

An introduction to *adeget* 2.0.0

Thibaut Jombart *

Imperial College London

MRC Centre for Outbreak Analysis and Modelling

June 23, 2015

Abstract

This vignette provides an introductory tutorial to the *adeget* package [1] for the R software [2]. This package implements tools to handle, analyse and simulate genetic data. Originally developed for multiallelic, codominant markers such as microsatellites, *adeget* now also handles dominant markers, allows for any ploidy in the data, handles SNPs and sequence data, and implements a memory-efficient storage for genome-wide SNP data. This tutorial provides an overview of *adeget*'s basic functionalities. Since *adeget* 1.4-0, this tutorial is no longer distributed as a package vignette. Also note that *adeget* has undergone substantial changes with version 2.0.0, including a reform of the data structure and new accessors, all documented in this tutorial.

*tjombart@imperial.ac.uk

Contents

1	Introduction	3
2	Getting started	4
2.1	Installing the package - stable version	4
2.2	Installing the package - devel version	4
2.3	Getting help in R	5
2.4	Asking help on a forum	5
2.5	Bug report, feature requests, contributions: we are all one!	6
3	Object classes	8
3.1	genind objects	8
3.2	genpop objects	11
3.3	Using accessors	12
4	Importing/exporting data	15
4.1	Importing data from GENETIX, STRUCTURE, FSTAT, Genepop	15
4.2	Importing data from other software	15
4.3	Handling presence/absence data	17
4.4	SNPs data	21
4.5	Extracting polymorphism from DNA sequences	22
4.6	Extracting polymorphism from proteic sequences	26
4.7	Using genind/genpop constructors	30
4.8	Exporting data	32
5	Basics of data analysis	33
5.1	Manipulating the data	33
5.2	Using summaries	40
5.3	Testing for Hardy-Weinberg equilibrium	43
5.4	Measuring and testing population structure (a.k.a F statistics)	43
5.5	Estimating inbreeding	46
6	Multivariate analysis	50
6.1	General overview	50
6.2	Performing a Principal Component Analysis on genind objects	52
6.3	Performing a Correspondance Analysis on genpop objects	59
7	Spatial analysis	64
7.1	Isolation by distance	64
7.1.1	Testing isolation by distance	64
7.1.2	Cline or distant patches?	66
7.2	Using Monmonier's algorithm to define genetic boundaries	68
8	Simulating hybridization	78

1 Introduction

This tutorial introduces some basic functionalities of the *adeigenet* package for R [2]. The purpose of this package is to provide tools for handling, analysing and simulating genetic data, with an emphasis on multivariate approaches and exploratory methods. Standard multivariate analyses are implemented in the *ade4* package [3], of which *adeigenet* was originally an extension. However, the package has since grown methods of its own such as the Discriminant Analysis of Principal Components (DAPC, [4]), the spatial Principal Components Analysis (sPCA, [5]), or the *SeqTrack* algorithm [6]. In this tutorial, we introduce the main data structures, show how to import data into *adeigenet*, and cover some basic population genetics and multivariate analysis.

Other tutorials are available via the command `adeigenetTutorial`:

- `adeigenetTutorial("spca")`: tutorial on the sPCA
- `adeigenetTutorial("dapc")`: tutorial on the DAPC
- `adeigenetTutorial("genomics")`: tutorial on handling large SNP datasets using `genlight` objects

2 Getting started

2.1 Installing the package - stable version

Before going further, we shall make sure that *adegenet* is well installed on the computer. The current version of the package is 2.0.0. Make sure you have a recent version of R ($\geq 3.0.3$) by typing:

```
R.version.string  
## [1] "R Under development (unstable) (2015-04-13 r68171)"
```

Then, install *adegenet* with dependencies using:

```
install.packages("adegenet", dep=TRUE)
```

We can now load the package alongside other useful packages:

```
library("ape")  
library("pegas")  
library("seqinr")  
library("ggplot2")
```

```
library("adegenet")
```

If at some point you are unsure about the version of the package, you can check it using:

```
packageDescription("adegenet", fields = "Version")  
## [1] "2.0.0"
```

adegenet version should read 2.0.0.

2.2 Installing the package - devel version

The development of *adegenet* is hosted on github:

<https://github.com/thibautjombart/adegenet>.

You can install this version using the package *devtools* and the following commands:

```
library("devtools")  
install_github("thibautjombart/adegenet")  
library("adegenet")
```

The development version may implement new features and fix known issues. However, it may also occasionally be broken, as this is our working copy of the project. Usual disclaimers apply here: this package is provided with no warranty, etc. If unsure, use the stable version.

2.3 Getting help in R

There are several ways of getting information about R in general, and about *adegenet* in particular. The function `help.search` is used to look for help on a given topic. For instance:

```
help.search("Hardy-Weinberg")
```

replies that there is a function `HWE.test.genind` in the *adegenet* package, and other similar functions in *genetics* and *pegas*. To get help for a given function, use `?foo` where `foo` is the function of interest. For instance (quotes and parentheses can be removed):

```
?spca
```

will open up the manpage of the spatial principal component analysis [5]. At the end of a manpage, an ‘example’ section often shows how to use a function. This can be copied and pasted to the console, or directly executed from the console using `example`. For further questions concerning R, the function `RSiteSearch` is a powerful tool for making online researches using keywords in R’s archives (mailing lists and manpages).

adegenet has a few extra documentation sources. Information can be found from the website (<http://adegenet.r-forge.r-project.org/>), in the ‘documents’ section, including several tutorials and a manual which compiles all manpages of the package, and a dedicated mailing list with searchable archives. To open the website from R, use:

```
adegenetWeb()
```

The same can be done for tutorials, using `adegenetTutorial` (see manpage to choose the tutorial to open). You will also find an overview of the main functionalities of the package typing:

```
?adegenet
```

Note that you can also browse help pages as html pages, using:

```
help.start()
```

To go to the *adegenet* page, click ‘packages’, ‘adegenet’, and ‘adegenet-package’.

2.4 Asking help on a forum

Several mailing lists are available to find different kinds of information on R. *adegenet* has its own dedicated forum/mailling list: adegenet-forum@lists.r-forge.r-project.org. To avoid spam, this list is filtered; subscription is recommended, and can be done at:

<https://lists.r-forge.r-project.org/cgi-bin/mailman/listinfo/adegenet-forum>

Posting questions on R forums can sometimes be a traumatic experience, and we are trying to avoid this as much as possible on the adegenet forum. To this end, the following points are worth keeping in mind:

- **read the doc first:** manpages and tutorials take an awful long time to write and maintain; make sure your answer is not in an obvious place before asking a question; pretending to have read all the available doc while you have not even looked at the basics tutorial is a clever, yet often unsuccessful strategy.
- **search the archives:** adegenet forum has searchable archives (see the adegenet website); your answer may be there already, so it is worth checking.
- **give us info:** you tried something, it is not working.. give us some information: what version of adegenet are you using, what commands did you enter and what was the output, etc.
- **avoid personal messages:** the adegenet forum has plenty of advantages: several people are likely to reply and participate in the conversation, answers are generally faster, and all of this is archived and searchable. Please do not email the developers directly, unless you need to discuss confidential matters.
- **short answers are okay:** some answers will be short. Do not take them as rude, or think people are upset: answering questions on a forum is a time-consuming activity and the reward for it is low. Sometimes the best answer will be pointing to relevant documentation, e.g. “Please check ?xvalDapc”. If you get this, we (most likely) still like you.

The adegenet forum is not the only forum that might be relevant. Others include:

- *R-sig-genetics*: genetics in R.
<https://stat.ethz.ch/mailman/listinfo/r-sig-genetics>
- *R-sig-phylo*: phylogenetics in R.
<https://stat.ethz.ch/mailman/listinfo/r-sig-phylo>
- *R-help*: general questions about R.
<https://stat.ethz.ch/mailman/listinfo/r-help>

Please avoid double-posting (it is often considered to be rude).

2.5 Bug report, feature requests, contributions: we are all one!

Free software evolves through interactions between communities of developers and users. This is especially true in R, where these two communities are very much intermingled. In short: **we are all contributors!** These contributions include:

- **asking a question:** see section above

- **asking for a new feature:** something useful is missing, and you think it will be useful to others? Say it! Post a feature request using github's *issues*:
<https://github.com/thibautjombart/adegetnet/issues>
- **reporting a possible bug:** bugs are rare, but if you think you found one, post it as an issue on githubL:
<https://github.com/thibautjombart/adegetnet/issues>
- **contributions:** github makes contributions very easy; fork the project, make the changes you want, and when you are happy and the package passes the checks, send a pull request; for more on this go to the github page:
<https://github.com/thibautjombart/adegetnet> And please, remember to add yourself as a contributor in the DESCRIPTION and relevant manpages!

3 Object classes

Two main classes of objects are used for storing genetic marker data, depending on the level at which the genetic information is considered: **genind** is used for individual genotypes, whereas **genpop** is used for alleles numbers counted by populations. Note that the term 'population', here and later, is employed in a broad sense: it simply refers to any grouping of individuals. The specific class **genlight** is used for storing large genome-wide SNPs data. See the *genomics* tutorial for more information on this topic.

3.1 genind objects

These objects store **genetic** data at an **individual** level, plus various meta-data (e.g. group membership). **genind** objects can be obtained by reading data files from other software, from a `data.frame` of genotypes, by conversion from a table of allele counts, or even from aligned DNA or proteic sequences (see 'importing data'). Here, we introduce this class using the dataset **nancycats**, which is already stored as a **genind** object:

```
data(nancycats)
is.genind(nancycats)

## [1] TRUE

nancycats

## /// GENIND OBJECT ///////////
##
## // 237 individuals; 9 loci; 108 alleles; size: 158.4 Kb
##
## // Basic content
##   @tab: 237 x 108 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 8-18)
##   @loc.fac: locus factor for the 108 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: genind(tab = truenames(nancycats)$tab, pop = truenames(nancycats)$pop)
##
## // Optional content
##   @pop: population of each individual (group size range: 9-23)
##   @strata: - empty -
##   @hierarchy: - empty -
##   @other: a list containing: xy
```

A **genind** object is formal S4 object with several slots, accessed using the '@' operator (see `class?genind`). Note that the '\$' is also implemented for adegenet objects, so that

slots can be accessed as if they were components of a list.

genind objects possess the following slots:

- **tab**: a matrix of alleles counts (individuals in rows, alleles in columns).
- **loc.n.all**: the number of alleles in each marker.
- **loc.fac**: a factor indicating which columns in **@tab** correspond to which marker.
- **all.names**: a list of allele names for each locus.
- **ploidy**: a vector of integer indicating the ploidy of each individual; constant ploidy is assumed across different loci of a single individual, but different individuals can have different ploidy.
- **type**: a character string indicating whether the marker is codominant (**codom**) or presence/absence (**PA**).
- **call**: the matched call, i.e. command used to create the object.
- **pop**: (optional) a factor storing group membership of the individuals; when present, method needing a population information will automatically use this if nothing else is provided.
- **strata**: (optional) a data.frame storing hierarchical structure of individuals (see dedicated tutorial).
- **hierarchy**: (optional) a formula defining the hierarchical structure of individuals (see dedicated tutorial).
- **other**: (optional) a list storing optional information.

Slots can be accessed using **'@'** or **'\$'**, although it is recommended to use accessors to retrieve and change slot values (see section 'Using accessors').

The main slot of a **genind** is the table of allelic counts **@tab**, with individuals in rows and alleles in columns. For instance:

```
nancycats[10:18, loc="fca8"]@tab
```

##	fca8.117	fca8.119	fca8.121	fca8.123	fca8.127	fca8.129	fca8.131
## N224	0	0	0	0	0	0	0
## N7	0	0	0	0	0	0	0
## N141	0	0	0	0	0	1	0
## N142	0	0	0	0	0	1	0
## N143	0	0	0	0	0	0	0
## N144	0	0	0	0	0	0	1
## N145	0	0	0	0	0	1	0

## N146	0	0	0	0	0	1	0
## N147	0	0	0	0	0	1	0
##	fca8.133	fca8.135	fca8.137	fca8.139	fca8.141	fca8.143	fca8.145
## N224	0	2	0	0	0	0	0
## N7	0	0	1	0	1	0	0
## N141	1	0	0	0	0	0	0
## N142	1	0	0	0	0	0	0
## N143	2	0	0	0	0	0	0
## N144	0	1	0	0	0	0	0
## N145	0	1	0	0	0	0	0
## N146	1	0	0	0	0	0	0
## N147	0	1	0	0	0	0	0
##	fca8.147	fca8.149					
## N224	0	0					
## N7	0	0					
## N141	0	0					
## N142	0	0					
## N143	0	0					
## N144	0	0					
## N145	0	0					
## N146	0	0					
## N147	0	0					

Individual 'N224' is an homozygote for the allele 135 at the locus 'fca8', while N141" is an heterozygote with alleles 129/133. Note that as of *adegenet* 2.0.0, this table is no storing data as relative frequencies. These can still be obtained using the accessor 'tab'. The particular case of presence/absence data is described in a dedicated section (see 'Handling presence/absence data'). As of version 2.0.0 of *adegenet*, the slots @strata and @hierarchy implement hierarchical population structures. See dedicated tutorial for more on this topic.

Note that objects can be regenerated using the matched call stored in *genind* objects, *i.e.* the instruction that created it. For instance:

```
obj <- read.genetix(system.file("files/nancycats.gtx",package="adegenet"))

##
## Converting data from GENETIX to a genind object...
##
## ...done.

obj$call

## read.genetix(file = system.file("files/nancycats.gtx", package = "adegenet"))

toto <- eval(obj$call)
```

```
##
## Converting data from GENETIX to a genind object...
##
## ...done.

identical(obj,toto)

## [1] TRUE
```

3.2 genpop objects

These objects store **genetic** data at a **population** level, plus various meta-data. Their structure is nearly identical to **genind** objects, only simpler as they no longer store group membership for individuals. **genpop** objects are created from **genind** objects using **genind2genpop**:

```
data(nancycats)
catpop <- genind2genpop(nancycats)

##
## Converting data from a genind to a genpop object...
##
## ...done.

catpop

## /// GENPOP OBJECT ///////////
##
## // 17 populations; 9 loci; 108 alleles; size: 28.5 Kb
##
## // Basic content
##   @tab: 17 x 108 matrix of alleles counts
##   @loc.n.all: number of alleles per locus (range: 8-18)
##   @loc.fac: locus factor for the 108 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: genind2genpop(x = nancycats)
##
## // Optional content
##   @other: a list containing: xy
```

As in **genind** objects, data are stored as numbers of alleles, but this time for populations (here, cat colonies):

```
catpop$stab[1:5,1:10]
```

```
##      fca8.117 fca8.119 fca8.121 fca8.123 fca8.127 fca8.129 fca8.131
## P01         0         0         0         0         0         0         0
## P02         0         0         0         0         0        10         9
## P03         0         0         0         4         0         0         0
## P04         0         0         0         3         0         0         0
## P05         0         0         0         1         0         0         0
##      fca8.133 fca8.135 fca8.137
## P01         2         9         1
## P02         8        14         2
## P03         0         1        10
## P04         1         7        17
## P05         0         7        10
```

3.3 Using accessors

One advantage of formal (S4) classes is that they allow for interacting simply with possibly complex objects. This is made possible by using accessors, i.e. functions that extract information from an object, rather than accessing the slots directly. Another advantage of this approach is that as long as accessors remain identical on the user's side, the internal structure of an object may change from one release to another without generating errors in old scripts. Although **genind** and **genpop** objects are fairly simple, we recommend using accessors whenever possible to access or modify their content.

Available accessors are:

- **nInd**: returns the number of individuals in the object; only for **genind**.
- **nLoc**: returns the number of loci.
- **nAll**: returns the number of alleles for each locus.
- **nPop**: returns the number of loci (SNPs).
- **tab**: returns a table of allele numbers, or frequencies (if requested), with optional replacement of missing values; replaces the former accessor '**truenames**'.
- **indNames[†]**: returns/sets labels for individuals; only for **genind**.
- **locNames[†]**: returns/sets labels for loci.
- **alleles[†]**: returns/sets alleles.
- **ploidy[†]**: returns/sets ploidy of the individuals; when setting values, a single value can be provided, in which case constant ploidy is assumed.

- `pop†`: returns/sets a factor grouping individuals; only for `genind`.
- `strata†`: returns/sets data defining strata of individuals; only for `genind`.
- `hier†`: returns/sets hierarchical groups of individuals; only for `genind`.
- `other†`: returns/sets misc information stored as a list.

where [†] indicates that a replacement method is available using `<-`; for instance:

```
head(indNames(nancycats),10)

##      001      002      003      004      005      006      007      008      009      010
## "N215" "N216" "N217" "N218" "N219" "N220" "N221" "N222" "N223" "N224"

indNames(nancycats) <- paste("cat", 1:nInd(nancycats),sep=".")
head(indNames(nancycats),10)

##  [1] "cat.1" "cat.2" "cat.3" "cat.4" "cat.5" "cat.6" "cat.7"
##  [8] "cat.8" "cat.9" "cat.10"
```

Some accessors such as `locNames` may have specific options; for instance:

```
locNames(nancycats)

## [1] "fca8" "fca23" "fca43" "fca45" "fca77" "fca78" "fca90" "fca96" "fca37"
```

returns the names of the loci, while:

```
temp <- locNames(nancycats, withAlleles=TRUE)
head(temp, 10)

##  [1] "fca8.117" "fca8.119" "fca8.121" "fca8.123" "fca8.127" "fca8.129"
##  [7] "fca8.131" "fca8.133" "fca8.135" "fca8.137"
```

returns the names of the alleles in the form 'loci.allele'.

The slot 'pop' can be retrieved and set using `pop`:

```
obj <- nancycats[sample(1:50,10)]
pop(obj)

##  [1] P02 P02 P02 P04 P02 P02 P02 P02 P03 P02
## Levels: P02 P03 P04

pop(obj) <- rep("newPop",10)
pop(obj)

##  [1] newPop newPop newPop newPop newPop newPop newPop newPop newPop newPop
## Levels: newPop
```

Accessors make things easier. For instance, when setting new names for loci, the columns of `@tab` are renamed automatically:

```
head(colnames(tab(obj)),20)

## [1] "fca8.117" "fca8.119" "fca8.121" "fca8.123" "fca8.127"
## [6] "fca8.129" "fca8.131" "fca8.133" "fca8.135" "fca8.137"
## [11] "fca8.139" "fca8.141" "fca8.143" "fca8.145" "fca8.147"
## [16] "fca8.149" "fca23.128" "fca23.130" "fca23.132" "fca23.136"

locNames(obj)

## [1] "fca8" "fca23" "fca43" "fca45" "fca77" "fca78" "fca90" "fca96" "fca37"

locNames(obj)[1] <- "newLocusName"
locNames(obj)

## [1] "newLocusName" "fca23" "fca43" "fca45"
## [5] "fca77" "fca78" "fca90" "fca96"
## [9] "fca37"

head(colnames(tab(obj)),20)

## [1] "newLocusName.117" "newLocusName.119" "newLocusName.121"
## [4] "newLocusName.123" "newLocusName.127" "newLocusName.129"
## [7] "newLocusName.131" "newLocusName.133" "newLocusName.135"
## [10] "newLocusName.137" "newLocusName.139" "newLocusName.141"
## [13] "newLocusName.143" "newLocusName.145" "newLocusName.147"
## [16] "newLocusName.149" "fca23.128" "fca23.130"
## [19] "fca23.132" "fca23.136"
```

An additional advantage of using accessors is they are most of the time safer to use. For instance, `pop<-` will check the length of the new group membership vector against the data, and complain if there is a mismatch. It also converts the provided replacement to a factor, while the command:

```
obj@pop <- rep("newPop",10)

## Error in checkAtAssignment("genind", "pop", "character"): assignment of
## an object of class "character" is not valid for @'pop' in an object of class
## "genind"; is(value, "factorOrNULL") is not TRUE
```

generates an error (since replacement is not a factor).

4 Importing/exporting data

4.1 Importing data from GENETIX, STRUCTURE, FSTAT, Genepop

Data can be read from the software GENETIX (extension .gtx), STRUCTURE (.str or .stru), FSTAT (.dat) and Genepop (.gen) files, using the corresponding `read` function: `read.genetix`, `read.structure`, `read.fstat`, and `read.genepop`. These functions take as main argument the path (as a string of characters) to an input file, and produce a `genind` object. Alternatively, one can use the function `import2genind` which detects a file format from its extension and uses the appropriate routine. For instance:

```
obj1 <- read.genetix(system.file("files/nancycats.gtx",package="adegenet"))

##
## Converting data from GENETIX to a genind object...
##
## ...done.

obj2 <- import2genind(system.file("files/nancycats.gtx", package="adegenet"))

##
## Converting data from GENETIX to a genind object...
##
## ...done.

all.equal(obj1,obj2)

## [1] "Attributes: < Component \"call\": target, current do not match when deparsed >"
```

The only difference between `obj1` and `obj2` is their call (which is normal as they were obtained from different command lines).

4.2 Importing data from other software

Raw genetic markers data are often stored as tables with individuals in row and markers in column, where each entry is a character string coding the alleles possessed at one locus. Such data are easily imported into R as a `data.frame`, using for instance `read.table` for text files or `read.csv` for comma-separated text files. Then, the obtained `data.frame` can be converted into a `genind` object using `df2genind`.

