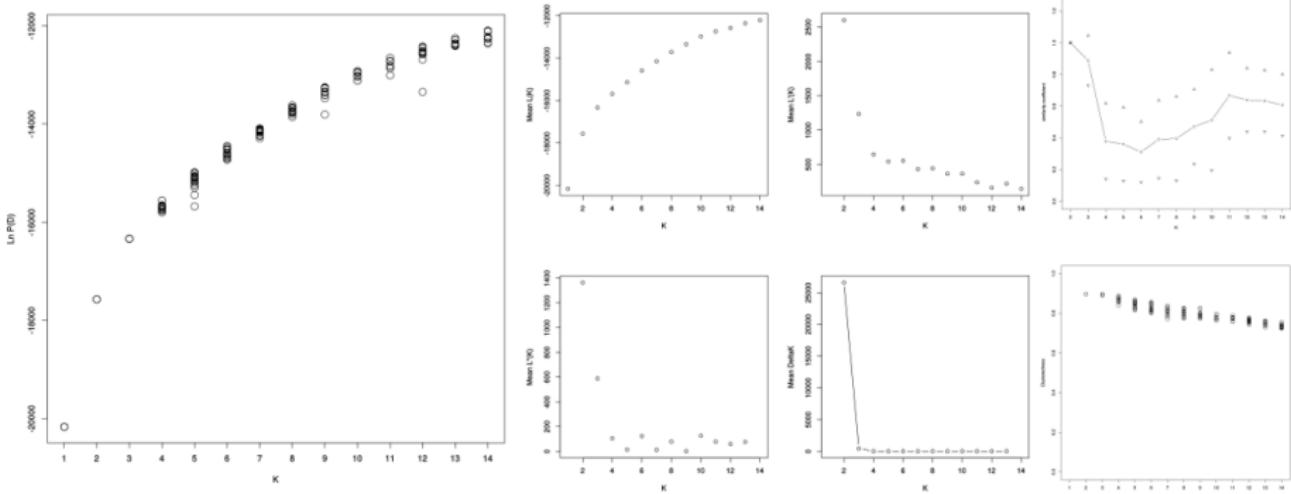
Run Structure-sum

```
1 # See documentation for details. Functions take as an argument
2 # list_k file and number of populations
3 Structure.table("list_k.txt", 5)
4 Structure.simil("list_k.txt", 5)
5 Structure.deltaK("list_k.txt", 5)
6 graphics.off() # Close graphics
7 Structure.cluster("list_k.txt", 5)
8 # Reordering ("alignment") of runs to get same clusters in same
9 # columns (prepare respective list_k files - one for each K)
10 Structure.order("list_k_02.txt", 5)
11 Structure.order("list_k_03.txt", 5)
12 Structure.order("list_k_04.txt", 5)
13 Structure.order("list_k_05.txt", 5)
14 Structure.order("list_k_06.txt", 5)
15 Structure.order("list_k_07.txt", 5)
16 # Continue with CLUMPP and distruct...
```

Details: <https://trapa.cz/en/structure-r-linux>

Outputs of Structure-sum — the best K is 2, may be 3 — stability of runs, good posterior probability

Results from different data set, not from our toy



Multiple sequence alignment — package muscle

```
1 # Muscle has its own functions read.fasta and write.fasta
2 library(muscle)
3 # Reads local file and writes R object or external file
4 # with alignment, you can use any command line argument
5 # available for muscle, see ?muscle and its manual
6 # Write result to external file
7 muscle(seqs="usflu.fasta", out="usflu.aln", quiet=FALSE)
8 # Write result to R object
9 usflu.muscle <- muscle(seqs="usflu.fasta", out=NULL, quiet=FALSE)
10 # Prints the alignment
11 print(usflu.muscle, from=1, to=usflu.muscle[["length"]])
12 class(usflu.muscle)
```

- Package **muscle** contains MUSCLE binary
- Read **MUSCLE documentation** to get correct parameters (if you are not fine with defaults)

Multiple sequence alignment

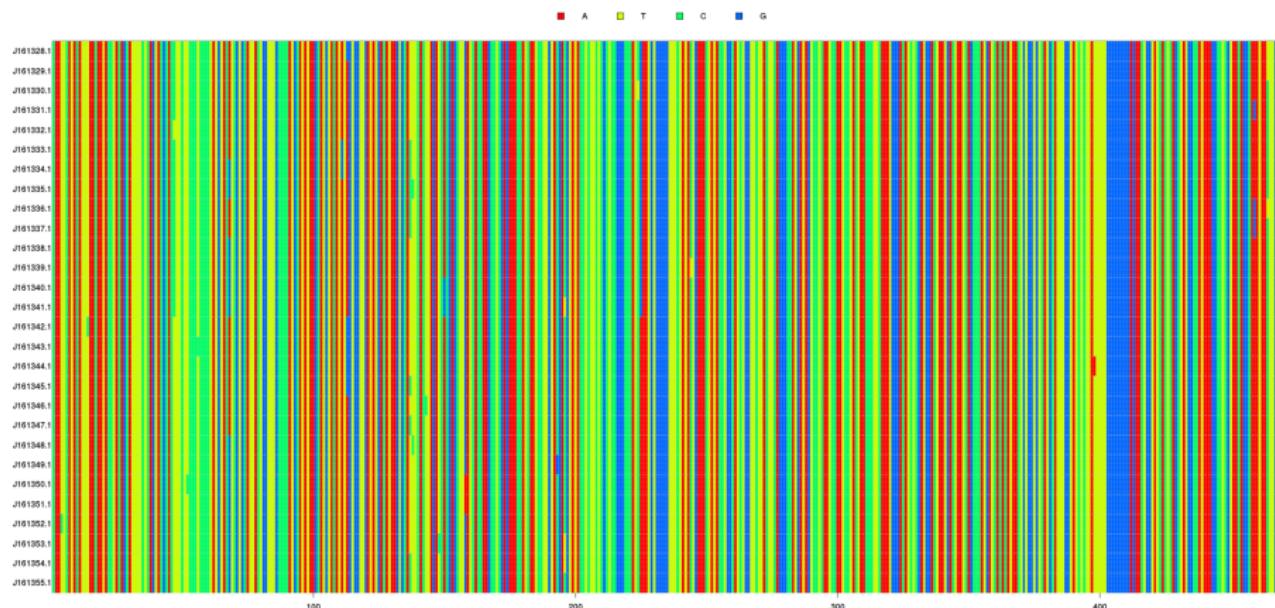
- MUSCLE is available in package **muscle** (above) and in package **ape**
 - first one reads FASTA files and writes external FASTA files or R objects of class “muscle”. The latter reads and writes object of class “DNAbin”
- ape also contains functions to use Clustal and T-Coffee — both read and write DNAbin
- MAFFT is available from (same author) in packages **ips** and **phyloch**
 - both read and write DNAbin

```
1 library(colorspace) # Libraries needed by phyloch/ips
2 library(XML)
3 library(phyloch) # Alignment with mafft, you can also try package ips
4 # Requires mafft binary on specific location
5 # you might be needed to copy it or make symlink or so
6 # read ?mafft and mafft's documentation
7 meles.mafft <- mafft(x=meles.dna, method="localpair", maxiterate=100)
8 meles.mafft
```