There are only a few pre-requisite the data should meet for this conversion to be possible. The easiest and clearest way of coding data is using a separator between alleles. For instance, "80/78", "80—78", or "80,78" are different ways of coding a genotype at a microsatellite locus with alleles '80' and '78'. Note that for haploid data, no separator shall be used. The only constraint when using a separator is that the same separator is used in all

the dataset. There are no constraints as to i) the type of separator used or ii) the ploidy of the data. These parameters can be set in `df2genind` through arguments `sep` and `ploidy`, respectively.

Alternatively, no separator may be used provided a fixed number of characters is used to code each allele. For instance, in a diploid organism, "0101" is an homozygote 1/1 while "1209" is a heterozygote 12/09 in a two-character per allele coding scheme. In a tetraploid system with one character per allele, "1209" will be understood as 1/2/0/9.

Here, we provide an example using randomly generated tetraploid data and no separator.

```
temp <- lapply(1:30, function(i) sample(1:9, 4, replace=TRUE))
temp <- sapply(temp, paste, collapse="")
temp <- matrix(temp, nrow=10, dimnames=list(paste("ind",1:10), paste("loc",1:3)))
temp

##          loc 1  loc 2  loc 3
## ind 1  "4814" "3958" "9539"
## ind 2  "3689" "1626" "5751"
## ind 3  "9688" "7134" "9962"
## ind 4  "4298" "7682" "6736"
## ind 5  "2716" "2912" "2128"
## ind 6  "2327" "6897" "6857"
## ind 7  "7729" "5979" "6314"
## ind 8  "7643" "2655" "1913"
## ind 9  "2155" "2116" "1365"
## ind 10 "8142" "9295" "8434"

obj <- df2genind(temp, ploidy=4, sep="")
obj

## /// GENIND OBJECT ///////////
##
## // 10 individuals; 3 loci; 27 alleles; size: 9.5 Kb
##
## // Basic content
##   @tab: 10 x 27 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 9-9)
##   @loc.fac: locus factor for the 27 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 4-4)
##   @type: codom
##   @call: df2genind(X = temp, sep = "", ploidy = 4)
##
## // Optional content
```



```
## @pop: - empty -
## @strata: - empty -
## @hierarchy: - empty -
## @other: - empty -
```

`obj` is a `genind` containing the same information, but recoded as a matrix of allele numbers (`$tab` slot). We can check that the conversion was exact by converting back the object into a table of character strings (function `genind2df`):

```
genind2df(obj, sep="|")

##          loc 1    loc 2    loc 3
## ind 1  4|4|8|1 3|9|5|8 9|9|5|3
## ind 2  8|3|6|9 1|6|6|2 5|5|7|1
## ind 3  8|8|6|9 3|1|7|4 9|9|6|2
## ind 4  4|8|9|2 8|6|2|7 3|7|6|6
## ind 5  1|6|2|7 9|1|2|2 1|2|2|8
## ind 6  3|2|2|7 9|8|6|7 5|7|6|8
## ind 7  9|2|7|7 9|9|5|7 3|1|6|4
## ind 8  4|3|6|7 5|5|6|2 9|3|1|1
## ind 9  1|2|5|5 1|1|6|2 5|3|1|6
## ind 10 4|8|1|2 9|9|5|2 3|8|4|4
```

4.3 Handling presence/absence data

adegenet was primarily designed to handle codominant, multiallelic markers like microsatellites. However, dominant markers like AFLP can be used as well. In such a case, only presence/absence of alleles can be deduced accurately from the genotypes. This has several consequences, like the inability to compute allele frequencies. Hence, some functionalities in *adegenet* won't be available for dominant markers.

From version 1.2-3 of *adegenet*, the distinction between both types of markers is made by the slot `@type` of `genind` or `genpop` objects, which equals `codom` for codominant markers, and `PA` for presence/absence data. In the latter case, the `'tab'` slot of a `genind` object no longer contains allele frequencies, but only presence/absence of alleles in a genotype. Similarly, the `tab` slot of a `genpop` object no longer contains counts of alleles in the populations; instead, it contains the number of genotypes in each population possessing at least one copy of the concerned alleles. Moreover, in the case of presence/absence, the slots `'loc.n.all'`, `'loc.fac'`, and `'all.names'` become useless, and are thus all set to `NULL`.

Objects of type `'PA'` are otherwise handled like usual (type `'codom'`) objects. Operations that are not available for `PA` type will issue an appropriate error message.

Here is an example using a toy dataset `'AFLP.txt'` that can be downloaded from the *adegenet* website, section `'Documentation'`:

```
dat <- read.table(system.file("files/AFLP.txt",package="adeigenet"), header=TRUE)
dat
```

```
##      loc1 loc2 loc3 loc4
## indA    1    0    1    1
## indB    0    1    1    1
## indC    1    1    0    1
## indD    0   NA    1   NA
## indE    1    1    0    0
## indF    1    0    1    1
## indG    0    1    1    0
```

The function `df2genind` is used to obtain a `genind` object:

```
obj <- df2genind(dat, ploidy=1, type="PA")
obj

## /// GENIND OBJECT ///////////
##
## // 7 individuals; 4 loci; 4 alleles; size: 4 Kb
##
## // Basic content
##   @tab: 7 x 4 matrix of allele counts
##   @loc.n.all: NULL
##   @loc.fac: NULL
##   @all.names: NULL
##   @ploidy: ploidy of each individual (range: 1-1)
##   @type: PA
##   @call: df2genind(X = dat, ploidy = 1, type = "PA")
##
## // Optional content
##   @pop: - empty -
##   @strata: - empty -
##   @hierarchy: - empty -
##   @other: - empty -
```

```
tab(obj)
```

```
##      loc1 loc2 loc3 loc4
## indA    1    0    1    1
## indB    0    1    1    1
## indC    1    1    0    1
## indD    0   NA    1   NA
## indE    1    1    0    0
## indF    1    0    1    1
## indG    0    1    1    0
```

One can see that for instance, the summary of this object is more simple (no numbers of alleles per locus, no heterozygosity):

```
pop(obj) <- rep(c('a','b'),4:3)
summary(obj)

##
## # Total number of genotypes: 7
##
## # Population sample sizes:
## a b
## 4 3
##
## # Percentage of missing data:
## [1] 7.142857
```

But we can still perform basic manipulation, like converting our object into a genpop:

```
obj2 <- genind2genpop(obj)

##
## Converting data from a genind to a genpop object...
##
## ...done.

obj2

## /// GENPOP OBJECT ///////////
##
## // 2 populations; 4 loci; 4 alleles; size: 2.8 Kb
##
## // Basic content
## @tab: 2 x 4 matrix of alleles counts
## @loc.n.all: NULL
## @loc.fac: NULL
## @all.names: NULL
## @ploidy: ploidy of each individual (range: 1-1)
## @type: PA
## @call: genind2genpop(x = obj)
##
## // Optional content
## @other: - empty -

tab(obj2)

## loc1 loc2 loc3 loc4
## a 2 2 3 3
## b 2 2 2 1
```

To continue with the toy example, we can perform a simple PCA. Allele presence absence data are extracted and NAs replaced using `tab`:

```
X <- tab(obj, NA.method="mean")
```

Now the PCA is performed and plotted:

```
## make PCA
pca1 <- dudi.pca(X,scannf=FALSE,scale=FALSE)
temp <- as.integer(pop(obj))
myCol <- transp(c("blue","red"),.7)[temp]
myPch <- c(15,17)[temp]

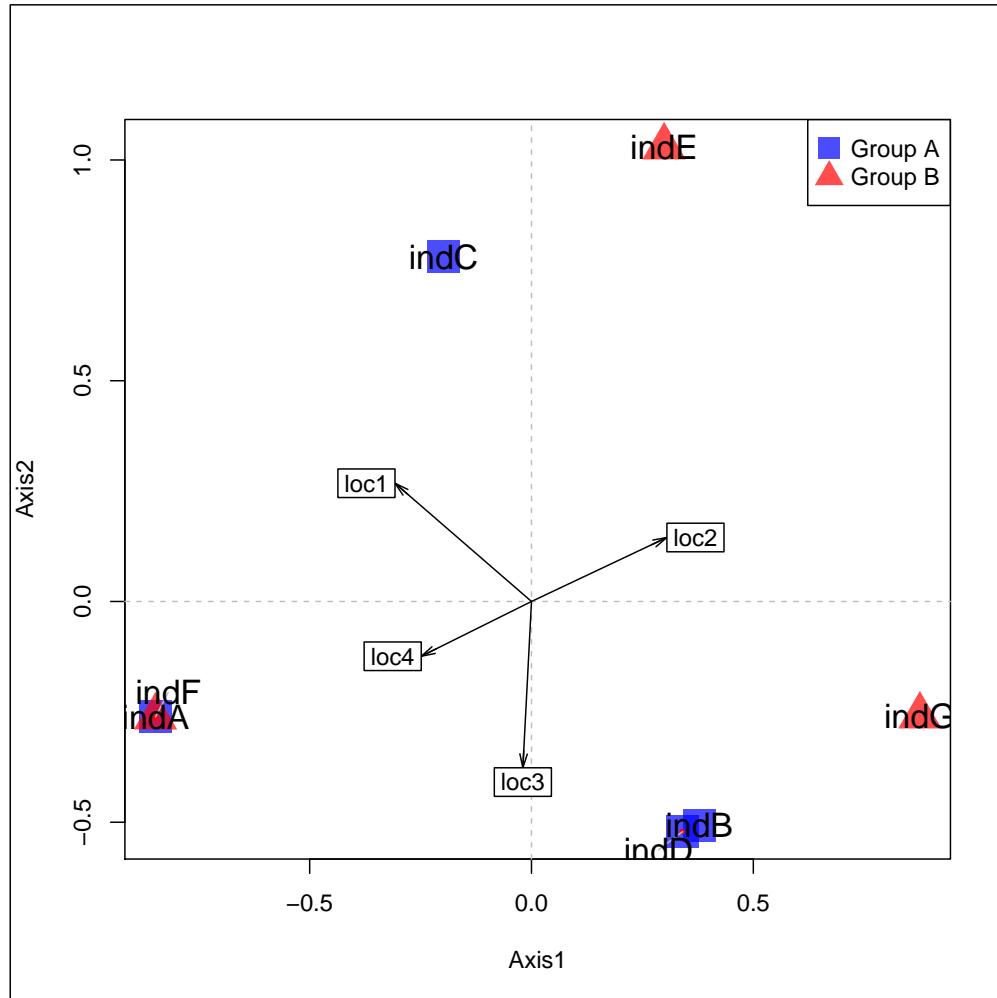
## basic plot
plot(pca1$li, col=myCol, cex=3, pch=myPch)

## use wordcloud for non-overlapping labels
library(wordcloud)

## Loading required package: RColorBrewer

textplot(pca1$li[,1], pca1$li[,2], words=rownames(X), cex=1.4, new=FALSE)

## legend the axes by adding loadings
abline(h=0,v=0,col="grey",lty=2)
s.arrow(pca1$c1*.5, add.plot=TRUE)
legend("topright", pch=c(15,17), col=transp(c("blue","red"),.7),
      leg=c("Group A","Group B"), pt.cex=2)
```



More generally, multivariate analyses from `ade4`, `sPCA` (`spca`), `DAPC` (`dapc`), the global and local tests (`global.rtest`, `local.rtest`), or the Monmonier's algorithm (`monmonier`) will work just fine with presence/absence data. However, it is clear that the usual Euclidean distance (used in PCA and sPCA), as well as many other distances, is not as accurate to measure genetic dissimilarity using presence/absence data as it is when using allele frequencies. The reason for this is that in presence/absence data, a part of the information is simply hidden. For instance, two individuals possessing the same allele will be considered at the same distance, whether they possess one or more copies of the allele. This might be especially problematic in organisms having a high degree of ploidy.

4.4 SNPs data

In *ade4*, SNP data can be handled in two different ways. For relatively small datasets (up to a few thousand SNPs) SNPs can be handled as usual codominant markers such as microsatellites using `genind` objects. In the case of genome-wide SNP data (from hundreds of thousands to millions of SNPs), `genind` objects are no longer efficient representation of the data. In this case, we use `genlight` objects to store and handle information with maximum efficiency and minimum memory requirements. See the tutorial *genomics* for more information. Below, we introduce only the case of SNPs handled using `genind` objects.

The most convenient way to convert SNPs into a **genind** is using **df2genind**, which is described in the previous section. Let **dat** be an input matrix, as can be read into R using **read.table** or **read.csv**, with genotypes in row and SNP loci in columns.

```
dat <- matrix(sample(c("a","t","g","c"), 15, replace=TRUE),nrow=3)
rownames(dat) <- paste("genot.", 1:3)
colnames(dat) <- 1:5
dat

##           1    2    3    4    5
## genot. 1 "c" "t" "t" "t" "a"
## genot. 2 "a" "t" "a" "c" "t"
## genot. 3 "g" "c" "g" "t" "a"

obj <- df2genind(dat, ploidy=1)
tab(obj)

##           1.c 1.a 1.g 2.t 2.c 3.t 3.a 3.g 4.t 4.c 5.a 5.t
## genot. 1    1    0    0    1    0    1    0    0    1    0    1    0
## genot. 2    0    1    0    1    0    0    1    0    0    1    0    1
## genot. 3    0    0    1    0    1    0    0    1    1    0    1    0
```

obj is a **genind** containing the SNPs information, which can be used for further analysis in **adegenet**.

4.5 Extracting polymorphism from DNA sequences

This section only covers the cases of relatively small datasets which can be handled efficiently using **genind** objects. For bigger (near full-genomes) datasets, SNPs can be extracted from *fasta* files into a **genlight** object using **fasta2genlight**. See the tutorial on *genomics* for more information.

DNA sequences can be read into R using *ape*'s **read.dna** (fasta and clustal formats), or *adegenet*'s **fasta2DNABin** for a RAM-greedy implementation (fasta alignments only). Other options include reading data directly from GenBank using **read.GenBank**, or from other public databases using the *seqinr* package and transforming the **alignment** object into a **DNABin** using **as.DNABin**. Here, we illustrate this approach by re-using the example of **read.GenBank**. A connection to the internet is required, as sequences are read directly from a distant database.

```
library(ape)
ref <- c("U15717", "U15718", "U15719", "U15720",
        "U15721", "U15722", "U15723", "U15724")
myDNA <- read.GenBank(ref)
myDNA
```

```
## 8 DNA sequences in binary format stored in a list.
##
## All sequences of same length: 1045
##
## Labels: U15717 U15718 U15719 U15720 U15721 U15722 ...
##
## Base composition:
##      a      c      g      t
## 0.267 0.351 0.134 0.247

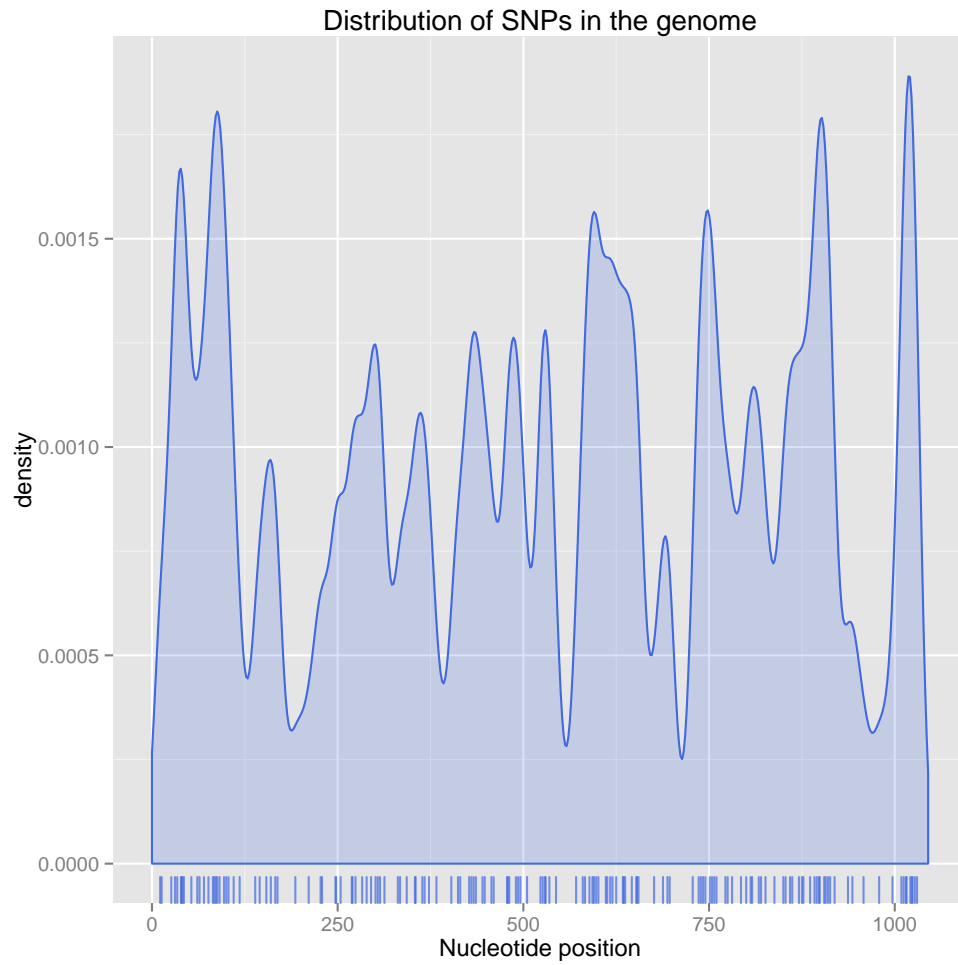
class(myDNA)

## [1] "DNABin"

myDNA <- as.matrix(myDNA)
```

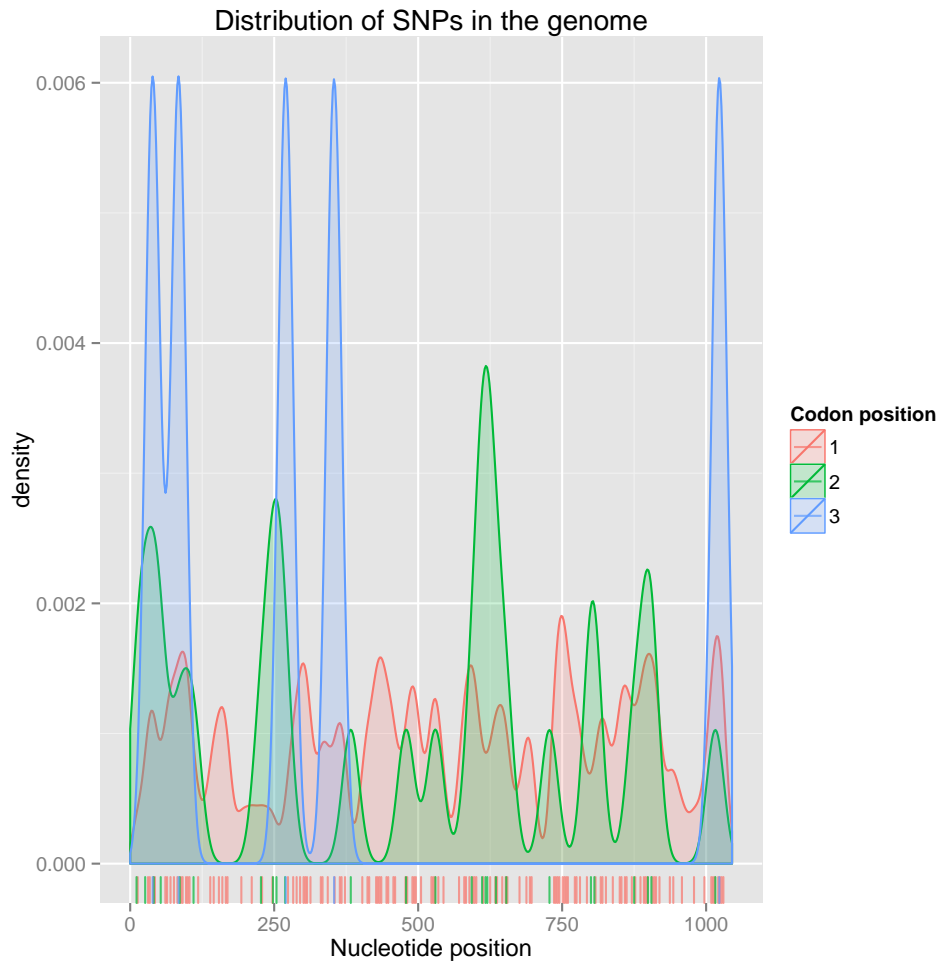
Polymorphism can be characterized using `snpposi.plot` and `snpposi.test`: the first plots SNP density along the alignment, the second tests whether these SNPs are randomly distributed. For instance:

```
snpposi.plot(myDNA, codon=FALSE)
```



By default, the function differentiates nucleotide positions:

```
snpposi.plot(myDNA)
```

In *adegenet*, only polymorphic loci are conserved to form a *genind* object. This conversion is achieved by `DNABin2genind`. This function allows one to specify a threshold for polymorphism; for instance, one could retain only SNPs for which the second largest allele frequency is greater than 1% (using the `polyThres` argument). This is achieved using:

```
obj <- DNABin2genind(myDNA, polyThres=0.01)
obj

## /// GENIND OBJECT ///////////
##
## // 8 individuals; 155 loci; 318 alleles; size: 83.4 Kb
##
## // Basic content
##   @tab: 8 x 318 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 2-4)
##   @loc.fac: locus factor for the 318 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 1-1)
##   @type: codom
```

```
## @call: DNABin2genind(x = myDNA, polyThres = 0.01)
##
## // Optional content
## @pop: - empty -
## @strata: - empty -
## @hierarchy: - empty -
## @other: - empty -
```

Here, out of the 1,045 nucleotides of the sequences, 318 SNPs were extracted and stored as a **genind** object. Positions of the SNPs are stored as names of the loci:

```
head(locNames(obj))

## [1] "11" "13" "26" "31" "34" "39"
```

4.6 Extracting polymorphism from proteic sequences

Alignments of proteic sequences can be exploited in *adegenet* in the same way as DNA sequences (see section above). Alignments are scanned for polymorphic sites, and only those are retained to form a **genind** object. Loci correspond to the position of the residue in the alignment, and alleles correspond to the different amino-acids (AA). Aligned proteic sequences are stored as objects of class **alignment** in the *seqinr* package [9]. See `?as.alignment` for a description of this class. The function extracting polymorphic sites from **alignment** objects is **alignment2genind**.