Clustal, MUSCLE and T-Coffee from ape

```
1 class(meles.mafft)
2 # read ?clustal and documentation of Clustal, Muscle and T-Coffee
3 # when using them to set correct parameters
4 clustal(x=meles.dna, pw.gapopen=10, pw.gapext=0.1, gapopen=10,
5   gapext=0.2, quiet=FALSE, original.ordering=TRUE)
6 meles.muscle <- muscle(x=meles.dna, exec="muscle",
7   quiet=FALSE, original.ordering=TRUE)
8 meles.muscle
9 class(meles.muscle)
10 # Plot the alignment - you can select which bases to plot
11 # and/or modify colours
12 image(x=meles.muscle, c("a", "t", "c", "g", "n"), col=rainbow(5))
13 # Add grey dotted grid
14 grid(nx=ncol(meles.muscle), ny=nrow(meles.muscle), col="lightgrey")
15 # Remove gaps from alignment - destroy it
16 meles.nogaps <- del.gaps(meles.muscle)
```

Multiple sequence alignment with MUSCLE



Read and write tree and drop tips

```
1 # Read trees in NEWICK format - single or multiple tree(s)
2 oxalis.trees <- read.tree
3 ("http://trapa.cz/sites/default/rcourse/oxalis.nwk")
4 summary(oxalis.trees)
5 length(oxalis.trees)
6 names(oxalis.trees)
7 # Export trees in NEWICK format
8 write.tree(phy=oxalis.trees, file="trees.nwk")
9 # Drop a tip from multiPhylo
10 plot.multiPhylo(x=oxalis.trees)
11 oxalis.trees.drop <- lapply(X=oxalis.trees, FUN=drop.tip, "TaxaOut")
12 class(oxalis.trees.drop) <- "multiPhylo"
13 plot.multiPhylo(x=oxalis.trees.drop)
14 # Drop a tip from single tree
15 plot.phylo(hauss.nj)
16 hauss.nj.drop <- drop.tip(phy=hauss.nj, tip=47)
17 plot.phylo(hauss.nj.drop)
```

Extract clades from trees and drop extinct tips

```
1 # Interactively extract tree
2 # Plot source tree
3 plot.phylo(hauss.nj)
4 # Select clade to extract by clicking on it
5 hauss.nj.extracted <- extract.clade(phy=hauss.nj, interactive=TRUE)
6 # See new extracted tree
7 plot.phylo(hauss.nj.extracted)
8 # Non-interactively extract tree
9 hauss.nj.extracted <- extract.clade(phy=hauss.nj, node=60,
10   interactive=FALSE)
11 # See new extracted tree
12 plot.phylo(hauss.nj.extracted)
13 # Drop "extinct" tips - those who don't reach end the tree
14 # tolerance is respective to the used metrics
15 plot.phylo(hauss.nj)
16 axisPhylo()
17 hauss.nj.fossil <- drop.fossil(phy=hauss.nj, tol=0.25)
18 plot.phylo(hauss.nj.fossil)
```

Join two trees, rotate tree

```
1 # Bind two trees into one
2 hauss.nj.bind <- bind.tree(x=hauss.nj.fossil, y=hauss.nj.extracted,
3   where="root", position=0, interactive=FALSE)
4 plot.phylo(hauss.nj.bind)
5 # Bind two trees interactively
6 # Plot tree receiving the new one
7 plot.phylo(hauss.nj.fossil)
8 # Select where to bind new tree to
9 hauss.nj.bind <- bind.tree(x=hauss.nj.fossil, y=hauss.nj.extracted,
10   interactive=TRUE)
11 plot.phylo(hauss.nj.bind)
12 # Rotate tree
13 plot.phylo(hauss.nj)
14 hauss.nj.rotated <- rotate(phy=hauss.nj, node="70")
15 plot.phylo(hauss.nj.rotated)
```

Ladderize and (un)root the tree

```
1 # Ladderize the tree
2 plot.phylo(hauss.nj)
3 hauss.nj.ladderized <- ladderize(hauss.nj)
4 plot.phylo(hauss.nj.ladderized)
5 # Root the tree
6 plot.phylo(hauss.nj)
7 print.phylo(hauss.nj)
8 hauss.nj.rooted <- root(phy=hauss.nj, outgroup="10")
9 print.phylo(hauss.nj.rooted)
10 plot.phylo(hauss.nj.rooted)
11 # Root the tree interactive
12 plot.phylo(hauss.nj)
13 hauss.nj.rooted <- root(phy=hauss.nj, interactive=TRUE)
14 plot.phylo(hauss.nj.rooted)
15 # unroot the tree
16 unroot()
17 # Check if it is rooted
18 is.rooted()
```

Check tree and compute branch lengths and times

```
1 # Check if the tree is ultrametric - is variance of distances
2 # of all tips to node 0? It is required for some analysis
3 is.ultrametric()
4 # Make tree ultrametric
5 chronos()
6 ?chronos # Check it for mode how to calculate the lengths
7 # chronos has more uses - it is mainly used for dating
8 # Compute branch lengths for trees without branch lengths
9 ?compute.brlen # Check it for mode how to calculate the lengths
10 compute.brlen()
11 # Computes the branch lengths of a tree giving its branching
12 # times (aka node ages or heights)
13 compute.brtimes()
14 ?compute.brtimes # Check it for mode how to calculate the lengths
```

Class “multiPhylo” is just a list of “phylo” objects to store multiple trees — you can perform most of analysis on it as on phylo, commonly using **lapply** function (afterwards use `class(x) <- "multiPhylo"` to ensure other functions will see it as multiPhylo object).

Topographical distances among trees I

We have plenty of trees. How much are their topologies different?

```
1 library(gplots)
2 library(corrplot)
3 library(phytools)
4 # Prepare matrix for distances
5 oxalis.trees.d <- matrix(nrow=length(oxalis.trees),
6   ncol=length(oxalis.trees))
7 # Calculate pairwise topographic distances
8 for (i in 1:length(oxalis.trees)) {
9   for (j in i:length(oxalis.trees)) {
10     print(c(i,j))
11     oxalis.trees.d[i,j] <- dist.topo(oxalis.trees[[i]],
12       oxalis.trees[[j]])
13   }
14 }
15 # Basic information about the distance matrix
16 dim(oxalis.trees.d)
17 head.matrix(oxalis.trees.d)
```

Topographical distances among trees II

Postprocess the matrix and plot it. There are several methods for calculating distance matrices among the trees — some take branch lengths into account, some only topology.

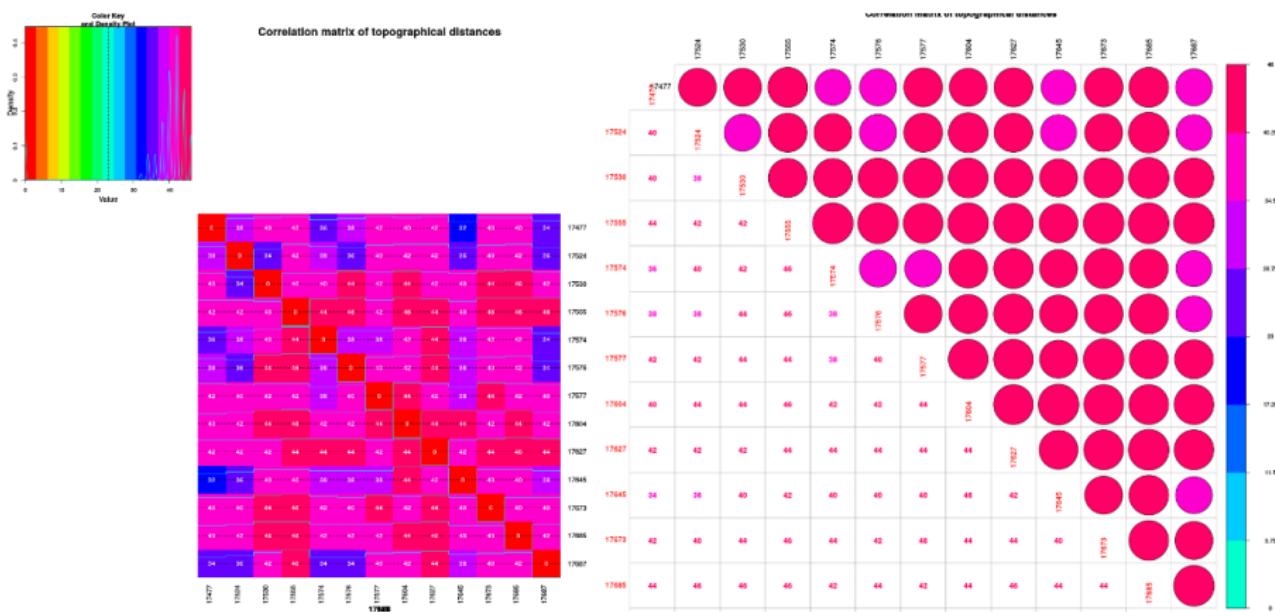
```
1 # Add names of columns and rows
2 colnames(oxalis.trees.d) <- names(oxalis.trees)
3 rownames(oxalis.trees.d) <- names(oxalis.trees)
4 # Make matrix symmetric
5 oxalis.trees.d[lower.tri(oxalis.trees.d)] =
6 t(oxalis.trees.d)[lower.tri(oxalis.trees.d)]
7 # Create heatmaps using heatmap.2 function from gplots package
8 heatmap.2(x=oxalis.trees.d, Rowv=FALSE, Colv="Rowv", dendrogram="none",
9 symm=TRUE, scale="none", na.rm=TRUE, revC=FALSE, col=rainbow(15),
10 cellnote=oxalis.trees.d, notecex=1, notecol="white", trace="row",
11 linecol="black", labRow=names(oxalis.trees),
12 labCol=names(oxalis.trees), key=TRUE, keyszie=2,
13 density.info="density", symkey=FALSE, main="Correlation matrix of
14 topographical distances", xlab=names(oxalis.trees),
15 ylab=names(oxalis.trees))
```

Topographical distances among trees III

Calculate Robinsons-Foulds distance (allow compare trees created by different algorithms) matrix among trees and plot it.

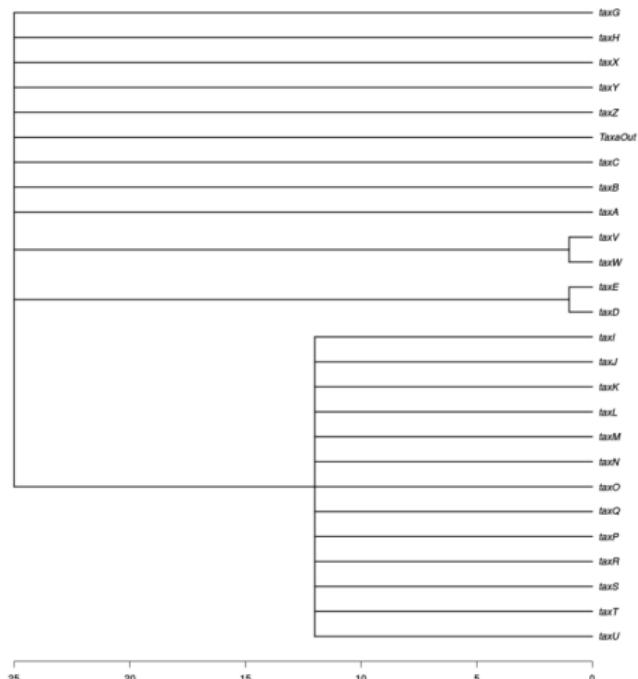
```
1 # Robinsons-Foulds distance
2 oxalis.trees.d.rf <- multiRF(oxalis.trees)
3 # Add names of columns and rows
4 colnames(oxalis.trees.d.rf) <- names(oxalis.trees)
5 rownames(oxalis.trees.d.rf) <- names(oxalis.trees)
6 # Create heatmap using corrplot function from corrplot package
7 corrplot(corr=oxalis.trees.d.rf, method="circle", type="upper",
8   col=rainbow(15), title="Correlation matrix of topographical
9   distances", is.corr=FALSE, diag=FALSE, outline=TRUE,
10  order="alphabet", tl.pos="lt", tl.col="black")
11 corrplot(corr=oxalis.trees.d.rf, method="number", type="lower",
12  add=TRUE, col=rainbow(15), title="Correlation matrix of
13  topographical distances", is.corr=FALSE, diag=FALSE,
14  outline=FALSE, order="alphabet", tl.pos="ld", cl.pos="n")
```