Its use is fairly simple. It is here illustrated using a small dataset of aligned proteic sequences:

```
library(seqinr)
mase.res <- read.alignment(file=system.file("sequences/test.mase",
                                             package="seqinr"), format = "mase")
mase.res

## $nb
## [1] 6
##
## $nam
## [1] "Langur" "Baboon" "Human" "Rat" "Cow" "Horse"
##
## $seq
## $seq[[1]]
## [1] "-kifercelartlkklgldgykgvslanwvclakwesgynteatnynpgdestdygifqinsrywcnngkpgavdachis"
##
## $seq[[2]]
## [1] "-kifercelartlkrldgyrgislanwvclakwesdyntqatnynpgdqstdygifqinshywcndgkpgavnachis"
```

```

##
## $seq[[3]]
## [1] "-kvfercelartlkrldgyrgislanwmclakwesgyntnatnynagdrstdygifqinsrywncndgkpgavnachls
##
## $seq[[4]]
## [1] "-ktyercefartlkrngmsgyygvsladwvclaquesnyntqarnydpqdqstdygifqinsrywncndgkpraknacgip
##
## $seq[[5]]
## [1] "-kvfercelartlkkldgldgykgvslanwlccltkwessyntkatnynpssestdygifqinskwncndgkpnavgchvs
##
## $seq[[6]]
## [1] "-kvfskcelahklkaqemdggfgygyslanwvcmaeyesnfntrafngknangssdyglfqlnnkwckdnkrsssnacnim
##
##
## $com
## [1] ";empty description\n" ";\\n" ";\\n"
## [4] ";\\n" ";\\n" ";\\n"
##
## attr("class")
## [1] "alignment"

x <- alignment2genind(mase.res)
x

## /// GENIND OBJECT ///////////
##
## // 6 individuals; 82 loci; 212 alleles; size: 51.1 Kb
##
## // Basic content
## @tab: 6 x 212 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-5)
## @loc.fac: locus factor for the 212 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 1-1)
## @type: codom
## @call: alignment2genind(x = mase.res)
##
## // Optional content
## @pop: - empty -
## @strata: - empty -
## @hierarchy: - empty -
## @other: a list containing: com

```

The six aligned protein sequences (`mase.res`) have been scanned for polymorphic sites, and these have been extracted to form the `genind` object `x`. Note that several settings

such as the characters corresponding to missing values (i.e., gaps) and the polymorphism threshold for a site to be retained can be specified through the function's arguments (see `?alignment2genind`).

The names of the loci directly provides the indices of polymorphic sites:

```
head(locNames(x))

## [1] "3" "4" "5" "6" "9" "11"
```

The table of polymorphic sites can be reconstructed easily by:

```
tabAA <- genind2df(x)
dim(tabAA)

## [1] 6 82

tabAA[, 1:20]

##           3 4 5 6 9 11 12 15 16 17 18 19 21 22 24 28 30 32 33 34
## Langur i f e r l r t k l g l d y k v n v l a k
## Baboon i f e r l r t r l g l d y r i n v l a k
## Human v f e r l r t r l g m d y r i n m l a k
## Rat t y e r f r t r n g m s y y v d v l a q
## Cow v f e r l r t k l g l d y k v n l l t k
## Horse v f s k l h k a q e m d f g y n v m a e
```

The global AA composition of the polymorphic sites is given by:

```
table(unlist(tabAA))

##
## a d e f g h i k l m n p q r s t v w y
## 35 38 16 9 33 13 27 28 31 8 44 10 26 47 36 20 42 6 23
```

Now that polymorphic sites have been converted into a **genind** object, simple distances can be computed between the sequences. Note that *adegenet* does not implement specific distances for protein sequences, we only use the simple Euclidean distance. Fancier protein distances are implemented in R; see for instance `dist.alignment` in the *seqinr* package, and `dist.ml` in the *phangorn* package.

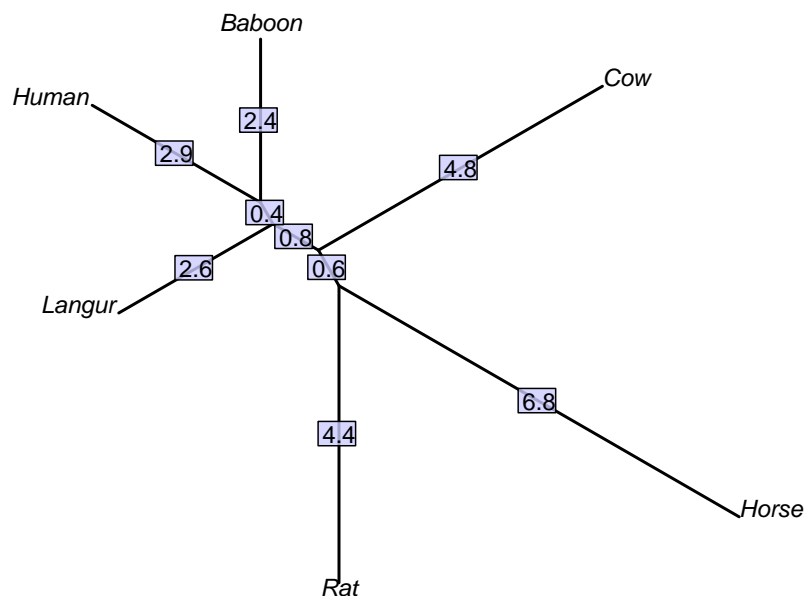
```
D <- dist(tab(x))
D

##           Langur      Baboon      Human      Rat      Cow
## Baboon  5.291503
```

```
## Human    6.000000  5.291503
## Rat      8.717798  8.124038  8.602325
## Cow      7.874008  8.717798  8.944272 10.392305
## Horse    11.313708 11.313708 11.224972 11.224972 11.747340
```

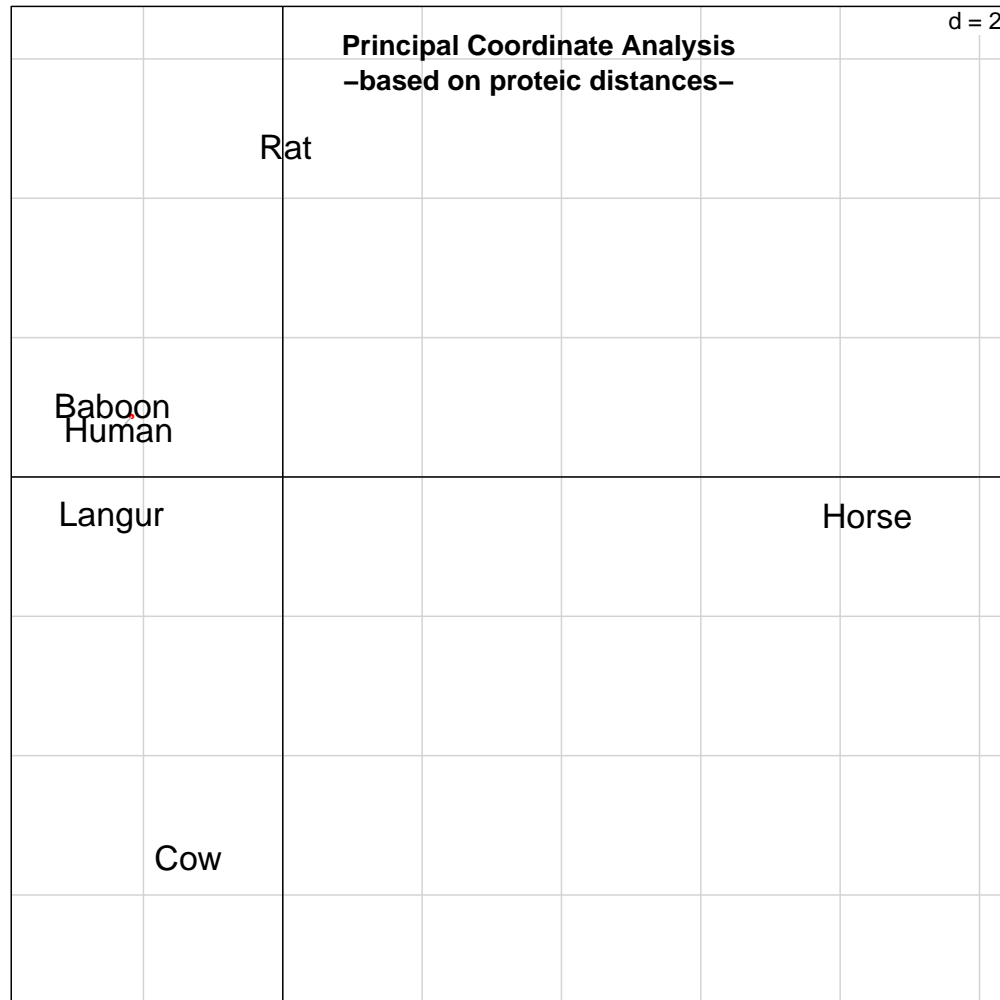
This matrix of distances is small enough for one to interpret the raw numbers. However, it is also very straightforward to represent these distances as a tree or in a reduced space. We first build a Neighbor-Joining tree using the *ape* package:

```
library(ape)
tre <- nj(D)
par(xpd=TRUE)
plot(tre, type="unrooted", edge.w=2)
edgelabels(tex=round(tre$edge.length,1), bg=rgb(.8,.8,1,.8))
```



The best possible planar representation of these Euclidean distances is achieved by Principal Coordinate Analyses (PCoA), which in this case will give identical results to PCA of the original (centred, non-scaled) data:

```
pco1 <- dudi.pco(D, scannf=FALSE,nf=2)
s.label(pco1$li*1.1, clab=0, pch="")
textplot(pco1$li[,1], pco1$li[,2], words=rownames(pco1$li),
         cex=1.4, new=FALSE, xpd=TRUE)
title("Principal Coordinate Analysis\n-based on proteic distances-")
```



4.7 Using genind/genpop constructors

genind or **genpop** objects are best obtained using the procedures described above. However, one can also build new **genind** or **genpop** objects using the constructor **new()**. In this case, the function must be given as main input an object corresponding to the **@tab** slot: a matrix of integers with individuals in row and alleles in columns, with columns being named as '[marker].[allele]'. Here is an example for **genpop** using a dataset from *ade4*:

```
data(microsatt)
microsatt$tab[10:15,12:15]

##          INRA32.168 INRA32.170 INRA32.174 INRA32.176
```

## Mtbeliard	0	0	0	1
## NDama	0	0	0	12
## Normand	1	0	0	2
## Parthenais	8	5	0	3
## Somba	0	0	0	20
## Vosgienne	2	0	0	0

`microsat$tab` contains alleles counts per populations, and can therefore be used to make a `genpop` object. Moreover, column names are set as required, and row names are unique. It is therefore safe to convert these data into a `genpop` using the constructor:

```
toto <- new("genpop", tab=microsat$tab)
toto

## /// GENPOP OBJECT //////////
##
## // 18 populations; 9 loci; 112 alleles; size: 28 Kb
##
## // Basic content
##   @tab: 18 x 112 matrix of alleles counts
##   @loc.n.all: number of alleles per locus (range: 8-17)
##   @loc.fac: locus factor for the 112 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: .local(.Object = .Object, tab = ..1)
##
## // Optional content
##   @other: - empty -

summary(toto)

##
## # Number of populations: 18
##
## # Number of alleles per locus:
## INRA5 INRA32 INRA35 INRA63 INRA72 ETH152 ETH225 INRA16 INRAK
##      8      15      11      10      17      10      14      15      12
##
## # Number of alleles per population:
##      Baoule      Borgou      BPN      Charolais      Holstein      Jersey
##      39          69          51          59          52          41
## Lagunaire      Limousin      MaineAnjou      Mtbeliard      NDama      Normand
##      34          48          46          47          43          56
## Parthenais      Somba      Vosgienne      ZChoa      ZMbororo      Zpeul
```

```
##           57           52           49           64           56           67
##
## # Percentage of missing data:
## [1] 0
```

4.8 Exporting data

The `genind` class tends to become a standard in population genetics packages. As of *adegenet* 2.0.0, export functions towards *hierfstat* have been removed, as the package now uses `genind` objects as a native class. Similarly, export towards the package *genetics* have been removed, as *adegenet* now relies on *pegas* for basic population genetics.

A generic way to export data is to produce a `data.frame` of genotypes coded by character strings. This is done by `genind2df`:

```
obj <- genind2df(nancycats)
obj[1:5,1:5]

##      pop   fca8  fca23  fca43  fca45
## cat.1 P01   <NA> 136146 139139 116120
## cat.2 P01   <NA> 146146 139145 120126
## cat.3 P01 135143 136146 141141 116116
## cat.4 P01 133135 138138 139141 116126
## cat.5 P01 133135 140146 141145 126126
```

This function is flexible; for instance, one can separate alleles by any character string:

```
genind2df(nancycats, sep="|")[1:5,1:5]

##      pop   fca8  fca23  fca43  fca45
## cat.1 P01   <NA> 136|146 139|139 116|120
## cat.2 P01   <NA> 146|146 139|145 120|126
## cat.3 P01 135|143 136|146 141|141 116|116
## cat.4 P01 133|135 138|138 139|141 116|126
## cat.5 P01 133|135 140|146 141|145 126|126
```

Note that tabulations can be obtained as follows using `'\t'` character.

5 Basics of data analysis

5.1 Manipulating the data

Data manipulation is meant to be particularly flexible in *adeigenet*. First, as **genind** and **genpop** objects are basically formed by a data matrix (the **@tab** slot), it is natural to subset these objects like it is done with a matrix. The **[** operator does this, forming a new object with the retained genotypes/populations and alleles:

```
data(microbov)
toto <- genind2genpop(microbov)

##
## Converting data from a genind to a genpop object...
##
## ...done.

toto

## /// GENPOP OBJECT ///
##
## // 15 populations; 30 loci; 373 alleles; size: 88.9 Kb
##
## // Basic content
## @tab: 15 x 373 matrix of alleles counts
## @loc.n.all: number of alleles per locus (range: 5-22)
## @loc.fac: locus factor for the 373 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: genind2genpop(x = microbov)
##
## // Optional content
## @other: a list containing: coun breed spe

popNames(toto)

## [1] "Borgou" "Zebu" "Lagunaire"
## [4] "NDama" "Somba" "Aubrac"
## [7] "Bazadais" "BlondeAquitaine" "BretPieNoire"
## [10] "Charolais" "Gascon" "Limousin"
## [13] "MaineAnjou" "Montbeliard" "Salers"

titi <- toto[1:3,]
popNames(titi)

## [1] "Borgou" "Zebu" "Lagunaire"
```

The object `toto` has been subsetted, keeping only the first three populations. Of course, any subsetting available for a matrix can be used with `genind` and `genpop` objects. In addition, we can subset loci directly using indices or logicals, in which case they refer to the output of `locNames`:

```
nAll(titi)

##  INRA63  INRA5  ETH225  ILSTS5  HEL5  HEL1  INRA35  ETH152  INRA23
##      9      7      12      5      11      9      7      12      13
##  ETH10  HEL9  CSSM66  INRA32  ETH3  BM2113  BM1824  HEL13  INRA37
##      9     13     16     14     14     14     10     10     19
##  BM1818  ILSTS6  MM12  CSRM60  ETH185  HAUT24  HAUT27  TGLA227  TGLA126
##     11     13     17     12     16     13     12     15     8
## TGLA122  TGLA53  SPS115
##     22     21     9

tata <- titi[,loc=c(1,3)]
tata

## /// GENPOP OBJECT ///////////
##
## // 3 populations; 2 loci; 21 alleles; size: 18.1 Kb
##
## // Basic content
## @tab: 3 x 21 matrix of alleles counts
## @loc.n.all: number of alleles per locus (range: 9-12)
## @loc.fac: locus factor for the 21 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: .local(x = x, i = i, j = j, loc = ..1, drop = drop)
##
## // Optional content
## @other: a list containing: coun breed spe

nAll(tata)

## INRA63 ETH225
##      9      12
```

Alternatively, one can subset loci using their explicit name:

```
locNames(titi)

## [1] "INRA63" "INRA5" "ETH225" "ILSTS5" "HEL5" "HEL1" "INRA35"
## [8] "ETH152" "INRA23" "ETH10" "HEL9" "CSSM66" "INRA32" "ETH3"
```

```
## [15] "BM2113" "BM1824" "HEL13" "INRA37" "BM1818" "ILSTS6" "MM12"
## [22] "CSRM60" "ETH185" "HAUT24" "HAUT27" "TGLA227" "TGLA126" "TGLA122"
## [29] "TGLA53" "SPS115"

hel5 <- titi[,loc="HEL5"]
hel5

## /// GENPOP OBJECT ///////////
##
## // 3 populations; 1 locus; 11 alleles; size: 16.5 Kb
##
## // Basic content
## @tab: 3 x 11 matrix of alleles counts
## @loc.n.all: number of alleles per locus (range: 11-11)
## @loc.fac: locus factor for the 11 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: .local(x = x, i = i, j = j, loc = "HEL5", drop = drop)
##
## // Optional content
## @other: a list containing: coun breed spe

locNames(hel5)

## [1] "HEL5"
```

To simplify the task of separating data by marker systematically, the function `seploc` can be used. It returns a list of objects (optionnaly, of data matrices), each corresponding to a marker:

```
data(nancycats)
sepCats <- seploc(nancycats)
class(sepCats)

## [1] "list"

names(sepCats)

## [1] "fca8" "fca23" "fca43" "fca45" "fca77" "fca78" "fca90" "fca96" "fca37"

sepCats$fca45

## /// GENIND OBJECT ///////////
##
## // 237 individuals; 1 locus; 9 alleles; size: 31.3 Kb
```

```
##
## // Basic content
## @tab: 237 x 9 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 9-9)
## @loc.fac: locus factor for the 9 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: .local(x = x)
##
## // Optional content
## @pop: population of each individual (group size range: 9-23)
## @strata: - empty -
## @hierarchy: - empty -
## @other: a list containing: xy

identical(tab(sepCats$fca45), tab(nancycats[,loc="fca45"]))

## [1] TRUE
```

The object `sepCats$fca45` only contains data of the marker `fca45`.

Following the same idea, `seppop` allows one to separate genotypes in a `genind` object by population. For instance, we can separate genotype of cattles in the dataset `microbov` by breed:

```
data(microbov)
obj <- seppop(microbov)
class(obj)

## [1] "list"

names(obj)

## [1] "Borgou" "Zebu" "Lagunaire"
## [4] "NDama" "Somba" "Aubrac"
## [7] "Bazadais" "BlondeAquitaine" "BretPieNoire"
## [10] "Charolais" "Gascon" "Limousin"
## [13] "MaineAnjou" "Montbeliard" "Salers"

obj$Borgou

## /// GENIND OBJECT ///
##
## // 50 individuals; 30 loci; 373 alleles; size: 135.4 Kb
```

```
##
## // Basic content
## @tab: 50 x 373 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 5-22)
## @loc.fac: locus factor for the 373 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: .local(x = x, i = i, j = j, treatOther = ..1, quiet = ..2, drop = drop)
##
## // Optional content
## @pop: population of each individual (group size range: 50-50)
## @strata: - empty -
## @hierarchy: - empty -
## @other: a list containing: coun breed spe
```

The returned object `obj` is a list of `genind` objects each containing genotypes of a given breed.

A last, rather vicious trick is to separate data by population and by marker. This is easy using `lapply`; one can first separate population then markers, or the contrary. Here, we separate markers inside each breed in `obj`:

```
obj <- lapply(obj, seploc)
names(obj)

## [1] "Borgou" "Zebu" "Lagunaire"
## [4] "NDama" "Somba" "Aubrac"
## [7] "Bazadais" "BlondeAquitaine" "BretPieNoire"
## [10] "Charolais" "Gascon" "Limousin"
## [13] "MaineAnjou" "Montbeliard" "Salers"

class(obj$Borgou)

## [1] "list"

names(obj$Borgou)

## [1] "INRA63" "INRA5" "ETH225" "ILSTS5" "HEL5" "HEL1" "INRA35"
## [8] "ETH152" "INRA23" "ETH10" "HEL9" "CSSM66" "INRA32" "ETH3"
## [15] "BM2113" "BM1824" "HEL13" "INRA37" "BM1818" "ILSTS6" "MM12"
## [22] "CSRM60" "ETH185" "HAUT24" "HAUT27" "TGLA227" "TGLA126" "TGLA122"
## [29] "TGLA53" "SPS115"

obj$Borgou$INRA63
```

```
## /// GENIND OBJECT ///////////
##
## // 50 individuals; 1 locus; 9 alleles; size: 12.9 Kb
##
## // Basic content
##   @tab: 50 x 9 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 9-9)
##   @loc.fac: locus factor for the 9 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: .local(x = x)
##
## // Optional content
##   @pop: population of each individual (group size range: 50-50)
##   @strata: - empty -
##   @hierarchy: - empty -
##   @other: a list containing: coun breed spe
```

For instance, `obj$Borgou$INRA63` contains genotypes of the breed Borgou for the marker INRA63.