Topographical distances among trees IV — the matrices

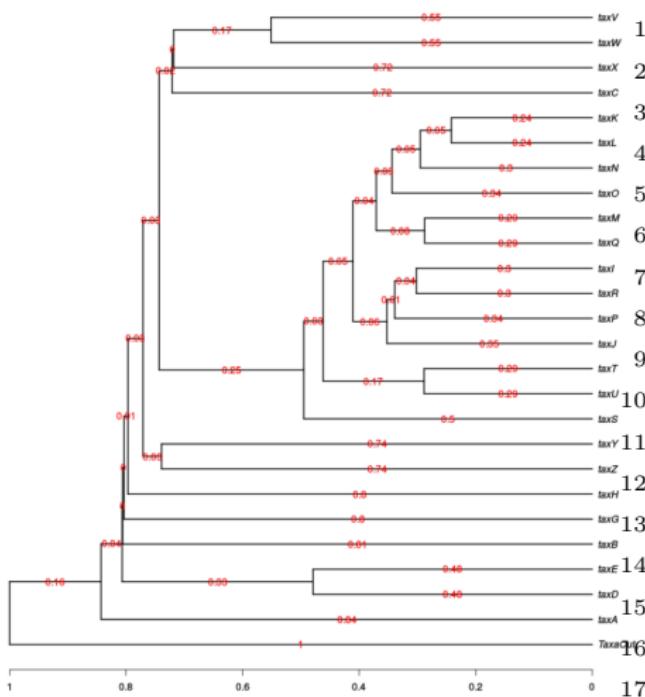


Consensus tree

```
1 # Root all trees
2 oxalis.trees.rooted <- lapply
3   (X=oxalis.trees, FUN=root,
4   "TaxaOut")
5 class(oxalis.trees.rooted) <-
6   "multiPhylo"
7 # Consensus tree (50 \% rule)
8 oxalis.tree.con <- consensus
9   (oxalis.trees.rooted, p=0.5,
10  check.labels=TRUE)
11 print.phylo(oxalis.tree.con)
12 # Plot the tree
13 plot.phylo(oxalis.tree.con,
14   edge.width=2, label.offset=0.3)
15 axisPhylo(side=1)
16 # What a nice tree... :-P
```



Species tree — all trees must be ultrametric

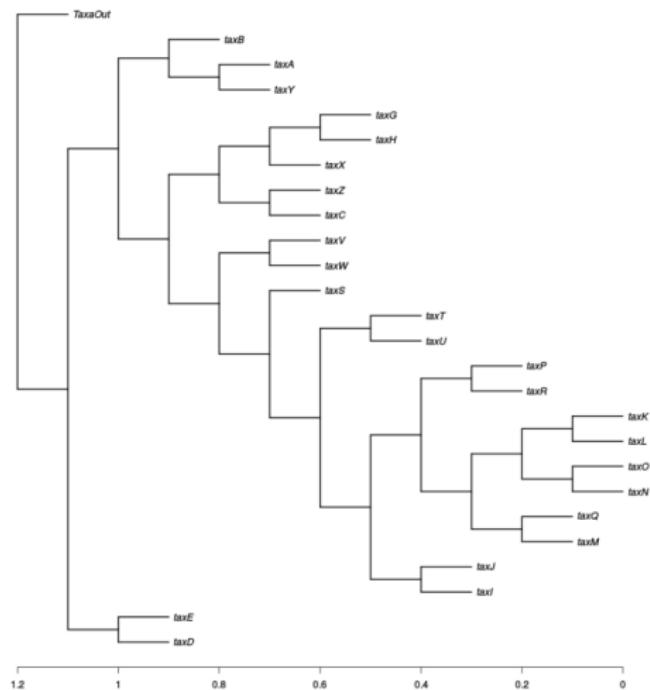


```
# Chronos scale trees
oxalis.trees.ultra <- lapply
(X=oxalis.trees.rooted,
FUN=chronos, model="correlated")
class(oxalis.trees.ultra) <-
"multiPhylo"
# Mean distances
oxalis.tree.sp.mean <- speciesTree
(oxalis.trees.ultra, mean)
# Plot the tree
plot.phylo(oxalis.tree.sp.mean,
edge.width=2, label.offset=0.01)
edgelabels(text=round(oxalis.
tree.sp.mean[["edge.length"]]),
digits=2), frame="none",
col="red", bg="none")
axisPhylo(side=1)
```

Parsimony super tree

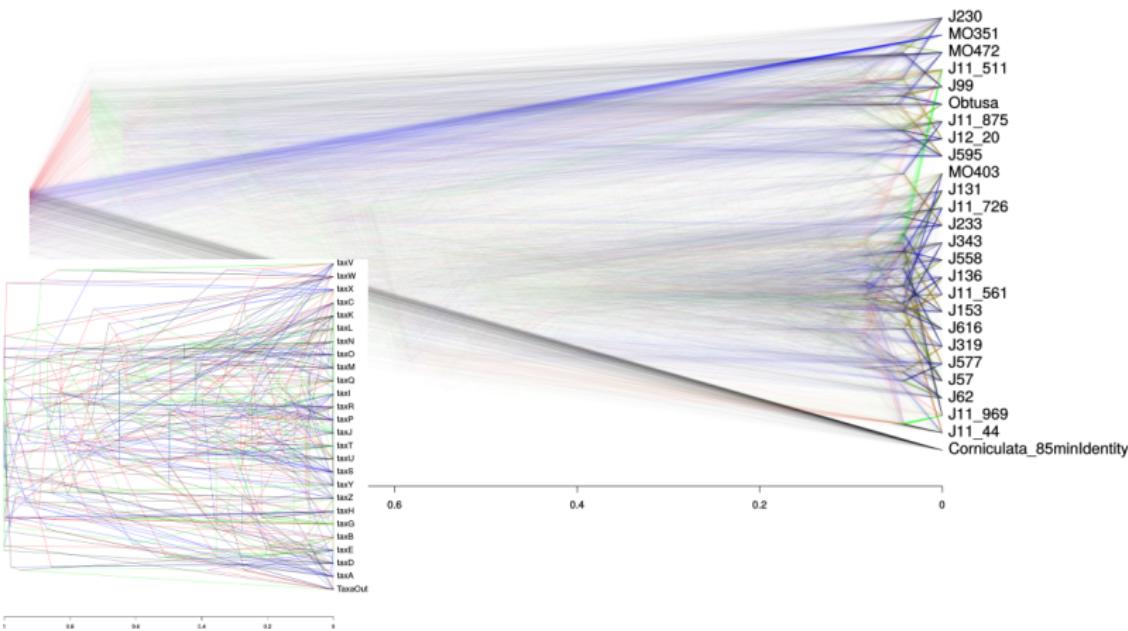
```
1 library(phangorn)
2 oxalis.tree.sp <- superTree(tree=
3   oxalis.trees.rooted, method=
4     "optim.parsimony", rooted=TRUE)
5 print.phylo(oxalis.tree.sp)
6 plot.phylo(oxalis.tree.sp,
7   edge.width=2, label.offset=0.01)
8 axisPhylo(side=1)
```

Currently, there is no package for coalescence based approach for creation of species tree from multiple gene trees. This function was in package `phybase`, but it is abandoned and doesn't install in R 3 without hacks