Lastly, one may want to pool genotypes in different datasets, but having the same markers, into a single dataset. This is more than just merging the `@tab` components of all datasets, because alleles can differ (they almost always do) and markers are not necessarily sorted the same way. The function `repool` is designed to avoid these problems. It can merge any `genind` provided as arguments as soon as the same markers are used. For instance, it can be used after a `seppop` to retain only some populations:

```
obj <- seppop(microbov)
names(obj)

##   [1] "Borgou"           "Zebu"             "Lagunaire"
##   [4] "NDama"            "Somba"            "Aubrac"
##   [7] "Bazadais"         "BlondeAquitaine" "BretPieNoire"
##  [10] "Charolais"        "Gascon"           "Limousin"
##  [13] "MaineAnjou"       "Montbeliard"      "Salers"

newObj <- repool(obj$Borgou, obj$Charolais)
newObj

## /// GENIND OBJECT ///////////
##
## // 105 individuals; 30 loci; 295 alleles; size: 175.8 Kb
##
```

```
## // Basic content
## @tab: 105 x 295 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 4-17)
## @loc.fac: locus factor for the 295 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: repool(obj$Borgou, obj$Charolais)
##
## // Optional content
## @pop: population of each individual (group size range: 50-55)
## @strata: - empty -
## @hierarchy: - empty -
## @other: - empty -

popNames(newObj)

## [1] "Borgou" "Charolais"
```

Done !

Note that the content of `@other` can be processed during the conversion from `genind` to `genpop` if the argument `process.other` is set to `TRUE`. Only vectors of a length, or matrices with a number of rows matching the number individuals will be processed. The way they are processed is defined by a function passed as the argument `other.action` (defaulting to 'mean'). Let us illustrate this using `sim2pop`:

```
data(sim2pop)
sim2pop

## /// GENIND OBJECT ///////////
##
## // 130 individuals; 20 loci; 241 alleles; size: 184.5 Kb
##
## // Basic content
## @tab: 130 x 241 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 7-17)
## @loc.fac: locus factor for the 241 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: old2new(object = sim2pop)
##
## // Optional content
## @pop: population of each individual (group size range: 30-100)
```

```
##      @strata: - empty -
##      @hierarchy: - empty -
##      @other: a list containing: xy

nInd(sim2pop)

## [1] 130

head(other(sim2pop)$xy)

##           x           y
## [1,] 35.11291 99.595997
## [2,] 22.57033  6.682107
## [3,] 76.99371 51.900514
## [4,] 44.31948 18.037868
## [5,] 94.40902 82.948821
## [6,] 51.29493 25.007193

dim(other(sim2pop)$xy)

## [1] 130  2
```

The component `sim2pop@other$xy` contains spatial coordinates of individuals from 2 populations.

```
other(genind2genpop(sim2pop, process.other=TRUE))

##
## Converting data from a genind to a genpop object...
##
## ...done.
## $xy
##           x           y
## P01 58.43405 48.37065
## P02 19.06501 61.00044
```

In this case, numeric vectors with a length corresponding to the number of individuals will be averaged per groups; note that any other function than `mean` can be used by providing any function to the argument `other.action`. Matrices with a number of rows corresponding to the number of individuals are processed similarly.

5.2 Using summaries

Both `genind` and `genpop` objects have a summary providing basic information about data. Informations are both printed and invisibly returned as a list.


```

toto <- summary(nancycats)

##
## # Total number of genotypes: 237
##
## # Population sample sizes:
## P01 P02 P03 P04 P05 P06 P07 P08 P09 P10 P11 P12 P13 P14 P15 P16 P17
## 10 22 12 23 15 11 14 10 9 11 20 14 13 17 11 12 13
##
## # Number of alleles per locus:
## fca8 fca23 fca43 fca45 fca77 fca78 fca90 fca96 fca37
## 16 11 10 9 12 8 12 12 18
##
## # Number of alleles per population:
## P01 P02 P03 P04 P05 P06 P07 P08 P09 P10 P11 P12 P13 P14 P15 P16 P17
## 36 53 50 67 48 56 42 54 43 46 70 52 44 61 42 40 35
##
## # Percentage of missing data:
## [1] 2.344116
##
## # Observed heterozygosity:
## fca8 fca23 fca43 fca45 fca77 fca78 fca90
## 0.6682028 0.6666667 0.6793249 0.7083333 0.6329114 0.5654008 0.6497890
## fca96 fca37
## 0.6184211 0.4514768
##
## # Expected heterozygosity:
## fca8 fca23 fca43 fca45 fca77 fca78 fca90
## 0.8657224 0.7928751 0.7953319 0.7603095 0.8702576 0.6884669 0.8157881
## fca96 fca37
## 0.7603493 0.6062686

names(toto)

## [1] "N" "pop.eff" "loc.n.all" "pop.nall" "NA.perc" "Hobs"
## [7] "Hexp"

par(mfrow=c(2,2))

plot(toto$pop.eff, toto$pop.nall, xlab="Colonies sample size",
     ylab="Number of alleles", main="Alleles numbers and sample sizes",
     type="n")
text(toto$pop.eff, toto$pop.nall, lab=names(toto$pop.eff))

barplot(toto$loc.n.all, ylab="Number of alleles",

```

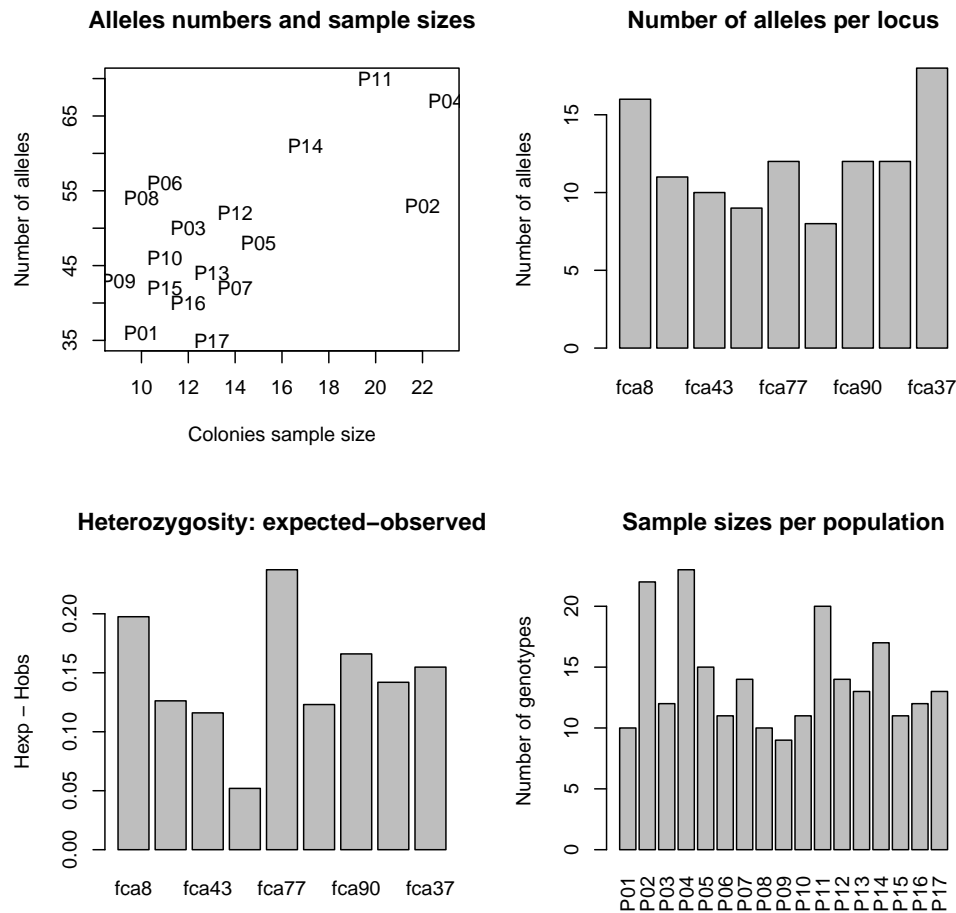
```

    main="Number of alleles per locus")

barplot(toto$Hexp-toto$Hobs, main="Heterozygosity: expected-observed",
        ylab="Hexp - Hobs")

barplot(toto$pop.eff, main="Sample sizes per population",
        ylab="Number of genotypes",las=3)

```



Is mean observed H significantly lower than mean expected H ?

```

bartlett.test(list(toto$Hexp,toto$Hobs))

##
## Bartlett test of homogeneity of variances
##
## data: list(toto$Hexp, toto$Hobs)
## Bartlett's K-squared = 0.046962, df = 1, p-value = 0.8284

t.test(toto$Hexp,toto$Hobs,pair=T,var.equal=TRUE,alter="greater")

```

```
##
## Paired t-test
##
## data: toto$Hexp and toto$Hobs
## t = 8.3294, df = 8, p-value = 1.631e-05
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
##  0.1134779      Inf
## sample estimates:
## mean of the differences
##                0.1460936
```

Yes, it is.

5.3 Testing for Hardy-Weinberg equilibrium

As of version 2.0.0, *adegenet* is designed to work alongside a number of other packages, especially *pegas* and *hierfstat* for a number of classical population genetics methods. The former function `HWE.test.genind` has consequently been removed, and replaced by *pegas*'s `hw.test`, which performs one test per locus:

```
library(pegas)
data(nancycats)
cats.hwt <- hw.test(nancycats, B=0)
cats.hwt

##           chi^2  df  Pr(chi^2 >)
## fca8  395.80006 120 0.000000e+00
## fca23 239.34221  55 0.000000e+00
## fca43 434.33397  45 0.000000e+00
## fca45  66.11849  36 1.622163e-03
## fca77 270.52066  66 0.000000e+00
## fca78 402.80002  28 0.000000e+00
## fca90 217.19836  66 0.000000e+00
## fca96 193.36764  66 1.965095e-14
## fca37 291.00731 153 1.209777e-10
```

Note that `B=0` is used for the parametric version; larger numbers will indicate the number of permutations to use for a Monte-Carlo version.

5.4 Measuring and testing population structure (a.k.a F statistics)

Population structure is traditionally measured and tested using F statistics, in particular F_{st} . As of version 2.0.0, *adegenet* relies on *hierfstat* and *pegas* for most F statistics. Note that some of these features might be part of the current devel version of *hierfstat*, which can be installed by:

```
library(devtools)
install_github("jgx65/hierfstat")
library("hierfstat")
```

Can we find any population structure in the cat colonies from Nancy? The basic F statistics are provided by:

```
fstat(nancycats)

## Error in eval(expr, envir, enclos): could not find function "fstat"
```

This table provides the three F statistics F_{st} (pop/total), F_{it} (Ind/total), and F_{is} (ind/pop). These are overall measures which take into account all genotypes and all loci.

For more detail, *pegas* provides estimates by locus:

```
library(pegas)
Fst(as.loci(nancycats))
```

##		Fit	Fst	Fis
##	fca8	0.2447420	0.10146648	0.159454807
##	fca23	0.1646295	0.06746762	0.104191391
##	fca43	0.1514487	0.06893755	0.088620458
##	fca45	0.1010807	0.09792456	0.003498722
##	fca77	0.2790495	0.10036588	0.198618075
##	fca78	0.1842490	0.07025915	0.122603911
##	fca90	0.2098744	0.09168833	0.130116240
##	fca96	0.2034755	0.10744024	0.107595351
##	fca37	0.2604033	0.06985321	0.204860244

Are these values significant? This question can be addressed using the G-statistic test [?]; it is implemented for *genind* objects and produces a *randtest* object (package *ade4*).

```
Gtest <- gstat.randtest(nancycats,nsim=99)

## Error in eval(expr, envir, enclos): could not find function
" gstat.randtest"

Gtest

## NULL

plot(Gtest)

## Error in plot.window(...): need finite 'xlim' values
```

Finally, pairwise Fst is frequently used as a measure of distance between populations. The function `pairwise.fst` computes Nei's estimator [11] of pairwise Fst , defined as:

$$Fst(A, B) = \frac{H_t - (n_A H_s(A) + n_B H_s(B)) / (n_A + n_B)}{H_t}$$

where A and B refer to the two populations of sample size n_A and n_B and respective expected heterozygosity $H_s(A)$ and $H_s(B)$, and H_t is the expected heterozygosity in the whole dataset. For a given locus, expected heterozygosity is computed as $1 - \sum p_i^2$, where p_i is the frequency of the i th allele, and the \sum represents summation over all alleles. For multilocus data, the heterozygosity is simply averaged over all loci. These computations are achieved for all pairs of populations by the function `pairwise.fst`; we illustrate this on a subset of individuals of `nancycats` (computations for the whole dataset would take a few tens of seconds):

```
matFst <- pairwise.fst(nancycats[1:50,])

## Error in eval(expr, envir, enclos): could not find function
## "pairwise.fst"

matFst
```

```
##           1           2           3
## 2 0.08018500
## 3 0.07140847 0.08200880
## 4 0.08163151 0.06512457 0.04131227
```

The resulting matrix is Euclidean when there are no missing values:

```
is.euclid(matFst)

## [1] TRUE
```

It can therefore be used in a Principal Coordinate Analysis (which requires Euclideanity), used to build trees, etc.

5.5 Estimating inbreeding

Inbreeding refers to an excess of homozygosity in a given individual due to the mating of genetically related parents. This excess of homozygosity is due to the fact that there are non-negligible chances of inheriting two identical alleles from a recent common ancestor. Inbreeding can be associated to a loss of fitness leading to "*inbreeding depression*". Typically, loss of fitness is caused by recessive deleterious alleles which have usually low frequency in the population, but for which inbred individuals are more likely to be homozygotes.

The inbreeding coefficient F is defined as the probability that at a given locus, two identical alleles have been inherited from a common ancestor. In the absence of inbreeding, the probability of being homozygote at one loci is (for diploid individuals) simply $\sum_i p_i^2$ where i indexes the alleles and p_i is the frequency of allele i . This can be generalized incorporating F as:

$$p(\text{homozygote}) = F + (1 - F) \sum_i p_i^2$$

and even more generally, for any ploidy π :

$$p(\text{homozygote}) = F + (1 - F) \sum_i p_i^\pi$$

This therefore allows for computing the likelihood of a given state (homozygote/heterozygote) in a given genotype (log-likelihood are summed across loci for more than one marker).

This estimation is achieved by `inbreeding`. Depending on the value of the argument `res.type`, the function returns a sample from the likelihood function (`res.type='sample'`) or the likelihood function itself, as a R function (`res.type='function'`). While likelihood functions are quickly obtained and easy to display graphically, sampling from the distributions is more computer intensive but useful to derive summary statistics of the distributions. Here, we illustrate `inbreeding` using the `microbov` dataset, which contains cattle breeds genotypes for 30 microsatellites; to focus on breed Salers only, we use `seppop`:

```

data(microbov)
sal <- seppop(microbov)$Salers
sal

## /// GENIND OBJECT ///////////
##
## // 50 individuals; 30 loci; 373 alleles; size: 135.4 Kb
##
## // Basic content
##   @tab: 50 x 373 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 5-22)
##   @loc.fac: locus factor for the 373 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: .local(x = x, i = i, j = j, treatOther = ..1, quiet = ..2, drop = drop)
##
## // Optional content
##   @pop: population of each individual (group size range: 50-50)
##   @strata: - empty -
##   @hierarchy: - empty -
##   @other: a list containing: coun breed spe

```

We first compute the mean inbreeding for each individual, and plot the resulting distribution:

```

temp <- inbreeding(sal, N=100)
class(temp)

## [1] "list"

head(names(temp))

## [1] "FRBTSAL9087" "FRBTSAL9088" "FRBTSAL9089" "FRBTSAL9090" "FRBTSAL9091"
## [6] "FRBTSAL9093"

head(temp[[1]],20)

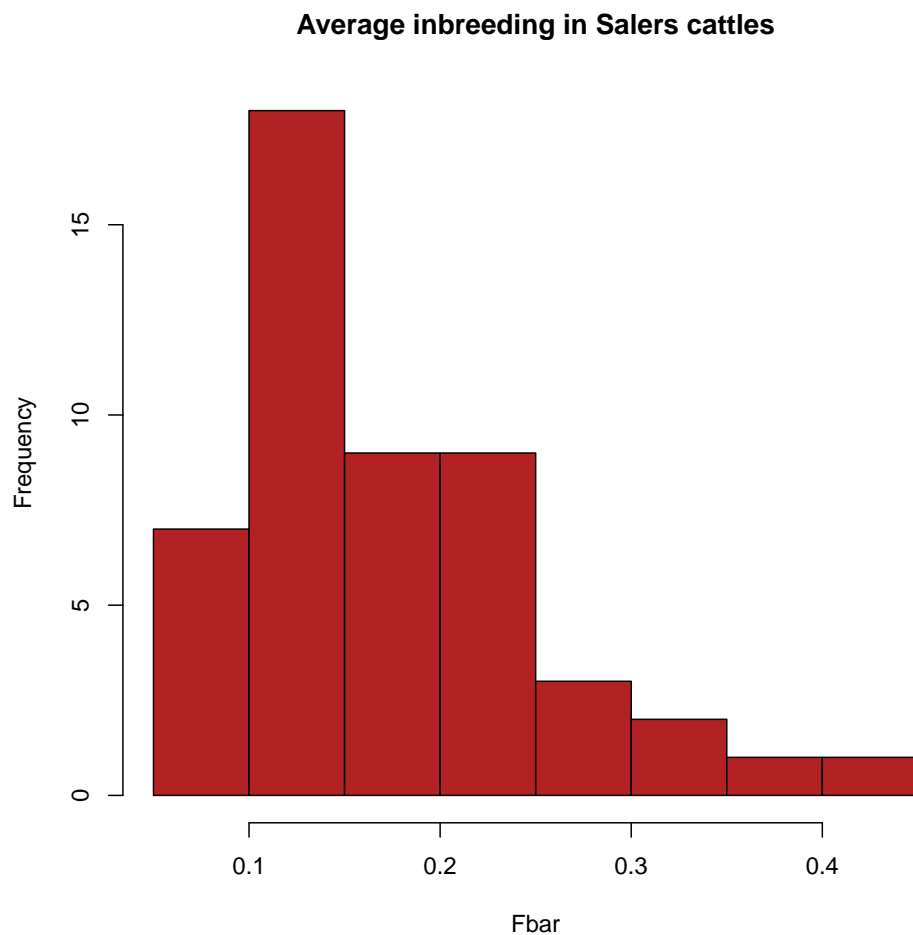
## [1] 0.17844511 0.12481989 0.07049219 0.06123628 0.11833790 0.10981626
## [7] 0.11134785 0.09905051 0.08405919 0.20097015 0.25144776 0.10228146
## [13] 0.22941423 0.10276143 0.07951652 0.32500351 0.08339947 0.19946721
## [19] 0.07212448 0.09216534

```

`temp` is a list of values sampled from the likelihood distribution of each individual; means values are obtained for all individuals using `sapply`:

```
Fbar <- sapply(temp, mean)
```

```
hist(Fbar, col="firebrick", main="Average inbreeding in Salers cattles")
```



We can see that some individuals (actually, a single one) have higher inbreeding (>0.4). We can recompute inbreeding for this individual, asking for the likelihood function to be returned:

```
which(Fbar>0.4)
```

```
## FRBTSAL9266
```

```
##          37
```

```
F <- inbreeding(sal, res.type="function")[which(Fbar>0.4)]
```

```
F
```

```
## $FRBTSAL9266
```

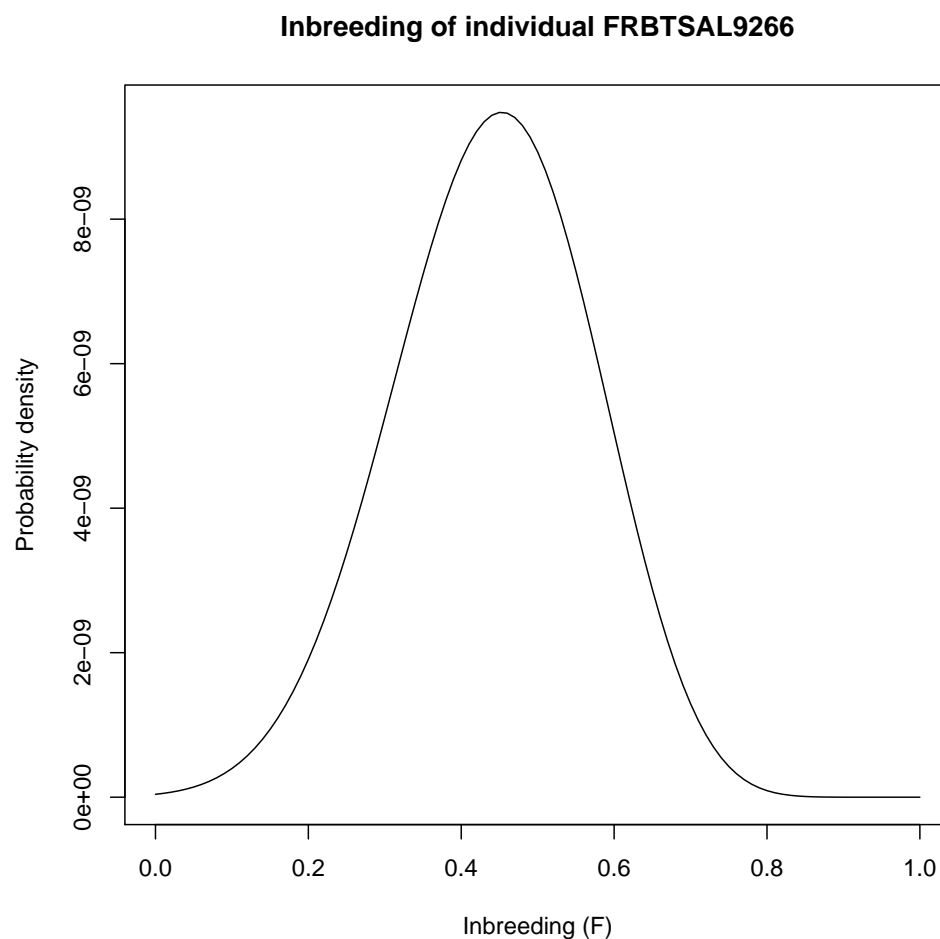
```
## function (F)
```



```
## {
##   args <- lapply(as.list(match.call())[-1L], eval, parent.frame())
##   names <- if (is.null(names(args)))
##     character(length(args))
##   else names(args)
##   dovec <- names %in% vectorize.args
##   do.call("mapply", c(FUN = FUN, args[dovec], MoreArgs = list(args[!dovec]),
##     SIMPLIFY = SIMPLIFY, USE.NAMES = USE.NAMES))
## }
## <environment: 0xafb4dd8>
```

The output object `F` can seem a bit cryptic: it is an function embedded within a hidden environment. This does not matter, however, since it is easily represented:

```
plot(F$FRBTSAL9266, main=paste("Inbreeding of individual", names(F)),
     xlab="Inbreeding (F)", ylab="Probability density")
```



Indeed, this individual shows subsequent inbreeding, with about 50% chances of being homozygote through inheritance from a common ancestor of its parents.

6 Multivariate analysis

6.1 General overview

Multivariate analysis consists in summarising a strongly multivariate information into a few synthetic variables. In genetics, such approaches are useful to get a simplified picture of the genetic diversity observed amongst individuals or populations. A review of multivariate analysis in population genetics can be found in [7]. Here, we aim at providing an overview of some applications using methods implemented in *ade4* and *adeigenet*.

Useful functions include:

- **scaleGen** (*adeigenet*): centre/scale allele frequencies and replaces missing data; useful, among other things, before running a principal component analysis (PCA).
- **dudi.pca** (*ade4*): implements PCA; can be used on transformed allele frequencies of individuals or populations.
- **dudi.ca** (*ade4*): implements Correspondance Analysis (CA); can be used on raw allele counts of populations (@tab slot in **genpop** objects).
- **dist.genpop** (*adeigenet*): implements 5 pairwise genetic distances between populations
- **pairwise.fst** (*adeigenet*): implements pairwise F_{ST} , which is also a Euclidean distance between populations.
- **dist** (*stats*): computes pairwise distances between multivariate observations; can be used on raw or transformed allele frequencies.
- **dudi.pco** (*ade4*): implements Principal Coordinates Analysis (PCoA); this method finds synthetic variables which summarize a Euclidean distance matrix as best as possible; can be used on outputs of **dist**, **dist.genpop**, and **pairwise.fst**.
- **is.euclid** (*ade4*): tests whether a distance matrix is Euclidean, which is a pre-requisite of PCoA.
- **cailliez** (*ade4*): renders a non-Euclidean distance matrix Euclidean by adding a constant to all entries.
- **dapc** (*adeigenet*): implements the Discriminant Analysis of Principal Components (DAPC [4]), a powerful method for the analysis of population genetic structures; see dedicated tutorial (*dapc*).
- **sPCA** (*adeigenet*): implements the spatial Principal Component Analysis (sPCA [5]), a method for the analysis of spatial genetic structures; see dedicated tutorial (*dapc*).
- **glPca** (*adeigenet*): implements PCA for genome-wide SNP data stored as **genlight** objects; see dedicated tutorial (*genomics*).

Besides the procedures themselves, graphic functions are also often of the utmost importance; these include:

- `scatter` (*ade4*, *adegenet*): generic function to display multivariate analyses; in practice, the most useful application for genetic data is the one implemented in *adegenet* for DAPC results.
- `s.label` (*ade4*): function used for basic display of principal components.
- `loadingplot` (*adegenet*): function used to display the loadings (i.e., contribution to a given structure) of alleles for a given principal component; annotates and returns the most contributing alleles.
- `s.class` (*ade4*): displays two quantitative variables with known groups of observations, using inertia ellipses for the groups; useful to represent principal components when groups are known.
- `s.chull` (*ade4*): same as `s.class`, except convex polygons are used rather than ellipses.
- `s.value` (*ade4*): graphical display of a quantitative variable distributed over a two-dimensional space; useful to map principal components or allele frequencies over a geographic area.
- `colorplot` (*adegenet*): graphical display of 1 to 3 quantitative variables distributed over a two-dimensional space; useful for combined representations of principal components over a geographic area. Can also be used to produce color versions of traditional scatterplots.
- `transp` (*adegenet*): auxiliary function making colors transparent.
- `num2col` (*adegenet*): auxiliary function transforming a quantitative variable into colors using a given palette.
- `assignplot` (*adegenet*): specific plot of group membership probabilities for DAPC; see dedicated tutorial (*dapc*).
- `compoplot` (*adegenet*): specific 'STRUCTURE-like' plot of group membership probabilities for DAPC; see dedicated tutorial (*dapc*).
- `add.scatter` (*ade4*): add inset plots to an existing figure.
- `add.scatter.eig` (*ade4*): specific application of `add.scatter` to add barplots of eigenvalues to an existing figure.

In the sections below, we briefly illustrate how these tools can be combined to extract information from genetic data.

6.2 Performing a Principal Component Analysis on `genind` objects

The tables contained in `genind` objects can be submitted to a Principal Component Analysis (PCA) to seek a summary of the genetic diversity among the sampled individuals. Such analysis is straightforward using *ade4* to prepare data and *ade4* for the analysis *per se*. One has first to replace missing data (NAs) and transform the allele frequencies in an appropriate way. These operations are achieved by `scaleGen`. NAs are replaced by the mean allele frequency; different scaling options are available (argument `method`), but in general centring is sufficient since allele frequencies have inherently comparable variances.

```
data(microbov)
sum(is.na(microbov$tab))

## [1] 6325
```

There are 6325 missing data. They will all be replaced by `scaleGen`:

```
X <- scaleGen(microbov, NA.method="mean")
class(X)

## [1] "matrix"

dim(X)

## [1] 704 373

X[1:5,1:5]

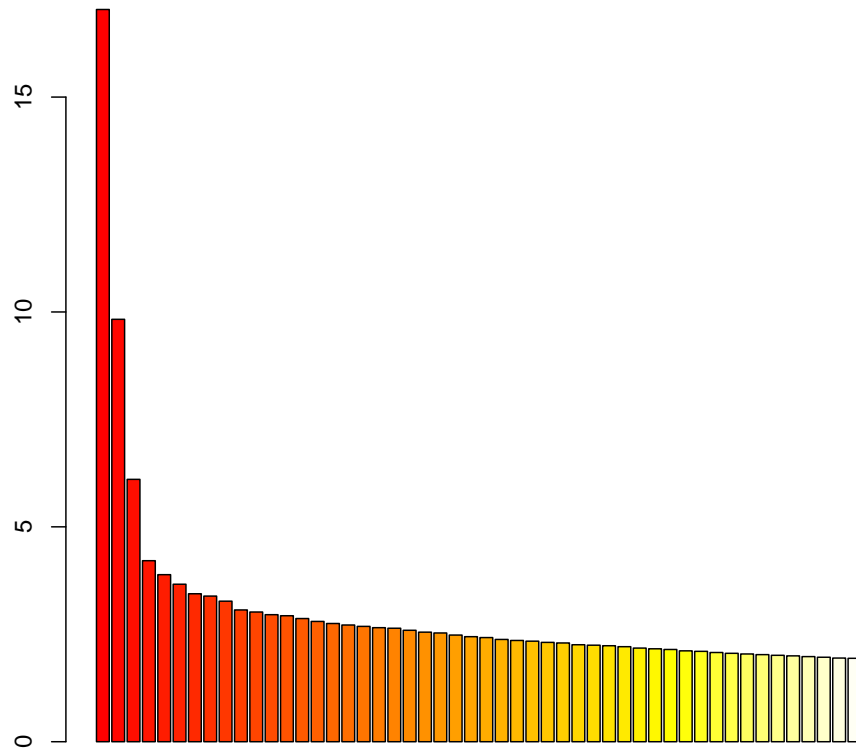
##           INRA63.167  INRA63.171  INRA63.173  INRA63.175  INRA63.177
## AFBIBOR9503 -0.03801312 -0.05379728  -0.101009  -1.061893  -0.8769237
## AFBIBOR9504 -0.03801312 -0.05379728  -0.101009  -1.061893  -0.8769237
## AFBIBOR9505 -0.03801312 -0.05379728  -0.101009  -1.061893   0.5498659
## AFBIBOR9506 -0.03801312 -0.05379728  -0.101009  -1.061893  -0.8769237
## AFBIBOR9507 -0.03801312 -0.05379728  -0.101009  -1.061893   0.5498659
```

Note that alternatively, we could have used `na.replace` to replace missing data, and then left the centring/scaling to `dudi.pca`.

The analysis can now be performed. We disable the scaling in `dudi.pca`, which would erase the scaling choice made earlier in `scaleGen`. Note: in practice, retained axes can be chosen interactively by removing the arguments `scannf=FALSE,nf=3`.

```
pca1 <- dudi.pca(X,cent=FALSE,scale=FALSE,scannf=FALSE,nf=3)
barplot(pca1$eig[1:50],main="PCA eigenvalues", col=heat.colors(50))
```

PCA eigenvalues



```
pca1

## Duality diagramm
## class: pca dudi
## $call: dudi.pca(df = X, center = FALSE, scale = FALSE, scannf = FALSE,
##      nf = 3)
##
## $nf: 3 axis-components saved
## $rank: 343
## eigen values: 17.04 9.829 6.105 4.212 3.887 ...
##   vector length mode   content
## 1 $cw      373      numeric column weights
## 2 $lw      704      numeric row weights
## 3 $eig     343      numeric eigen values
##
##   data.frame nrow ncol content
## 1 $tab        704   373 modified array
## 2 $li         704    3   row coordinates
```

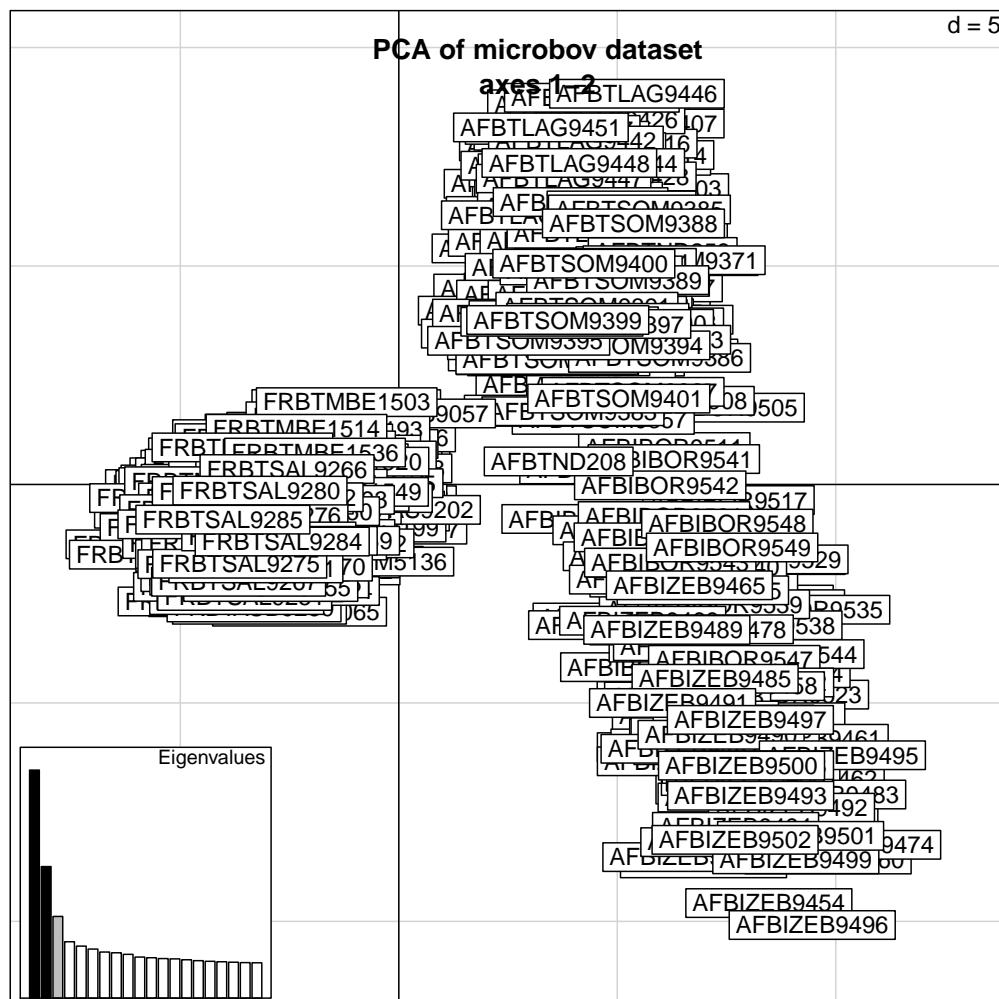
```
## 3 $l1      704 3    row normed scores
## 4 $co      373 3    column coordinates
## 5 $c1      373 3    column normed scores
## other elements: cent norm
```

The output object `pca1` is a list containing various information; of particular interest are:

- `$eig`: the eigenvalues of the analysis, indicating the amount of variance represented by each principal component (PC).
- `$li`: the principal components of the analysis; these are the synthetic variables summarizing the genetic diversity, usually visualized using scatterplots.
- `$c1`: the allele loadings, used to compute linear combinations forming the PCs; squared, they represent the contribution to each PCs.

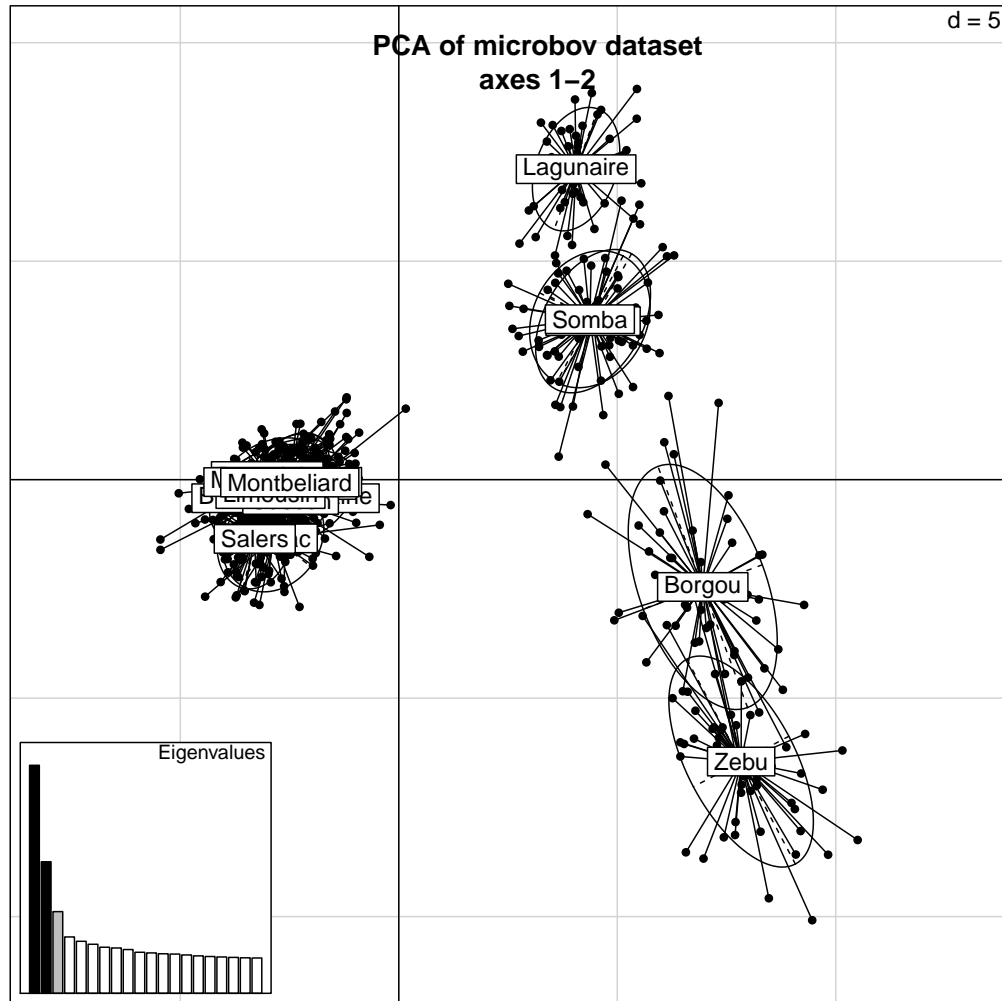
The basic scatterplot for this analysis can be obtained by:

```
s.label(pca1$l1)
title("PCA of microbov dataset\naxes 1-2")
add.scatter.eig(pca1$eig[1:20], 3,1,2)
```



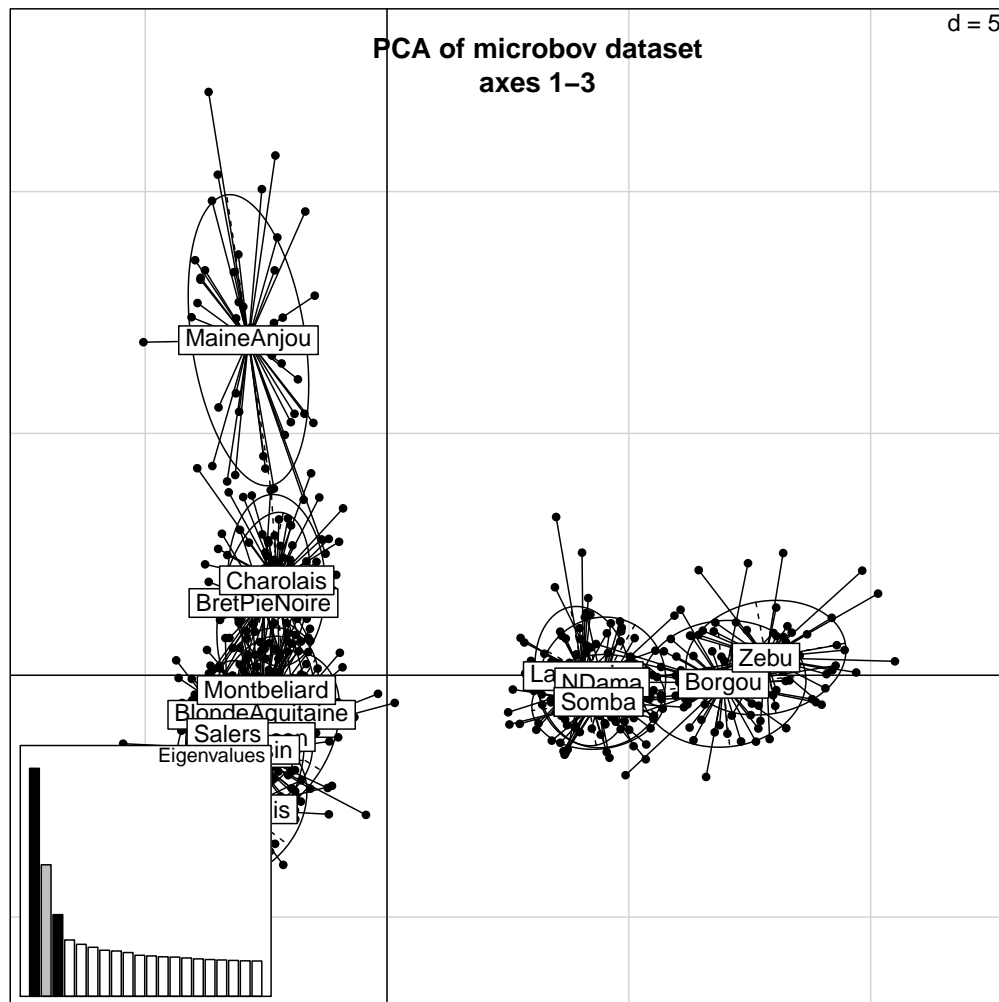
However, this figure can largely be improved. First, we can use `s.class` to represent both the genotypes and inertia ellipses for populations.

```
s.class(pca1$li, pop(microbov))
title("PCA of microbov dataset\naxes 1-2")
add.scatter.eig(pca1$eig[1:20], 3,1,2)
```



This plane shows that the main structuring is between African and French breeds, the second structure reflecting genetic diversity among African breeds. The third axis reflects the diversity among French breeds:

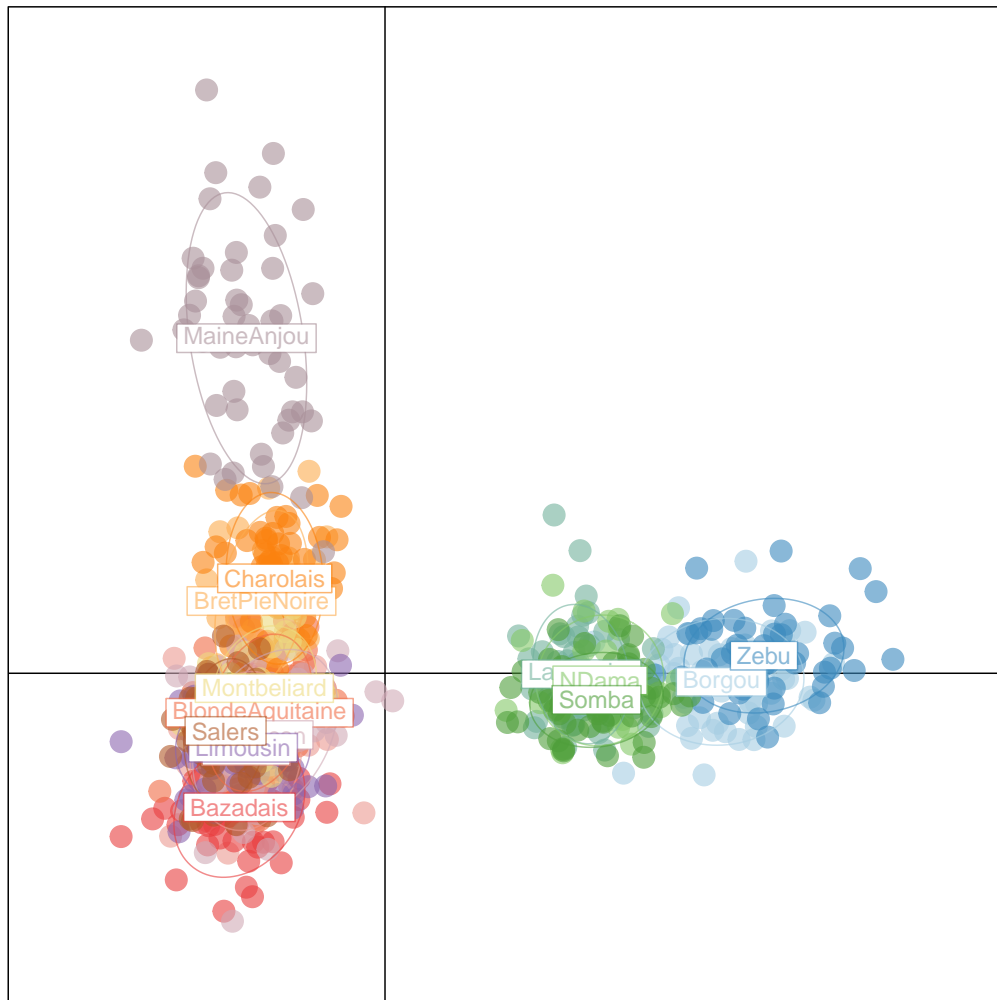
```
s.class(pca1$li, pop(microbov), xax=1, yax=3, sub="PCA 1-3", csub=2)
title("PCA of microbov dataset\naxes 1-3")
add.scatter.eig(pca1$eig[1:20], nf=3, xax=1, yax=3)
```



Overall, all breeds seem well differentiated.