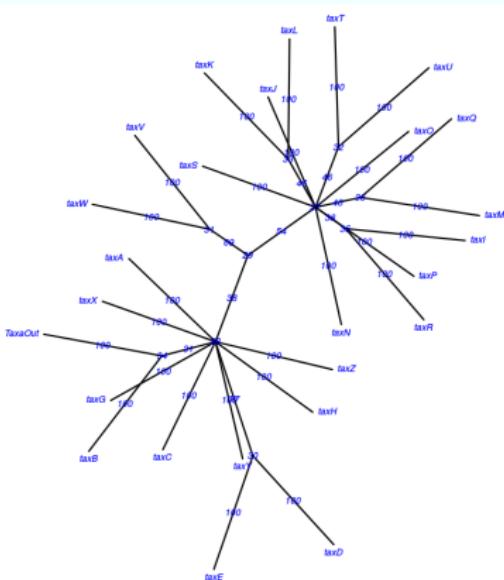
Density tree

```
1 densiTree(x=oxalis.trees.ultra, type="cladogram", alpha=0.5,
2   consensus=oxalis.tree.sp.mean, scaleX=TRUE, col=c("black",
3   "green", "blue", "red"), cex=1.5)
```



Networks

```
1 oxalis.tree.net <- consensusNet
2 (oxalis.trees.rooted, prob=0.25)
```



```
1 plot.network(x=oxalis.tree.net,
2 planar=FALSE, type="2D",
3 use.edge.length=TRUE,
4 show.tip.label=TRUE,
5 show.edge.label=TRUE,
6 show.node.label=TRUE,
7 show.nodes=TRUE,
8 edge.color="black",
9 tip.color="blue")
10 plot.network(x=oxalis.tree.net,
11 planar=FALSE, type="3D",
12 use.edge.length=TRUE,
13 show.tip.label=TRUE,
14 show.edge.label=TRUE,
15 show.node.label=TRUE,
16 show.nodes=TRUE, edge.color=
17 "black", tip.color="blue")
```

Kronoviz — see all trees on same scale

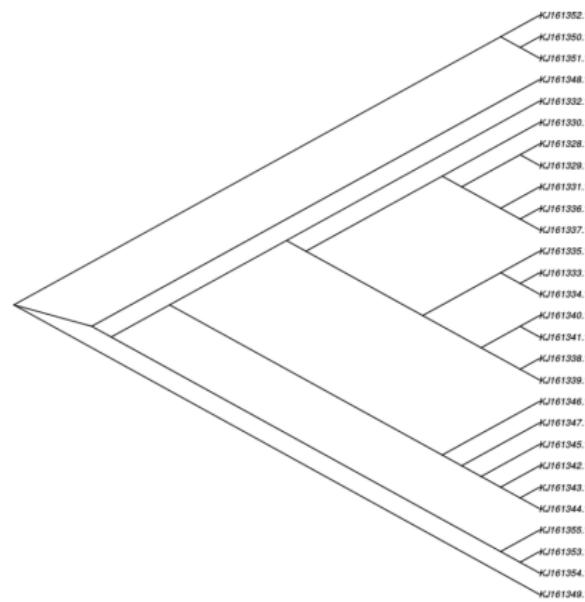


```
1 kronoviz(x=oxalis.trees.rooted,  
2   layout=length(oxalis.trees.  
3   rooted), horiz=TRUE)  
4 # Close graphical device to  
5 # cancel division of plotting  
6 # device  
7 dev.off()
```

The plot can be very long and it can be hard to see details. But one can get impression if all trees are more or less in same scale (have comparable length) or not.

Maximum parsimony

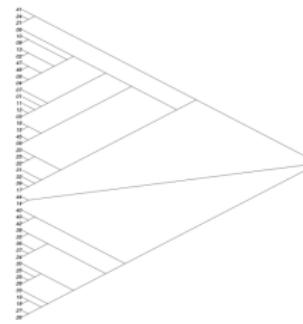
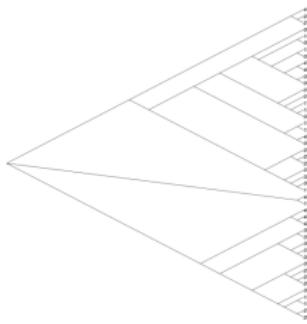
```
1 # Conversion to phyDat
2 meles.phydat <-
3   as.phyDat(meles.dna)
4 # Prepare starting tree
5 meles.tre.ini <- nj(dist.
6   dna(x=meles.dna, model="raw"))
7 # Parsimony: ?parsimony
8 # Maximum parsimony score
9 parsimony(tree=meles.tre.ini,
10   data=meles.phydat)
11 # Optimisation
12 # Maximum parsimony tree
13 meles.tre.pars <- optim.
14   parsimony(tree=meles.tre.
15   ini, data=meles.phydat)
16 # Draw a tree
17 plot.phylo(x=meles.tre.pars,
18   type="clad", edge.width=2)
```

Maximum-parsimony tree of *Meles*

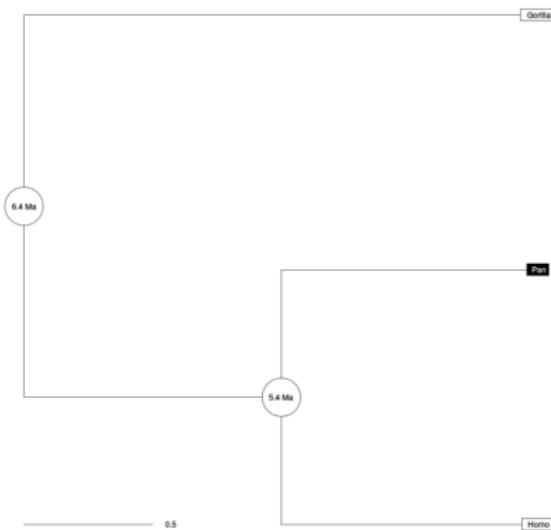
Change orientation of plots

plot.phylo() has plenty of possibilities to influence — check **?plot.phylo** and **?par**

```
1 ?plot.phylo # check it for various possibilities what to influence
2 par(mfrow=c(1, 2)) # Plot two plots in one row
3 plot.phylo(x=hauss.nj, type="cladogram", use.edge.length=FALSE,
4   direction="rightwards")
5 plot.phylo(x=hauss.nj, type="cladogram", use.edge.length=FALSE,
6   direction="leftwards")
7 dev.off() # Close graphical device to cancel par() settings
```



Highlighted labels



```
1 # Nice tiplabels and
2 # highlighted tiplabel
3 trape <- read.tree(text=
4   "((Homo, Pan), Gorilla);")
5 plot.phylo(x=trape,
6   show.tip.label=FALSE)
7 tiplabels(trape[["tip.label"]],
8   bg=c("white", "black",
9   "white"), col=c("black",
10  "white", "black"))
11 nodelabels(text=c("6.4 Ma",
12  "5.4 Ma"), frame="circle",
13  bg="white")
14 add.scale.bar()
```