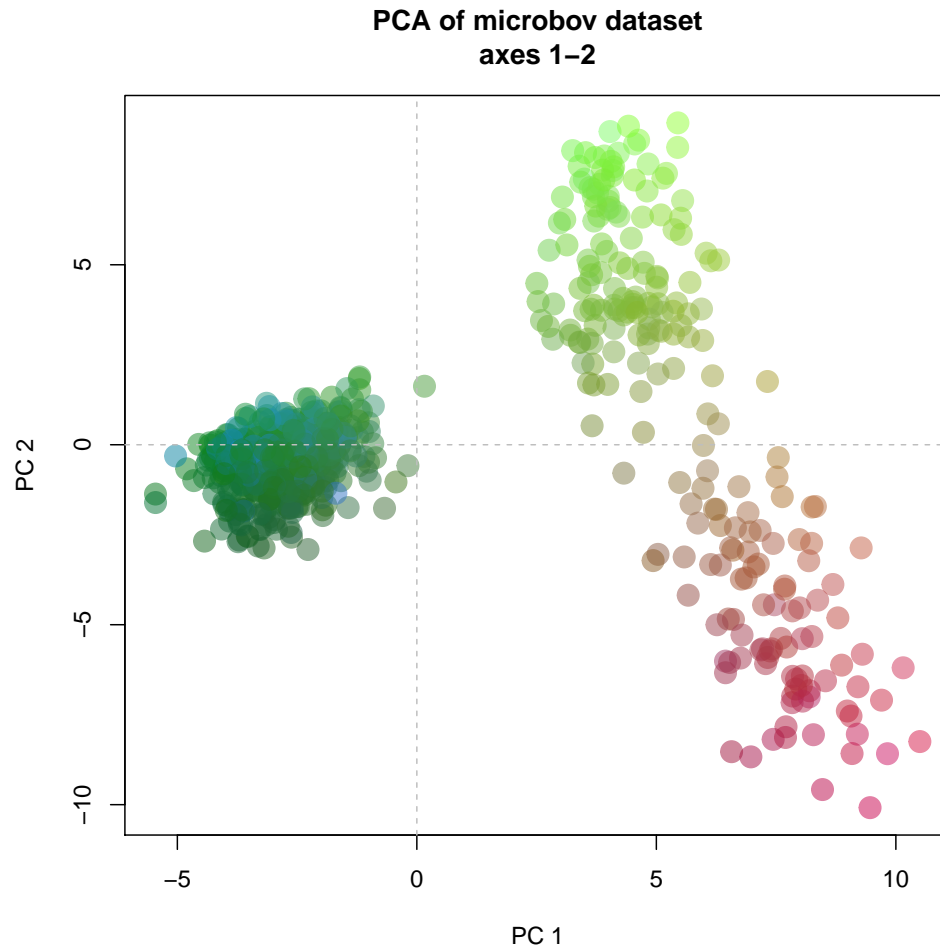
However, we can yet improve these scatterplots, which are fortunately easy to customize. For instance, we can remove the grid, choose different colors for the groups, use larger dots and transparency to better assess the density of points, and remove internal segments of the ellipses:

```
col <- funky(15)
s.class(pca1$li, pop(microbov), xax=1, yax=3, col=transp(col,.6), axesell=FALSE,
        cstar=0, cpoint=3, grid=FALSE)
```

Let us now assume that we ignore the group memberships. We can still use color in an informative way. For instance, we can recode the principal components represented in the scatterplot on the RGB scale:

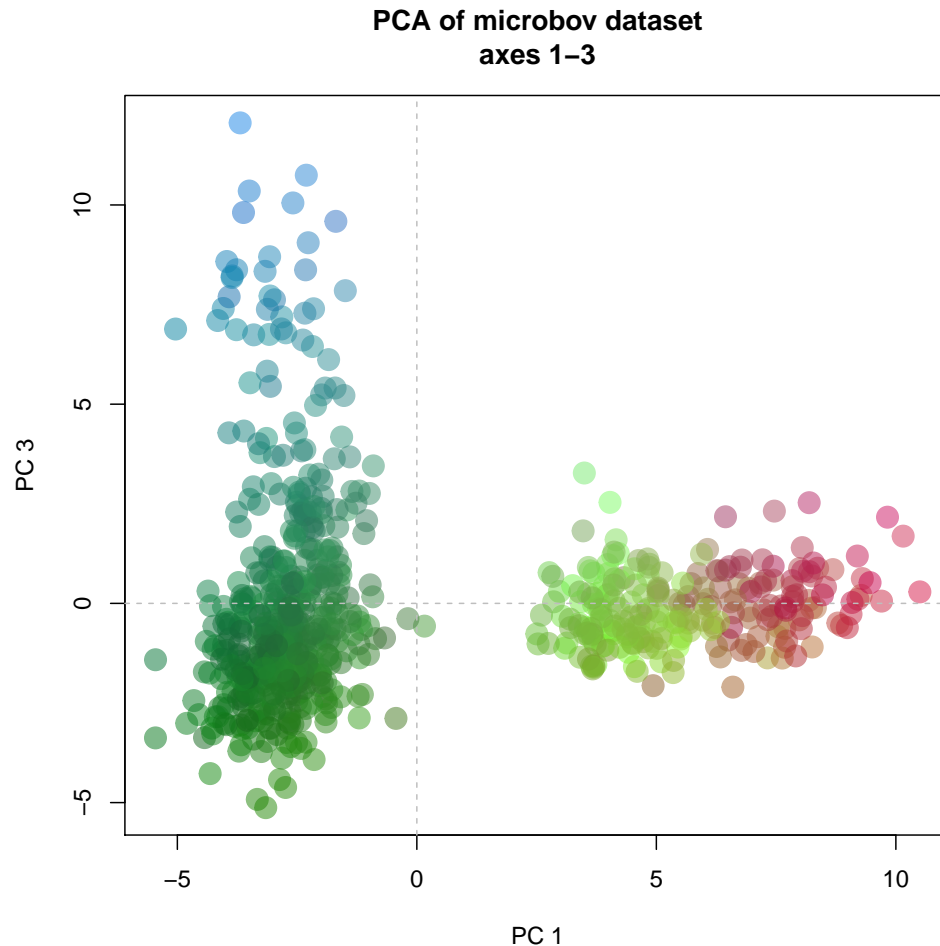
```
colorplot(pca1$li, pca1$li, transp=TRUE, cex=3, xlab="PC 1", ylab="PC 2")
title("PCA of microbov dataset\naxes 1-2")
abline(v=0,h=0,col="grey", lty=2)
```



Colors are based on the first three PCs of the PCA, recoded respectively on the red, green, and blue channel. In this figure, the genetic diversity is represented in two complementary ways: by the distances (further away = more genetically different), and by the colors (more different colors = more genetically different).

We can represent the diversity on the third axis similarly:

```
colorplot(pca1$li[c(1,3)], pca1$li, transp=TRUE, cex=3, xlab="PC 1", ylab="PC 3")
title("PCA of microbov dataset\naxes 1-3")
abline(v=0,h=0,col="grey", lty=2)
```



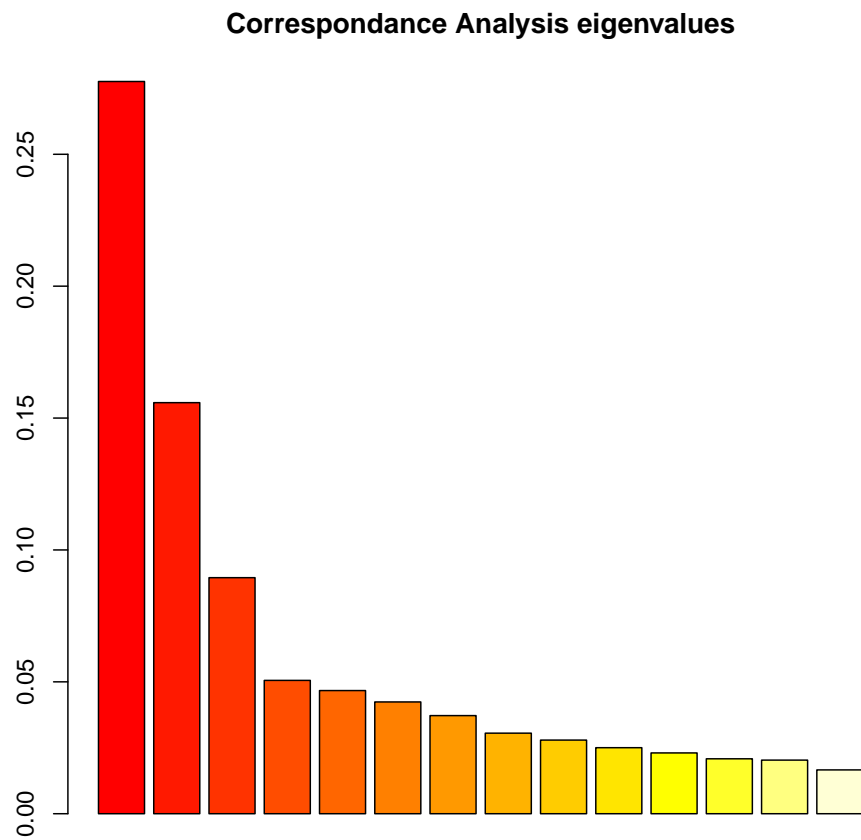
6.3 Performing a Correspondance Analysis on genpop objects

Being contingency tables, the `@tab` slot in `genpop` objects can be submitted to a Correspondance Analysis (CA) to seek a typology of populations. The approach is very similar to the previous one for PCA.

```
data(microbov)
obj <- genind2genpop(microbov)

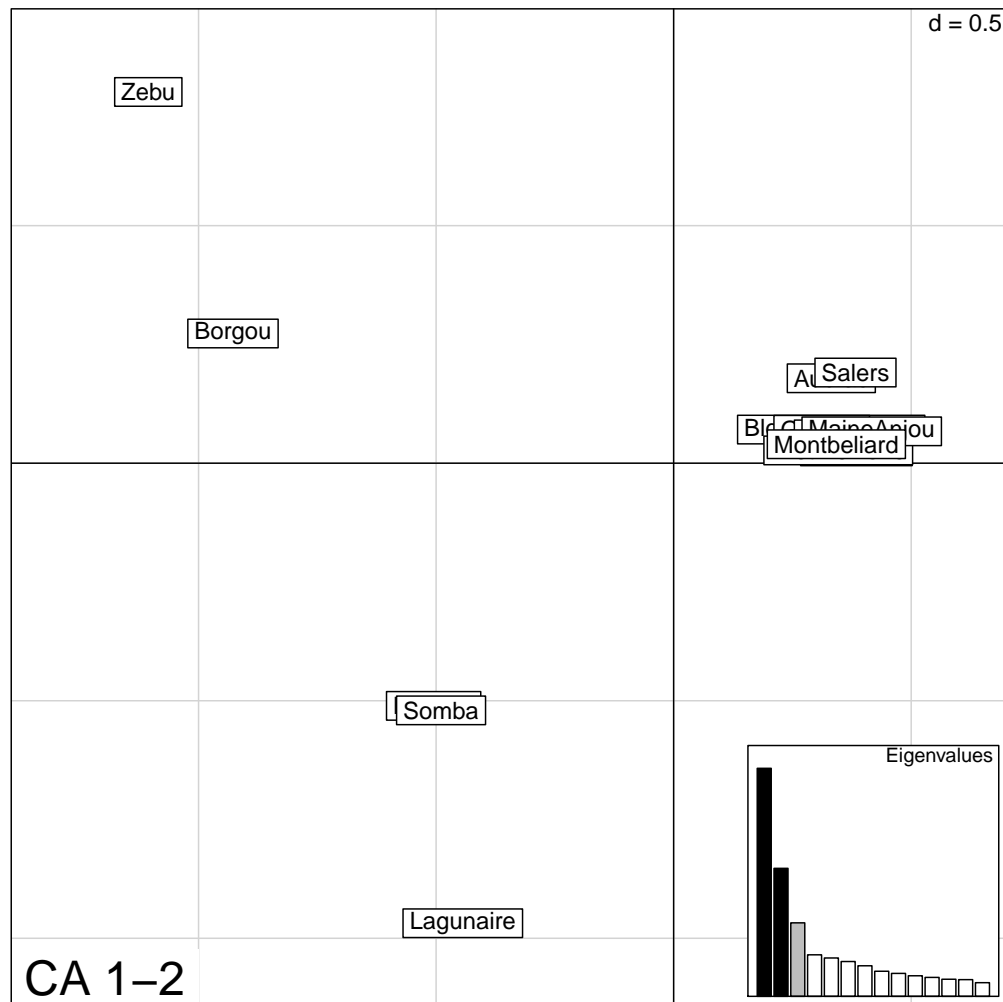
##
## Converting data from a genind to a genpop object...
##
## ...done.

ca1 <- dudi.coa(tab(obj),scannf=FALSE,nf=3)
barplot(ca1$eig,main="Correspondance Analysis eigenvalues",
        col=heat.colors(length(ca1$eig)))
```



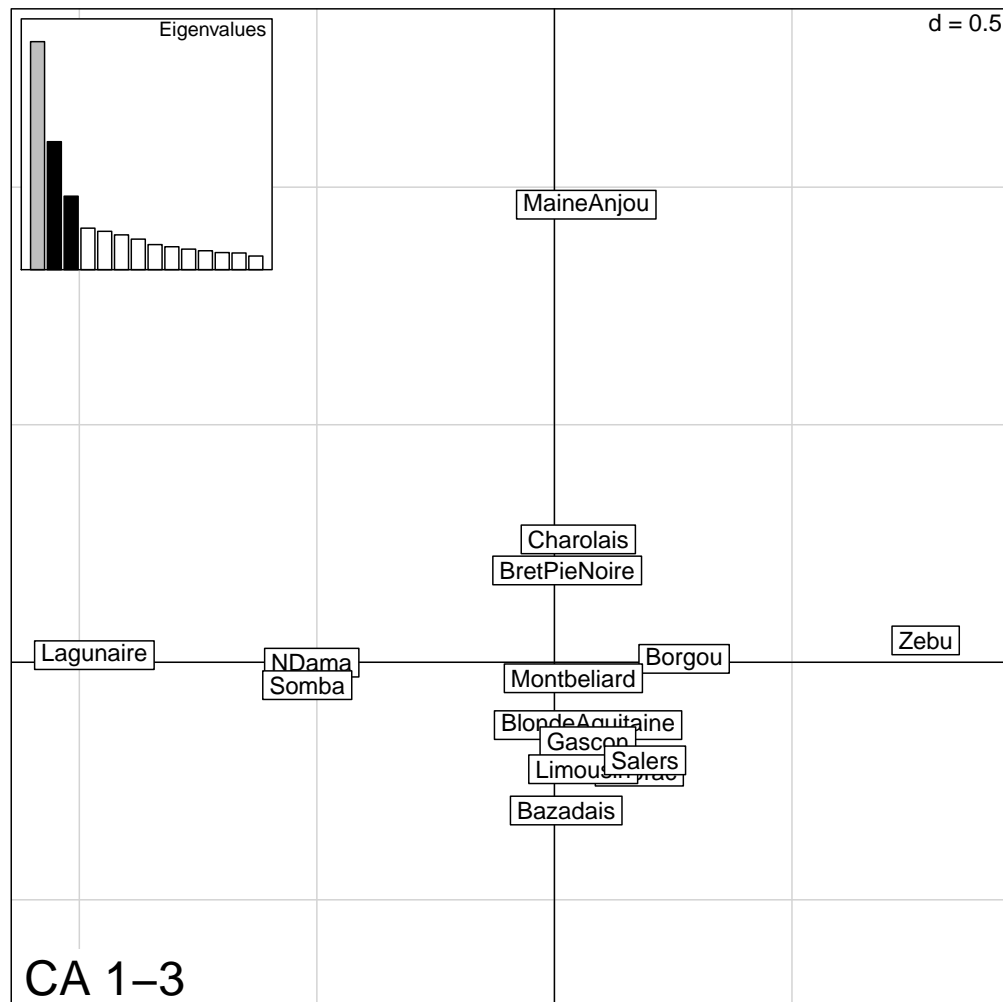
Now we display the resulting typology using a basic scatterplot:

```
s.label(ca1$li, sub="CA 1-2", csub=2)  
add.scatter.eig(ca1$eig, nf=3, xax=1, yax=2, posi="bottomright")
```



The same graph is derived for the axes 2-3:

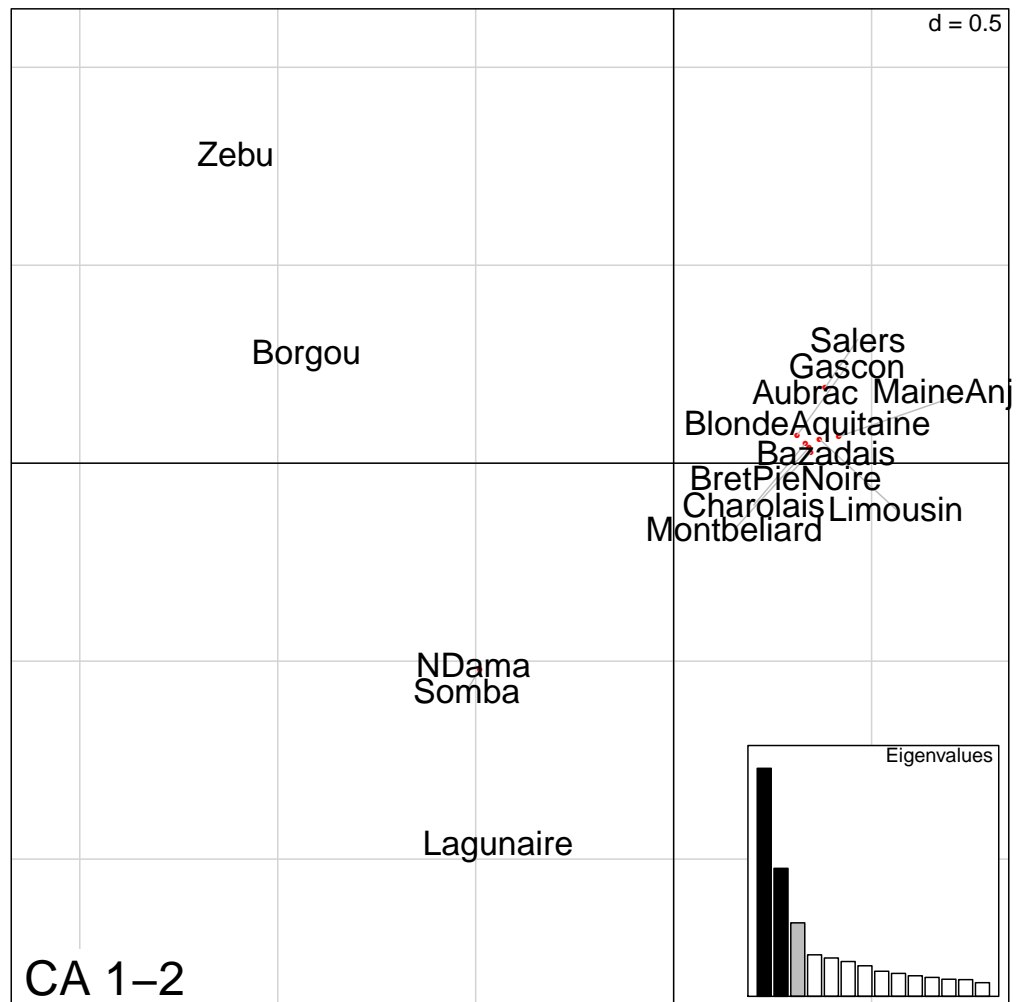
```
s.label(ca1$li,xax=2,yax=3,lab=popNames(obj),sub="CA 1-3",csub=2)
add.scatter.eig(ca1$eig,nf=3,xax=2,yax=3,posit="topleft")
```



As in the PCA above, axes are to be interpreted separately in terms of continental differentiation, and between-breeds diversity. Importantly, as in any analysis carried out at a population level, all information about the diversity within populations is lost in this analysis. See the tutorial on DAPC for an individual-based approach which is nonetheless optimal in terms of group separation (*dapc*).

Note that as an alternative, *wordcloud* can be used to avoid overlaps in labels:

```
library(wordcloud)
set.seed(1)
s.label(ca1$li*1.2, sub="CA 1-2", csub=2, clab=0, cpoint="")
textplot(ca1$li[,1], ca1$li[,2], words=popNames(obj),
         cex=1.4, new=FALSE, xpd=TRUE)
add.scatter.eig(ca1$eig, nf=3, xax=1, yax=2, posi="bottomright")
```



However, only general trends can be interpreted: labels positions are randomised to avoid overlap, so they no longer accurately position populations on the factorial axes.

7 Spatial analysis

The R software probably offers the largest collection of spatial methods among statistical software. Here, we briefly illustrate two methods commonly used in population genetics. Spatial multivariate analysis is covered in a dedicated tutorial; see *spca* tutorial for more information.

7.1 Isolation by distance

7.1.1 Testing isolation by distance

Isolation by distance (IBD) is tested using Mantel test between a matrix of genetic distances and a matrix of geographic distances. It can be tested using individuals as well as populations. This example uses cat colonies from the city of Nancy. We test the correlation between Edwards' distances and Euclidean geographic distances between colonies.

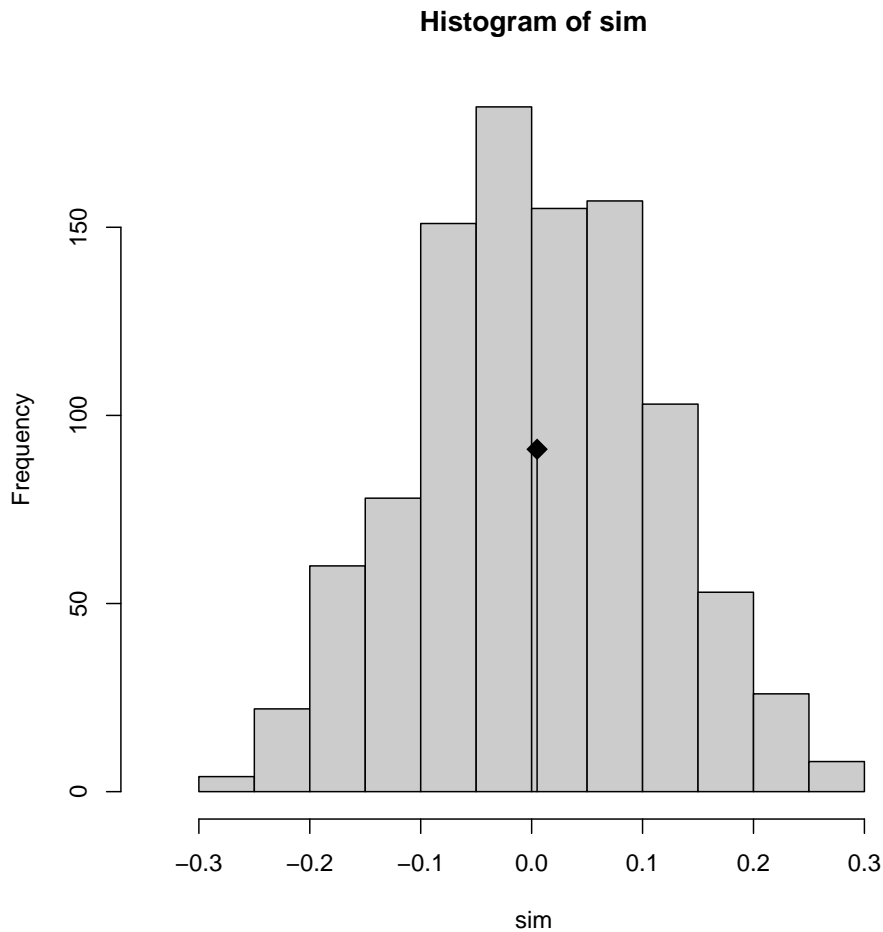
```
data(nancycats)
toto <- genind2genpop(nancycats)

##
## Converting data from a genind to a genpop object...
##
## ...done.

Dgen <- dist.genpop(toto,method=2)
Dgeo <- dist(nancycats$other$xy)
ibd <- mantel.randtest(Dgen,Dgeo)
ibd

## Monte-Carlo test
## Call: mantel.randtest(m1 = Dgen, m2 = Dgeo)
##
## Observation: 0.00492068
##
## Based on 999 replicates
## Simulated p-value: 0.488
## Alternative hypothesis: greater
##
##      Std.Obs Expectation      Variance
## 0.018023759 0.002993049 0.011438193
```

```
plot(ibd)
```

The original value of the correlation between the distance matrices is represented by the dot, while histograms represent permuted values (i.e., under the absence of spatial structure). Significant spatial structure would therefore result in the original value being out of the reference distribution. Here, isolation by distance is clearly not significant.

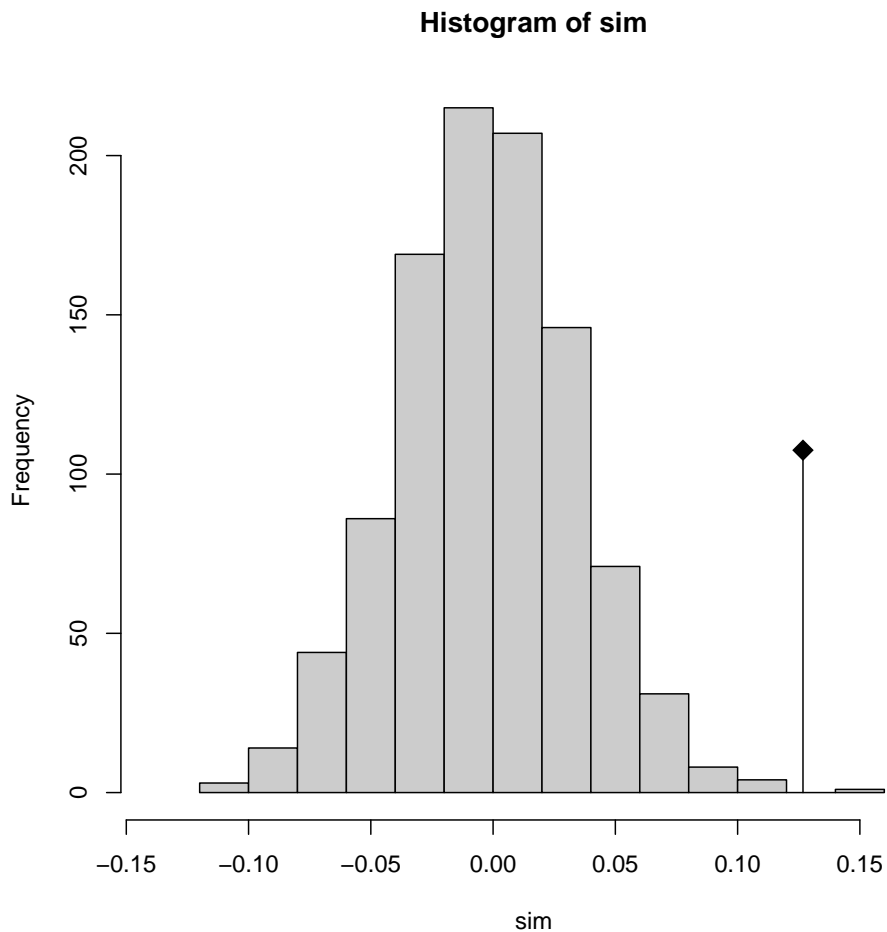
Let us provide another example using a dataset of individuals simulated under an IBD model:

```
data(spcaIllus)
x <- spcaIllus$dat2B
Dgen <- dist(x$tab)
Dgeo <- dist(other(x)$xy)
ibd <- mantel.randtest(Dgen,Dgeo)
ibd

## Monte-Carlo test
## Call: mantel.randtest(m1 = Dgen, m2 = Dgeo)
##
## Observation: 0.1267341
```

```
##
## Based on 999 replicates
## Simulated p-value: 0.002
## Alternative hypothesis: greater
##
##      Std.Obs  Expectation    Variance
## 3.533869077 -0.002448621  0.001336312
```

```
plot(ibd)
```

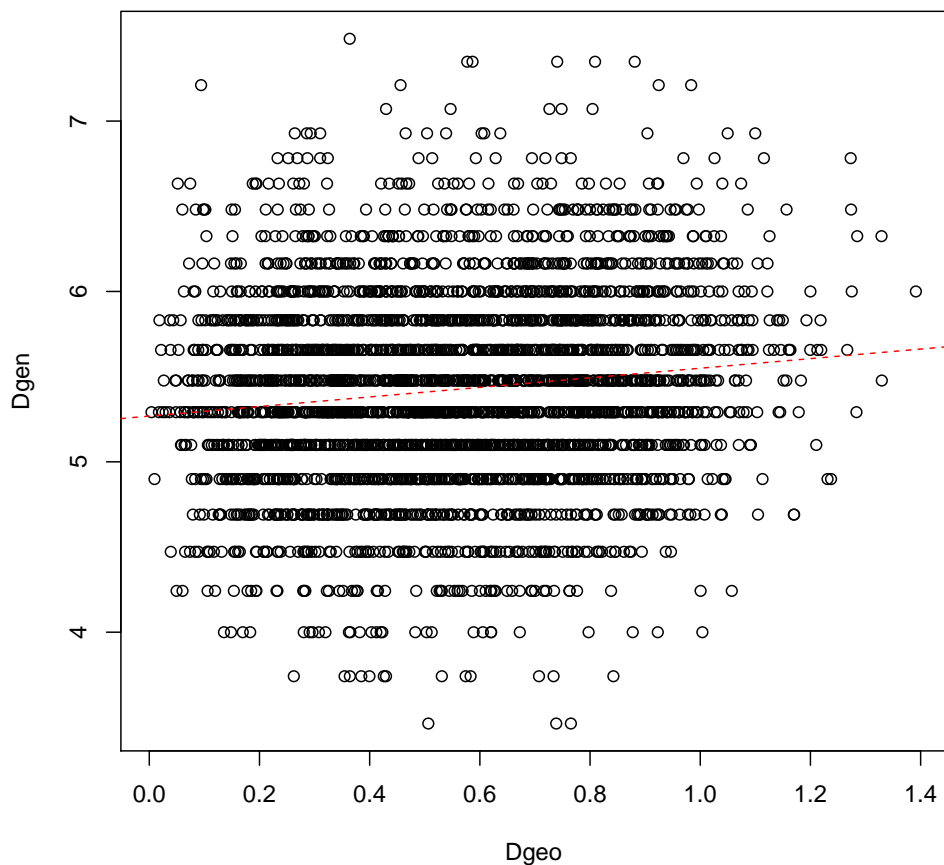


This time there is a clear isolation by distance pattern.

7.1.2 Cline or distant patches?

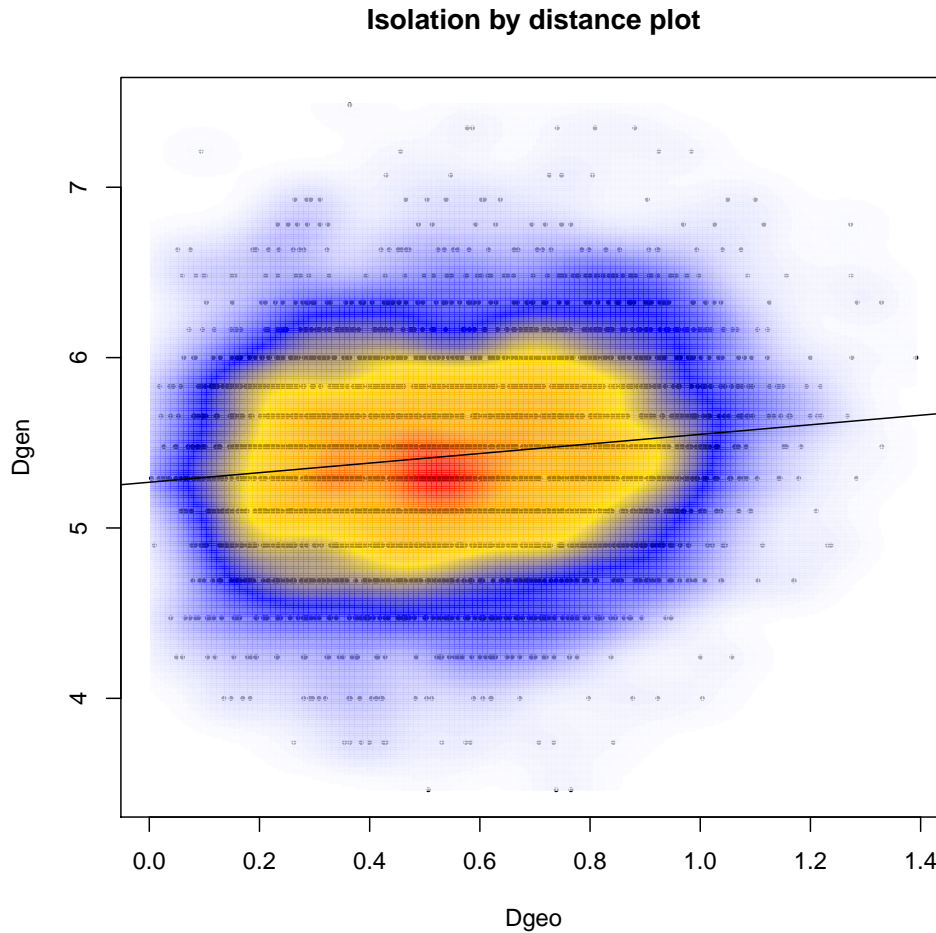
The correlation between genetic and geographic distances can occur under a range of different biological scenarios. Classical IBD would result in continuous clines of genetic differentiation and cause such correlation. However, distant and differentiated populations would also result in such a pattern. These are slightly different processes and we would like to be able to disentangle them. A very simple first approach is simply plotting both distances:

```
plot(Dgeo, Dgen)
abline(lm(Dgen~Dgeo), col="red", lty=2)
```



Most of the time, simple scatterplots fail to provide a good picture of the data as the density of points in the scatterplot is badly displayed. Colors can be used to provide better (and prettier) plots. Local density is measured using a 2-dimensional kernel density estimation (`kde2d`), and the results are displayed using `image`; `colorRampPalette` is used to generate a customized color palette:

```
library(MASS)
dens <- kde2d(Dgeo,Dgen, n=300)
myPal <- colorRampPalette(c("white","blue","gold", "orange", "red"))
plot(Dgeo, Dgen, pch=20,cex=.5)
image(dens, col=transp(myPal(300),.7), add=TRUE)
abline(lm(Dgen~Dgeo))
title("Isolation by distance plot")
```



The scatterplot clearly shows one single consistent cloud of point, without discontinuities which would have indicated patches. This is reassuring, since the data were actually simulated under an IBD (continuous) model.

7.2 Using Monmonier's algorithm to define genetic boundaries

Monmonier's algorithm [12] was originally designed to find boundaries of maximum differences between contiguous polygons of a tessellation. As such, the method was basically used in geographical analysis. More recently, [13] suggested that this algorithm could be employed to detect genetic boundaries among georeferenced genotypes (or populations). This algorithm is implemented using a more general approach than the initial one in *adegenet*.

Instead of using Voronoi tessellation as in the original version, the functions `monmonier` and `optimize.monmonier` can handle various neighbouring graphs such as Delaunay triangulation, Gabriel's graph, Relative Neighbours graph, etc. These graphs define spatial connectivity among locations (of genotypes or populations), with couple of locations being neighbours (if connected) or not. Another information is given by a set of markers which define genetic distances among these 'points'. The aim of Monmonier's algorithm is to find the path through the strongest genetic distances between neighbours. A more complete description of the principle of this algorithm will be found in the documentation of

monmonier. Indeed, the very purpose of this tutorial is simply to show how it can be used on genetic data.

Let's take the example from the function's manpage and detail it. The dataset used is `sim2pop`.

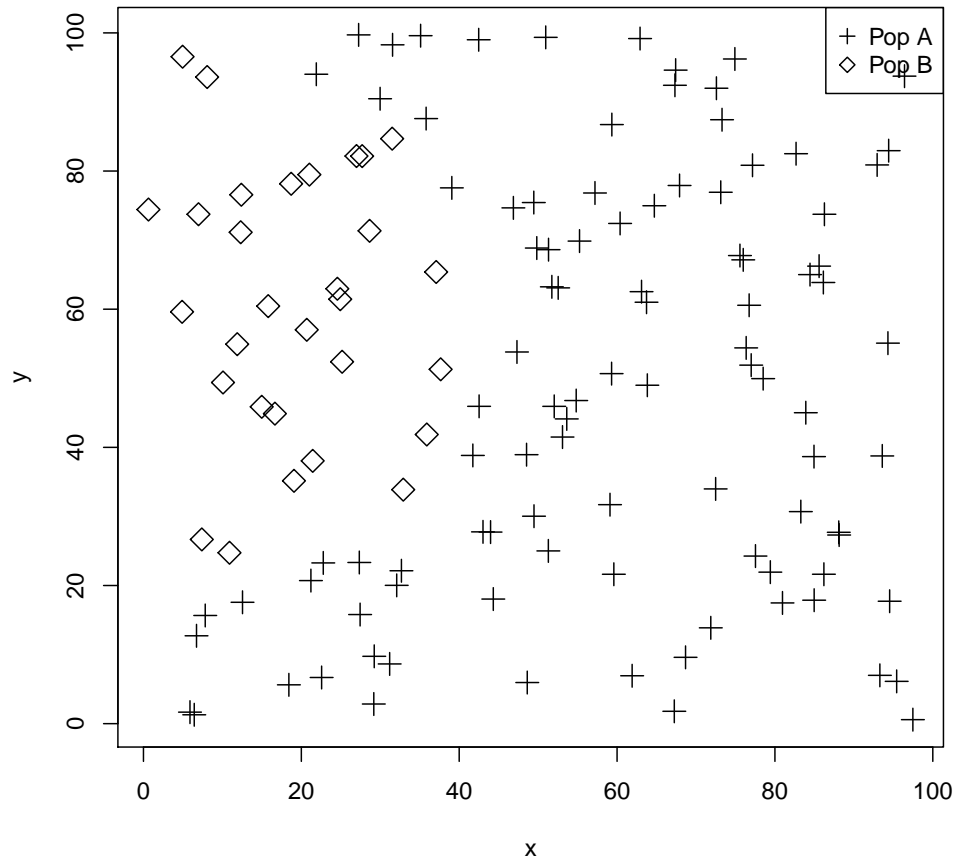
```
data(sim2pop)
sim2pop

## /// GENIND OBJECT ///////////
##
## // 130 individuals; 20 loci; 241 alleles; size: 184.5 Kb
##
## // Basic content
##   @tab: 130 x 241 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 7-17)
##   @loc.fac: locus factor for the 241 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: old2new(object = sim2pop)
##
## // Optional content
##   @pop: population of each individual (group size range: 30-100)
##   @strata: - empty -
##   @hierarchy: - empty -
##   @other: a list containing: xy

summary(sim2pop$pop)

## P01 P02
## 100 30

temp <- sim2pop$pop
levels(temp) <- c(3,5)
temp <- as.numeric(as.character(temp))
plot(sim2pop$other$xy,pch=temp,cex=1.5,xlab='x',ylab='y')
legend("topright",leg=c("Pop A", "Pop B"),pch=c(3,5))
```



There are two sampled populations in this dataset, with unequal sample sizes (100 and 30). Twenty microsatellite-like loci are available for all genotypes (no missing data). `monmonier` requires several arguments to be specified:

```
args(monmonier)

## function (xy, dist, cn, threshold = NULL, bd.length = NULL, nrun = 1,
##      skip.local.diff = rep(0, nrun), scanthres = is.null(threshold),
##      allowLoop = TRUE)
## NULL
```