Phylogenetic independent contrast

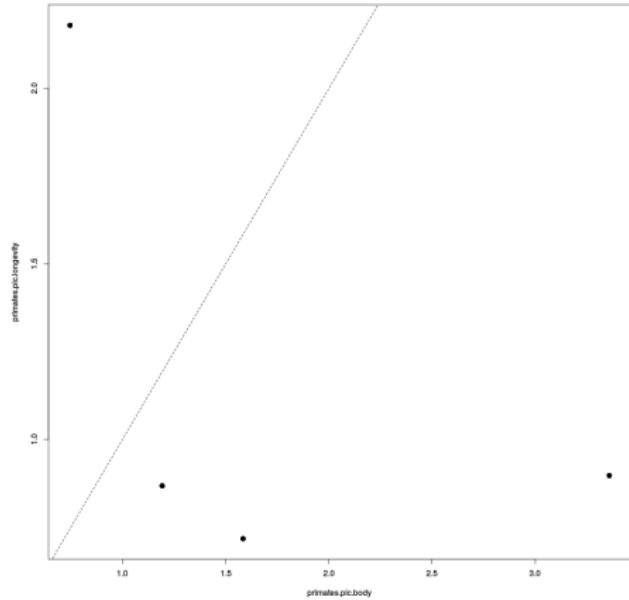
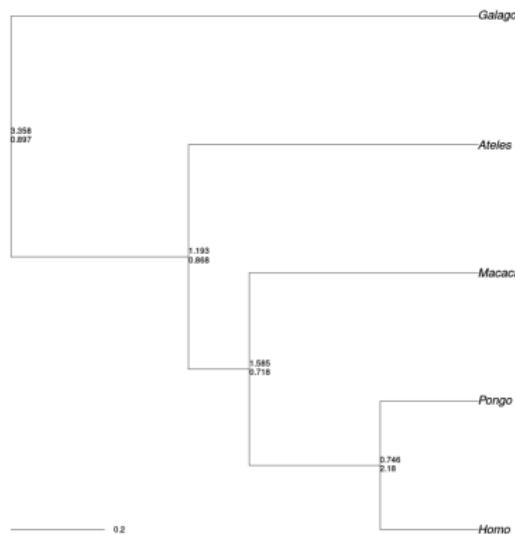
- When analysing comparative data takes phylogeny into account
- If we assume that a continuous trait evolves randomly in any direction (i.e., the Brownian motion model), then the “contrast” between two species is expected to have a normal distribution with mean zero, and variance proportional to the time since divergence

```
1 # Prepare the data # Body mass of primates
2 primates.body <- c(4.09434, 3.61092, 2.37024, 2.02815, 1.46968)
3 # Longevity of primates
4 primates.longevity <- c(4.74493, 3.3322, 3.3673, 2.89037, 2.30259)
5 # Add names to the values
6 names(primates.body) <- names(primates.longevity) <- c("Homo", "Pongo",
7 "Macaca", "Ateles", "Galago")
8 # Create a tree in Newick format
9 primates.tree <- read.tree(text="(((Homo:0.21, Pongo:0.21):0.28,
10 Macaca:0.49):0.13, Ateles:0.62):0.38, Galago:1.00);")
11 plot.phylo(primates.tree)
```

PIC and its plotting

```
1 primates.pic.body <- pic(x=primates.body, phy=primates.tree,
2   scaled=TRUE, var.contrasts=FALSE, rescaled.tree=FALSE)
3 primates.pic.longevity <- pic(x=primates.longevity, phy=primates.tree,
4   scaled=TRUE, var.contrasts=FALSE, rescaled.tree=FALSE)
5 # Plot a tree with PIC values
6 plot.phylo(x=primates.tree, lwd=2, cex=1.5, offset=0.2)
7 nodelabels(round(primates.pic.body, digits=3), adj=c(0, -0.5),
8   frame="none")
9 nodelabels(round(primates.pic.longevity, digits=3), adj=c(0, 1),
10  frame="none")
11 add.scale.bar()
12 # Plot PIC
13 plot(x=primates.pic.body, y=primates.pic.longevity, pch=16, cex=1.5)
14 abline(a=0, b=1, lty=2) # x=y line
15 # correlation coefficient of both PICs
16 cor(x=primates.pic.body, y=primates.pic.longevity, method="pearson")
17 [1] -0.5179156
```

Plot of PIC (on the tree)



Test it

```
1 lm(formula=primates.pic.longevity~primates.pic.body)
2 Coefficients:
3   (Intercept) primates.pic.body
4           1.6957          -0.3081
5 # Because PICs have expected mean zero - such linear regressions
6 # should be done through the origin (the intercept is set to zero)
7 lm(formula=primates.pic.longevity~primates.pic.body-1)
8 Coefficients:
9 primates.pic.body
10            0.4319
11 # Permutation procedure to test PIC
12 lmorigin(formula=primates.pic.longevity~primates.pic.body, nperm=1000)
13 Regression through the origin
14 Permutation method = raw data
15 Coefficients and parametric test results
16             Coefficient Std_error t-value Pr(>|t|) 
17 primates.pic.body     0.43193  0.28649  1.5077  0.2288
18 F-statistic: 2.273067 on 1 and 3 DF:
19 permutational p-value: 0.2377622
```

Phylogenetic autocorrelation

- Autocorrelation coefficient to quantify whether the distribution of a trait among a set of species is affected or not by their phylogenetic relationships
- In the absence of phylogenetic autocorrelation, the mean expected value of I and its variance are known - it is thus possible to test the null hypothesis of the absence of dependence among observations

```
1 # Let's choose weights as  $w_{ij} = 1/d_{ij}$ , where the  $d$ 's is the distances
2 # measured on the tree - cophenetic() calculates cophenetic distances
3 # can be just cophenetic(primates.tree) or some other transformation
4 primates.weights <- 1/cophenetic(primates.tree)
5 primates.weights # See it
6 class(primates.weights)
7 diag(primates.weights) <- 0 # Set diagonal to 0
```