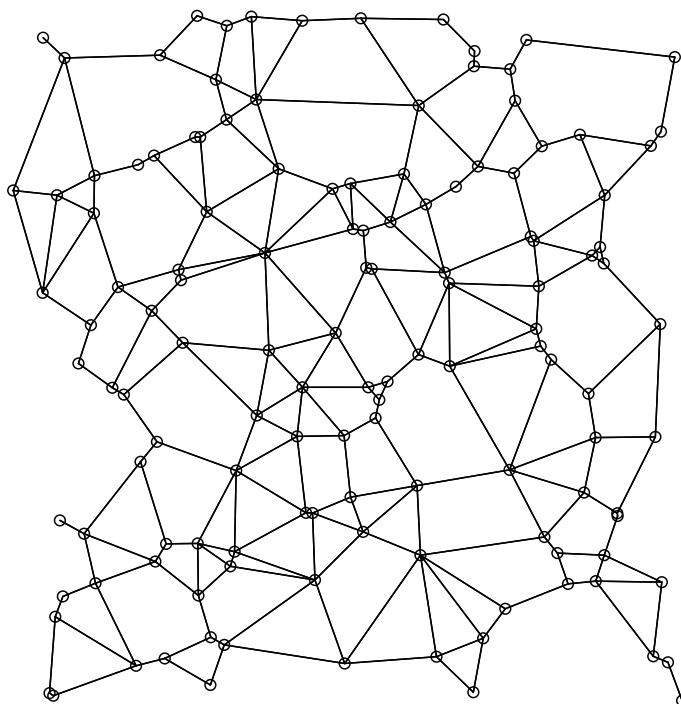
The first argument (`xy`) is a matrix of geographic coordinates, already stored in `sim2pop`. Next argument is an object of class `dist`, which is the matrix of pairwise genetic distances. For now, we will use the classical Euclidean distance between allelic profiles of the individuals. This is obtained by:

```
D <- dist(sim2pop$tab)
```

The next argument (`cn`) is a connection network. Routines for building such networks are scattered over several packages, but all made available through the function `chooseCN`. Here,

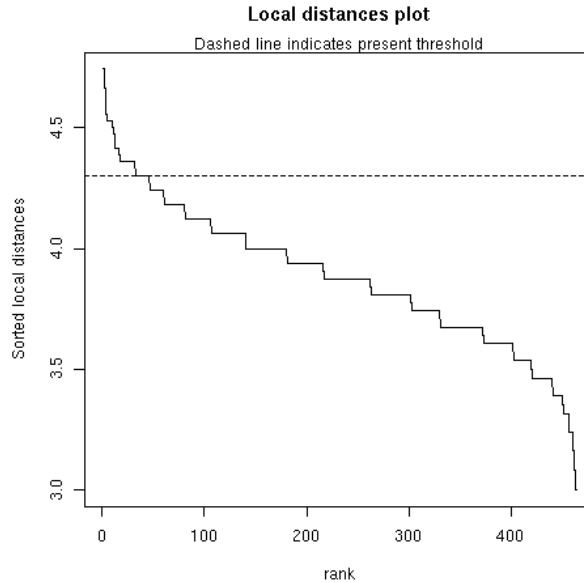
we disable the interactivity of the function (`ask=FALSE`) and select the second type of graph which is the graph of Gabriel (`type=2`).

```
gab <- chooseCN(sim2pop$other$xy,ask=FALSE,type=2)
```



The obtained network is automatically plotted by the function. It seems we are now ready to proceed to the algorithm.

```
mon1 <- monmonier(sim2pop$other$xy,D,gab)
```



This plot shows all local differences sorted in decreasing order. The idea behind this is that a significant boundary would cause local differences to decrease abruptly after the boundary. This should be used to choose the *threshold* difference for the algorithm to stop extending the boundary. Here, there is no indication of an actual boundary.

Why do the algorithm fail to find a boundary? Either because there is no genetic differentiation to be found, or because the signal differentiating both populations is too weak to overcome the random noise in genetic distances. What is the F_{st} between the two samples?

```
pairwise.fst(sim2pop)

## Error in eval(expr, envir, enclos): could not find function
"pairwise.fst"
```

This value would be considered as very weak differentiation ($F_{ST} = 0.023$). Is it significant? We can easily elaborate a permutation test of this F_{ST} value; to save computational time, we use only a small number of replicates to generate F_{ST} values in absence of population structure:

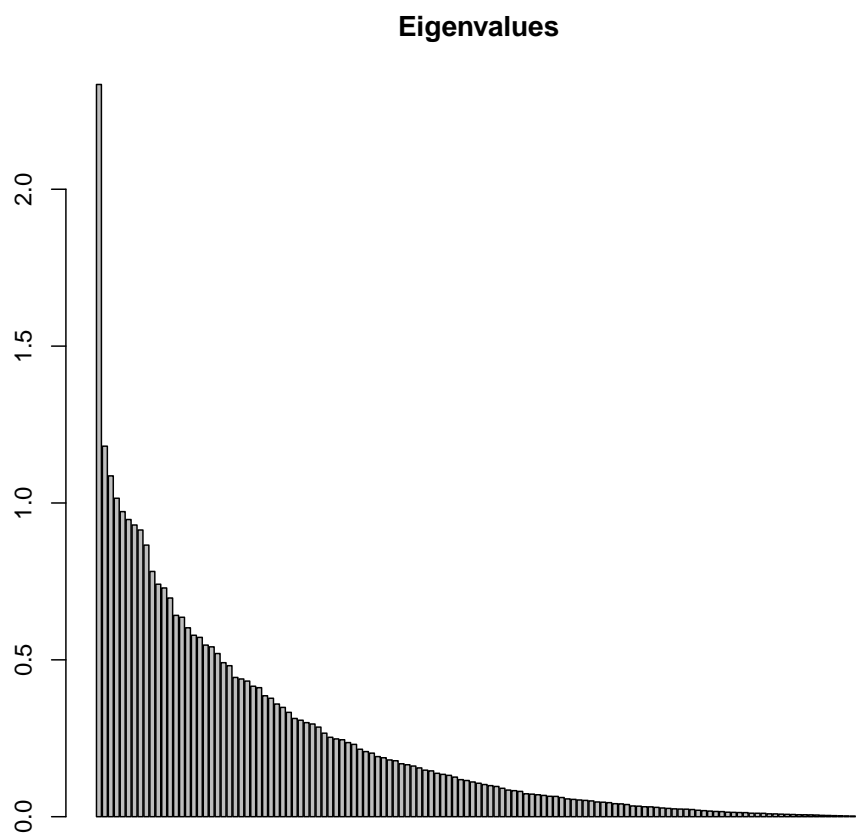
```
replicate(10, pairwise.fst(sim2pop, pop=sample(pop(sim2pop))))

## Error in FUN(X[[i]], ...): could not find function "pairwise.fst"
```

F_{ST} values in absence of population structure would be one order of magnitude lower (more replicate would give a very low p-value — just replace 10 by 200 in the above command). In fact, the two samples are indeed genetically differentiated.

Can Monmonier's algorithm find a boundary between the two populations? Yes, if we get rid of the random noise. This can be achieved using a simple ordination method such as Principal Coordinates Analysis.

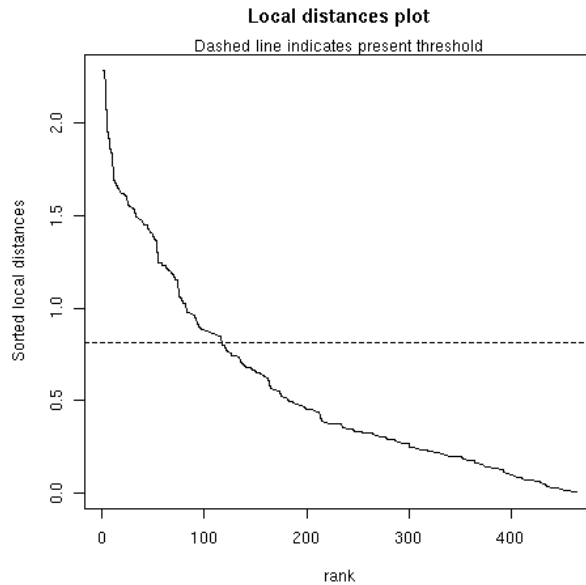
```
pco1 <- dudi.pco(D,scannf=FALSE,nf=1)
barplot(pco1$eig,main="Eigenvalues")
```



We retain only the first eigenvalue. The corresponding coordinates are used to redefine the genetic distances among genotypes. The algorithm is then re-run.

```
D <- dist(pco1$li)
```

```
mon1 <- monmonier(sim2pop$other$xy,D,gab)
```



```
##
## #####
## # List of paths of maximum differences between neighbours #
## #           Using a Monmonier based algorithm           #
## #####
##
## $call:monmonier(xy = sim2pop$other$xy, dist = D, cn = gab, scanthres = FALSE)
##
##      # Object content #
## Class: monmonier
## $nrun (number of successive runs): 1
## $run1: run of the algorithm
## $threshold (minimum difference between neighbours): 1.631
## $xy: spatial coordinates
## $cn: connection network
##
##      # Runs content #
## # Run 1
## # First direction
## Class: list
## $path:
##           x           y
## Point_1 14.98299 93.81162
##
## $values:
## 4.563555
## # Second direction
```

```
## Class: list
## $path:
##           x           y
## Point_1 14.98299 93.81162
## Point_2 30.74508 87.57724
## Point_3 33.66093 86.14115
## ...
##
## $values:
## 4.563555 3.23581 3.906439 ...
```

This may take some time... but never more than five minutes on an 'ordinary' personal computer. The object `mon1` contains the whole information about the boundaries found. As several boundaries can be sought at the same time (argument `nrun`), you have to specify about which run and which direction you want to get informations (values of differences or path coordinates). For instance:

```
names(mon1)

## [1] "run1"      "nrun"      "threshold" "xy"        "cn"        "call"

names(mon1$run1)

## [1] "dir1" "dir2"

mon1$run1$dir1

## $path
##           x           y
## Point_1 14.98299 93.81162
##
## $values
## [1] 4.563555
```

It can also be useful to identify which points are crossed by the barrier; this can be done using `coords.monmonier`:

```
coords.monmonier(mon1)

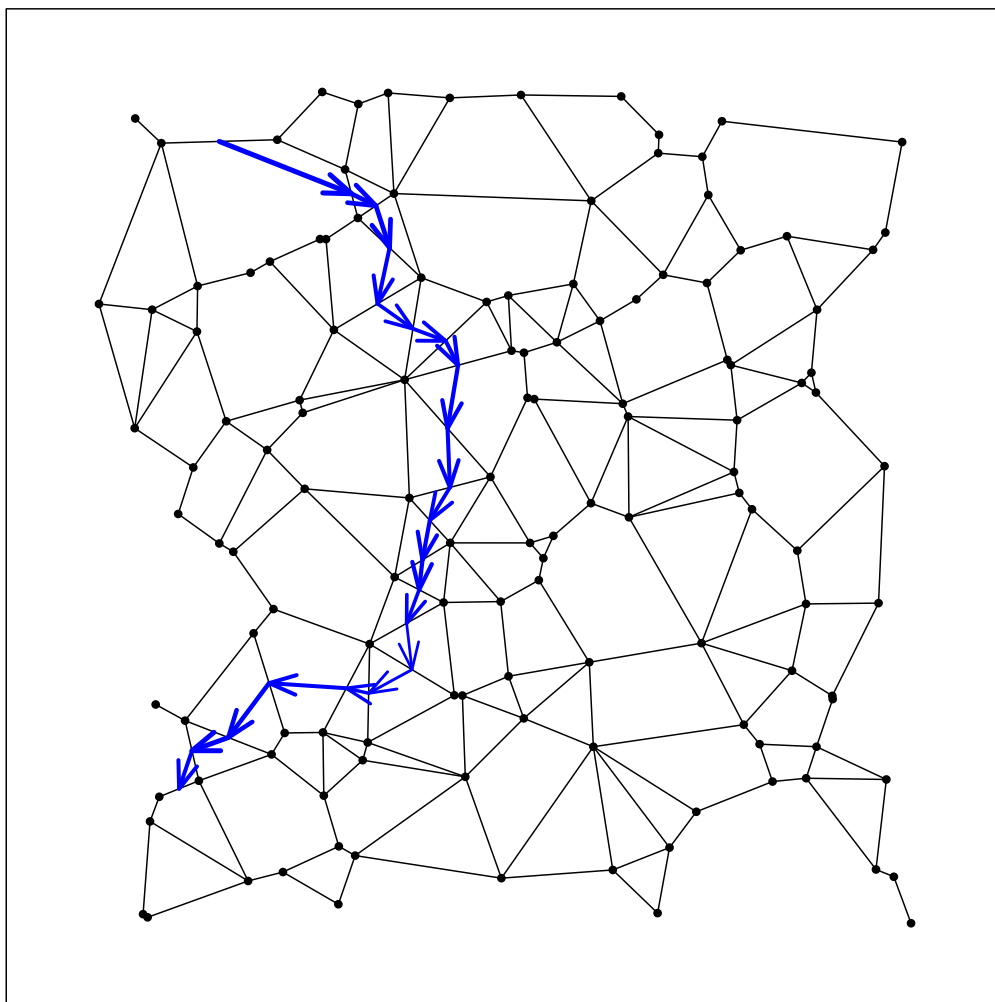
## $run1
## $run1$dir1
##           x.hw      y.hw first second
## Point_1 14.98299 93.81162    11    125
##
## $run1$dir2
```

##		x.hw	y.hw	first	second
##	Point_1	14.98299	93.81162	11	125
##	Point_2	30.74508	87.57724	44	128
##	Point_3	33.66093	86.14115	20	128
##	Point_4	35.28914	81.12578	68	128
##	Point_5	33.85756	74.45492	68	117
##	Point_6	38.07622	71.47532	68	122
##	Point_7	41.97494	70.02783	35	122
##	Point_8	43.45812	67.12026	69	122
##	Point_9	42.20206	59.59613	22	122
##	Point_10	42.48613	52.55145	22	124
##	Point_11	40.08702	48.61795	13	124
##	Point_12	39.20791	43.89978	13	127
##	Point_13	38.81236	40.34516	62	127
##	Point_14	37.32112	36.35265	62	130
##	Point_15	37.96426	30.82105	94	130
##	Point_16	32.79703	28.00517	16	130
##	Point_17	30.12832	28.60376	85	130
##	Point_18	20.92496	29.21211	63	119
##	Point_19	16.05811	22.72600	61	126
##	Point_20	11.72524	21.15519	89	126
##	Point_21	10.18696	16.61536	74	89

The returned dataframe contains, in this order, the x and y coordinates of the points of the barrier, and the identifiers of the two 'parent' points, that is, the points whose barycenter is the point of the barrier.

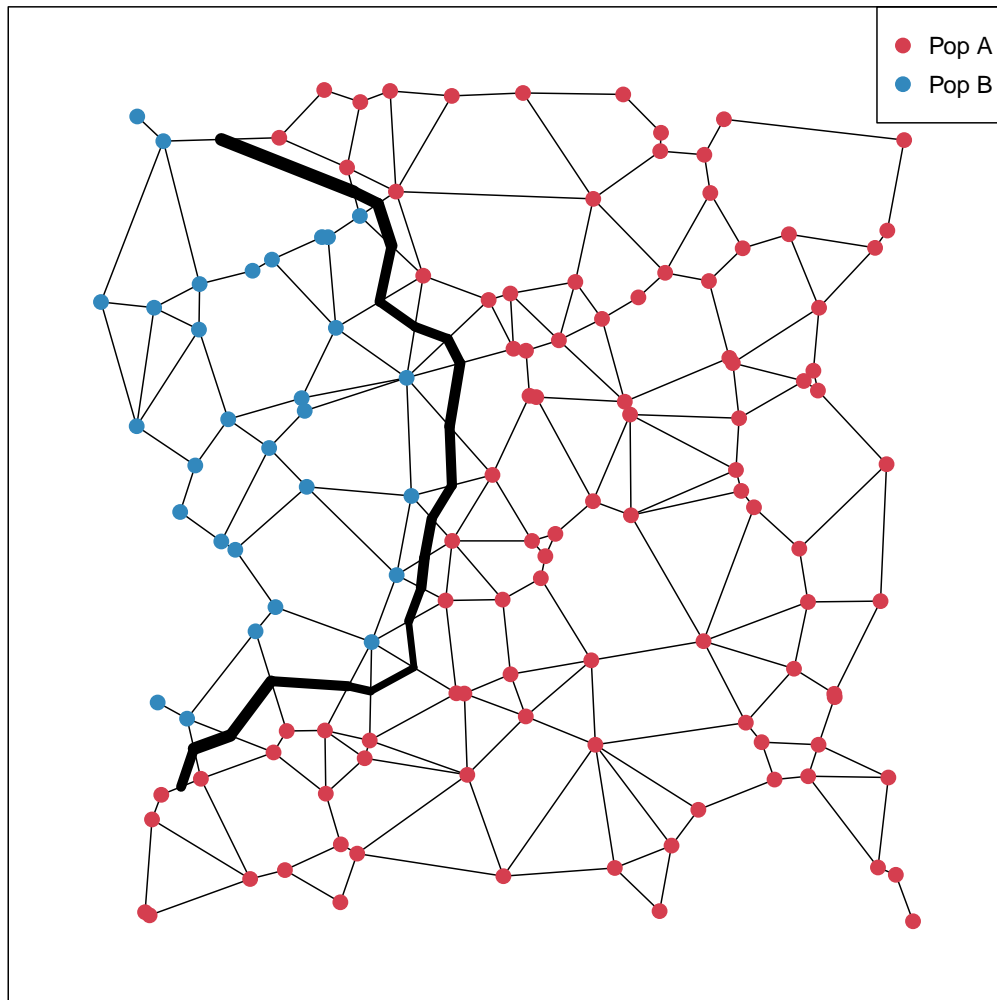
Finally, you can plot very simply the obtained boundary using the method `plot`:

```
plot(mon1)
```



see arguments in `?plot.monmonier` to customize this representation. Last, we can compare the inferred boundary with the actual distribution of populations:

```
plot(mon1, add.arrows=FALSE, bwd=10, col="black")
points(sim2pop$other$xy, cex=2, pch=20,
       col=fac2col(pop(sim2pop), col.pal=spectral))
legend("topright", leg=c("Pop A", "Pop B"), pch=c(20),
      col=spectral(2), pt.cex=2)
```



Not too bad...

8 Simulating hybridization

The function `hybridize` allows to simulate hybridization between individuals from two distinct genetic pools, or more broadly between two `genind` objects. Here, we use the example from the manpage of the function, to go a little further. Please have a look at the documentation, especially at the different possible outputs (outputs for the software STRUCTURE is notably available).

```
temp <- seppop(microbov)
names(temp)
```

```
## [1] "Borgou"      "Zebu"        "Lagunaire"
## [4] "NDama"      "Somba"      "Aubrac"
## [7] "Bazadais"   "BlondeAquitaine" "BretPieNoire"
## [10] "Charolais"  "Gascon"     "Limousin"
## [13] "MaineAnjou" "Montbeliard" "Salers"
```

```
salers <- temp$Salers
zebu <- temp$Zebu
zebler <- hybridize(salers, zebu, n=40, pop="zebler")
```

A first generation (F1) of hybrids 'zebler' is obtained. Is it possible to perform a backcross, say, with 'salers' population? Yes, here it is:

```
F2 <- hybridize(salers, zebler, n=40)
F3 <- hybridize(salers, F2, n=40)
F4 <- hybridize(salers, F3, n=40)
```

Finally, note that despite this example shows hybridization between diploid organisms, `hybridize` is not restrained to this case. In fact, organisms with any even level of ploidy can be used, in which case half of the genes is taken from each reference population. Ultimately, more complex mating schemes could be implemented... suggestion or (better) contributions are welcome!

References

- [1] Jombart, T. (2008) adegenet: a R package for the multivariate analysis of genetic markers. *Bioinformatics* 24: 1403-1405.
- [2] R Development Core Team (2014). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
- [3] Dray S and Dufour A-B (2007) The ade4 package: implementing the duality diagram for ecologists. *Journal of Statistical Software* 22: 1-20.
- [4] Jombart T, Devillard S and Balloux, F (2010). Discriminant analysis of principal components: a new method for the analysis of genetically structured populations. *BMC Genetics* 11: 94.
- [5] Jombart T, Devillard S, Dufour A-B and Pontier D (2008) Revealing cryptic spatial patterns in genetic variability by a new multivariate method. *Heredity* 101: 92-103.
- [6] Jombart T, Eggo RM, Dodd PJ and Balloux F (2010) Reconstructing disease outbreaks from genetic data: a graph approach. *Heredity* 106: 383-390.
- [7] Jombart T, Pontier D and Dufour A-B (2009) Genetic markers in the playground of multivariate analysis. *Heredity* 102: 330-341.
- [8] Paradis E, Claude J, and Strimmer K (2004) APE: analyses of phylogenetics and evolution in R language. *Bioinformatics*: 20, 289-290.
- [9] Charif D, and Lobry J (2007) SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis. *in* Structural approaches to sequence evolution: Molecules, networks, populations, *Springer Verlag*, 207-232.
- [10] Goudet J., Raymond, M., DeMeeus, T. and Rousset F. (1996) Testing differentiation in diploid populations. *Genetics* 144: 1933-1940.
- [11] Nei M (1973) Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci USA* 70: 3321-3323.
- [12] Monmonier M (1973) Maximum-difference barriers: an alternative numerical regionalization method. *Geographical Analysis* 3: 245-261.
- [13] Manni F, Guerard E and Heyer E (2004) Geographic patterns of (genetic, morphologic, linguistic) variation: how barriers can be detected by "Monmonier's algorithm". *Human Biology* 76: 173-190.