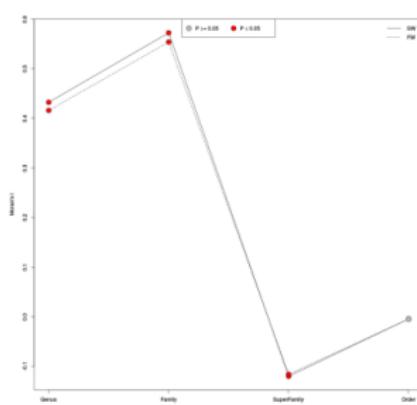
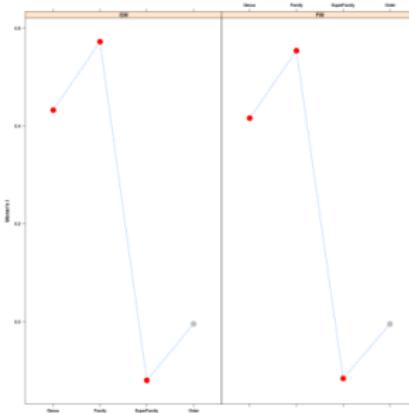
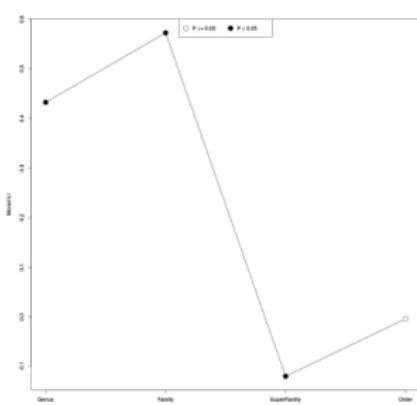
Testing of Moran's I

```
1 # Calculate Moran's I
2 # Slightly significant positive phylogenetic correlation among body mass
3 Moran.I(x=primates.body, weight=primates.weights,
4   alternative="greater")
5 # Positive, but non-significant
6 Moran.I(x=primates.longevity, weight=primates.weights,
7   alternative="greater")
8 # Test of Moran's with randomisation procedure
9 # Body is significant - nonrandom, longevity not (random)
10 gearymoran(bilis=primates.weights, X=data.frame(primates.body,
11   primates.longevity), nrepet=1000)
12 # Test of Abouheif designed to detect phylogenetic autocorrelation in a
13 # quantitative trait - in fact Moran's I test using a particular
14 # phylogenetic proximity between tips
15 library(adephylo)
16 abouheif.moran(x=cbind(primates.body, primates.longevity),
17   W=primates.weights, method="oriAbouheif", a=1, nrepet=1000,
18   alter="greater")
```

Correlogram to visualize results of phylogenetic autocorrelative analysis

```
1 data(carnivora) # Loads training data set
2 head(carnivora) # Look at the data
3 # Calculate the correlogram
4 carnivora.correlogram <- correlogram.formula
5   (formula=SW~Order/SuperFamily/Family/Genus, data=carnivora)
6 carnivora.correlogram # See results
7 # Calculate the correlogram - test for both body masses
8 carnivora.correlogram2 <- correlogram.formula
9   (formula=SW+FW~Order/SuperFamily/Family/Genus, data=carnivora)
10 carnivora.correlogram2 # See results
11 plot.correlogram(x=carnivora.correlogram, legend=TRUE,
12   test.level=0.05, col=c("white", "black")) # Plot it
13 # Plot it - test for both body masses - two or one graph(s)
14 plot.correlogramList(x=carnivora.correlogram2, lattice=TRUE,
15   legend=TRUE, test.level=0.05)
16 plot.correlogramList(x=carnivora.correlogram2, lattice=FALSE,
17   legend=TRUE, test.level=0.05)
```

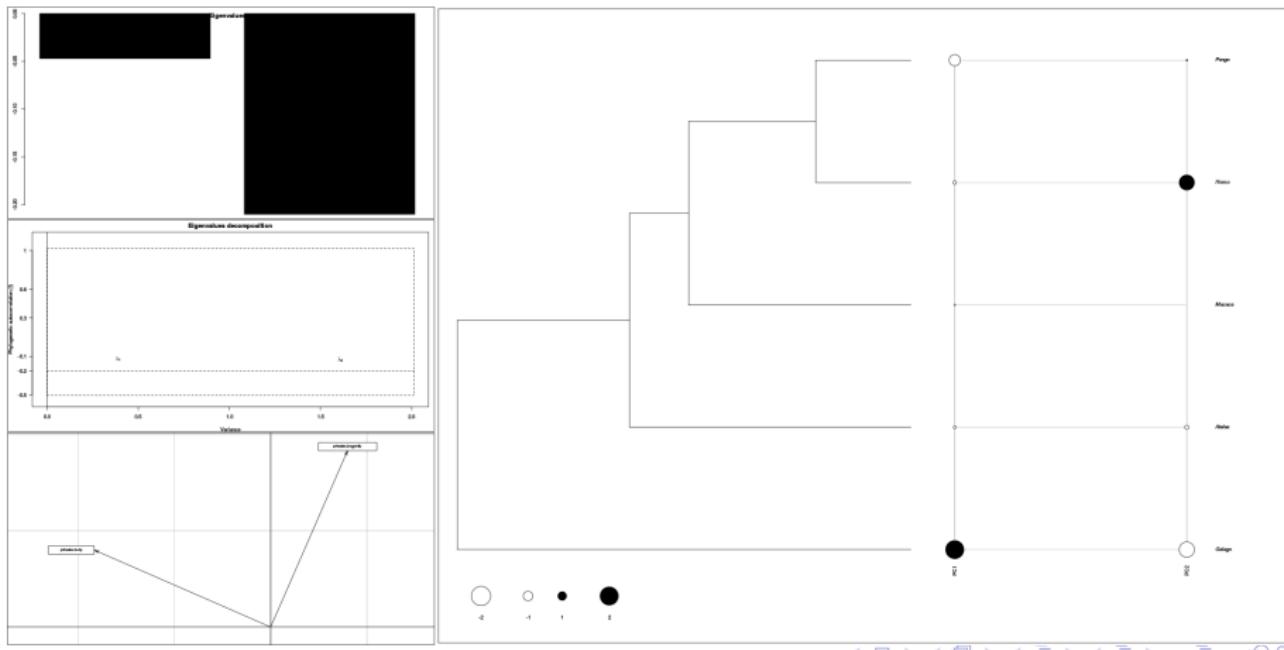
Correlograms of SW and SW+FW (in one or two graphs) depending on taxonomical level with marked significance



Phylogenetic principal component analysis

```
1 # Library needed to create phylo4d object required by pPCA
2 library(phylobase)
3 # Calculate pPCA
4 primates.pppca <- pppca(x=phylo4d(x=primates.tree, cbind(
5   primates.body, primates.longevity)), method="patristic",
6   a=1, center=TRUE, scale=TRUE, scannnf=TRUE, nffposi=1, nfnega=0)
7 # Print results
8 print.pppca(primates.pppca)
9 # See summary information
10 summary.pppca(primates.pppca)
11 # See PCA scores for variables on phylogenetic tree
12 scatter.pppca(primates.pppca)
13 # See decomposition of pPCA eigenvalues
14 screeplot.pppca(primates.pppca)
15 # Plot pPCA results - global vs. local structure, decomposition
16 # of pPCA eigenvalues, PCA plot of variables and PCA scores
17 # for variables on phylogenetic tree
18 plot.pppca(primates.pppca)
```

Plot pPCA results - global vs. local structure, decomposition of pPCA eigenvalues, PCA plot of variables and PCA scores for variables on phylogenetic tree



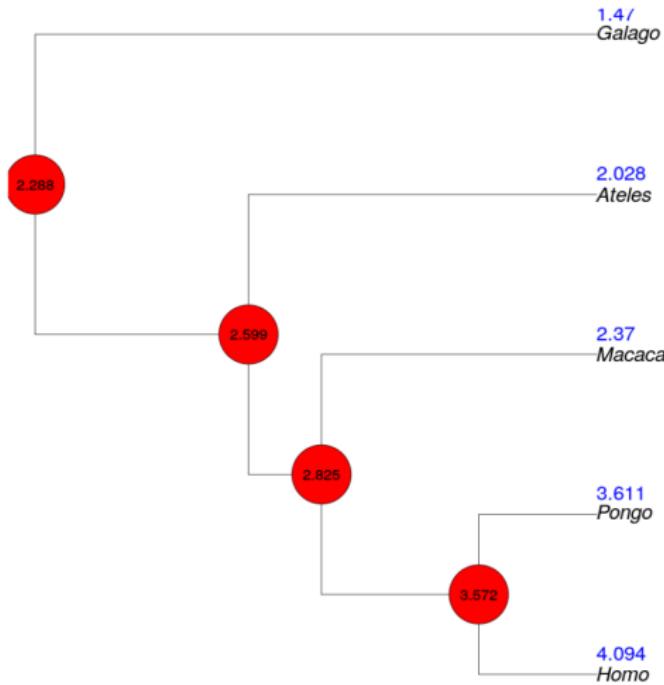
Ancestral state reconstruction

By default **ace** performs estimation for continuous characters assuming a Brownian motion model fit by maximum likelihood

```
1 # See ?ace for possible settings
2 primates.body.ace <- ace(x=primates.body, phy=primates.tree,
3   type="continuous", method="REML",
4   corStruct=corBrownian(value=1, phy=primates.tree))
5 # See result - reconstructions are in $ace
6 # To be plotted on nodes - 1st column are node numbers
7 primates.body.ace
8 # Plot it
9 plot(primates.tree, lwd=2, cex=2)
10 tiplabels(round(primates.body, digits=3), adj=c(0, -1),
11   frame="none", col="blue", cex=2)
12 nodelabels(round(primates.body.ace$ace, digits=3),
13   frame="circle", bg="red", cex=1.5)
```

Other implementations are available in packages geiger (functions fitContinuousMCMC and fitDiscrete), phangorn and more in ape (MPR)

Ancestral state reconstructions of primates body weights



Citations

- To correctly cite R launch **citation()** and see information there — it is slightly different for every version of R
- Cite used packages — launch **citation("PackageName")** — if this information is missing, go to its manual page and/or homepage and find the information there
- Packages/functions commonly provide various methods to calculate desired task — check function's help page (**?FunctionName**) and find references there and cite them accordingly

Where to look for the help I

- R phylogeny mailing list
<https://stat.ethz.ch/mailman/listinfo/r-sig-phylo>
- R genetics mailing list
<https://stat.ethz.ch/mailman/listinfo/r-sig-genetics>
- Bioconductor mailing list
<https://stat.ethz.ch/mailman/listinfo/bioconductor> and Bioconductor help pages <http://master.bioconductor.org/help/>
- R phylo wiki http://www.r-phylo.org/wiki/Main_Page
- R phylogenetics at CRAN
<http://cran.r-project.org/web/views/Phylogenetics.html>
- Little Book of R for Bioinformatics
<http://a-little-book-of-r-for-bioinformatics.readthedocs.org/en/latest/>

Where to look for the help II

- Adegenet web <http://adegenet.r-forge.r-project.org/> and help forum <https://lists.r-forge.r-project.org/cgi-bin/mailman/listinfo/adegenet-forum>
- APE home page <http://ape-package.ird.fr/>
- Phytools <http://phytools.org/> and its blog <http://blog.phytools.org/>
- Poppr documentation and forum
<http://grunwaldlab.cgrb.oregonstate.edu/poppr-r-package-population-genetics/manual-documentation-news-and-changes>
- ade4 home page and documentation
<http://pbil.univ-lyon1.fr/ade4/home.php?lang=eng>
- Phangorn resources
<http://cran.r-project.org/web/packages/phangorn/index.html>
- Information and manual about pegas
<http://ape-package.ird.fr/pegas.html>

Where to look for the help III

- R help mailing list <https://stat.ethz.ch/mailman/listinfo/r-help> (web interface <http://r.789695.n4.nabble.com/>)
- R announce mailing list
<https://stat.ethz.ch/mailman/listinfo/r-announce>
- R manuals <http://cran.r-project.org/manuals.html>
- Books about R <http://www.r-project.org/doc/bib/R-books.html>
- R at StackOverflow StackExchange
<https://stackoverflow.com/questions/tagged/r>
- R at CrossValidated StackExchange
<https://stats.stackexchange.com/questions/tagged/r>
- The R journal <http://journal.r-project.org/>
- R-bloggers — aggregation of R blogs <http://www.r-bloggers.com/>

Where to look for the help IV

- R on The Molecular Ecologist
<http://www.molecularecologist.com/category/software/r/>
- R tutorial <http://www.r-tutor.com/>
- Cookbook for R <http://www.cookbook-r.com/>
- spatial R <https://sites.google.com/site/spatialr/>
- Statistics with R http://zoonek2.free.fr/UNIX/48_R/all.html
- Springer R series
<https://www.springer.com/series/6991?detailsPage=titles>
- Biostars — general bioinformatics forum <https://www.biostars.org/>
- Biology — general forum about biology at StackExchange
<https://biology.stackexchange.com/>

Where to look for the help V

- Do not hesitate to ask on the forum or contact author of package with which you have problem, preferably through some public forum or mailing list, they usually respond quickly and helpfully...
- Uncle Google is your friend ("how to XXX in R")...
- R packages commonly contain vignettes (tutorials) — list them by **vignette()** and load selected by **vignette("VignetteName")**
- List available training datasets from various R packages by **data()** and load selected by **data(DatasetName)**

Further reading



Michael J. Crawley

The R Book, second edition

Wiley, 2012



Emmanuel Paradis

Analysis of Phylogenetics and Evolution with R, second edition

Springer, 2012

<http://ape-package.ird.fr/APER.html>



Paurush Praveen Sinha

Bioinformatics with R Cookbook

Packt Publishing, 2014

The end

Our course is over...

...I hope it was helpful for You...

...any feedback is welcomed...

...happy **R** hacking...

Typesetting using X_EL_AT_EX on openSUSE GNU/Linux