

Poppr 1.0.0: An R package for genetic analysis of populations with mixed (clonal/sexual) reproduction

Zhian N. Kamvar¹ and Niklaus J. Grünwald^{1,2}

1) Department of Botany and Plant Pathology, Oregon State University, Corvallis, OR

2) Horticultural Crops Research Laboratory, USDA-ARS, Corvallis, OR

June 6, 2013

Abstract

Poppr provides open-source, cross-platform tools for quick analysis of population genetic data enabling focus on data analysis and interpretation. While there are a plethora of packages for population genetic analysis, few are able to offer quick and easy analysis of populations with mixed reproductive modes. *Poppr*'s main advantage is the ease of use and integration with other packages such as *adegenet* and *vegan*, including support for novel methods such as clone correction, multilocus genotype analysis, calculation of Bruvo's distance and the index of association.



Contents

1	Introduction	3
1.1	Purpose	3
1.2	Installation	3
1.3	Quick start	4
1.4	Get out of my dreams and into my R {importing data into poppr}	6
1.4.1	Function: getfile	6
1.4.2	Function: read.genalex	9
1.4.3	Genalex formatting shortcuts	11
1.4.4	Other ways of importing data	12
1.4.5	Function: genind2genalex	12
1.5	Getting to know adegenet's genind object	14
1.5.1	The other slot	14
1.5.2	Setting the population factor {adegenet's function: pop}	16
2	Data Manipulation	17
2.1	Inside the golden days of missing data {replace or remove missing data}	18
2.1.1	Function: missingno	18
2.2	Can you take me hier(archy)? {population hierarchy construction}	21
2.2.1	Function: splitcombine	21
2.3	Divide (populations) and conquer (your analysis) {extract populations}	24
2.3.1	Function: popsub	25
2.4	Attack of the clone correction {clone-censor data sets}	26
2.4.1	Function: clonecorrect	27
2.5	Every day I'm shuffling (data sets) {permutations and bootstrap resampling}	27
2.5.1	Function: shufflepop	28
3	Multilocus Genotype Analysis	29
3.1	Just a peek {How many multilocus genotypes are in our data set?}	29
3.1.1	Function: mlg	30
3.2	Clone-ing around {MLGs across populations}	30
3.2.1	Function: mlg.crosspop	30
3.3	Bringing something to the table {producing MLG tables and graphs}	31
3.3.1	Function: mlg.table	32
3.4	Getting into the mix {combining MLG functions}	34
3.4.1	Function: mlg.vector	35
3.5	Do you see what I see? {alternative data visualization}	37
4	Index and Distance Calculations	37
4.1	The missing linkage disequilibrium {calculating the index of association, I_A and \bar{r}_d }	37
4.1.1	Function: ia	38
4.2	Going the distance {dissimilarity distance}	40
4.2.1	Function: diss.dist	40
4.3	Step by stepwise mutation {Bruvo's distance}	41
4.3.1	Function: bruvo.dist	41
4.4	See the forest for the trees {visualizing distances with dendrograms and networks}	43
4.4.1	Function: bruvo.boot	43
4.4.2	Function: greycurve	45
4.4.3	Function: bruvo.msn	47
4.4.4	Function: poppr.msn	49

5	I know what you did last summary table {diversity table}	51
5.1	Function: poppr	51
6	Appendix	54
6.1	Algorithmic Details	54
6.1.1	I_A and \bar{r}_d	54
6.1.2	Bruvo's distance	56
6.1.3	Special Cases of Bruvo's distance	57
6.2	Exporting Graphics	57
6.2.1	Basics	58
6.2.2	Image Editors	58
6.2.3	Exporting ggplot2 graphics	58
6.2.4	Exporting any graphics	59
6.3	Function calls	59

1 Introduction

1.1 Purpose

Poppr is an R package with convenient functions for analysis of genetic data with mixed modes of reproduction including sexual and clonal reproduction. While there are many R packages in CRAN and other repositories with tools for population genetic analyses, few are appropriate for populations with mixed modes of reproduction. There are several stand alone programs that can handle these types of data sets, but they are often platform specific and often only accept specific data types. Furthermore, a typical analysis often involves switching between many programs, and converting data to each specific format.

Poppr is designed to make analysis of populations with mixed reproductive modes more streamlined and user friendly so that the researcher using it can focus on data analysis and interpretation. *Poppr* allows analysis of haploid and diploid dominant/co-dominant marker data including microsatellites, Single Nucleotide Polymorphisms (SNP), and Amplified Fragment Length Polymorphisms (AFLP). To avoid creating yet another file format that is specific to a program, *poppr* was created on the backbone of the popular R package *adegenet* and can take all the file formats that *adegenet* can take (Genpop, Gentix, Fstat, and Structure) and newly introduces compatibility with GenAlEx formatted files (exported to CSV). This means that anything you can analyze in *adegenet* can be further analyzed with *poppr*.

The real power of *poppr* is in the data manipulation and analytic tools. *Poppr* has the ability to define multiple population hierarchies, clone-censor, and subset data sets. With *poppr* you can also quickly calculate Bruvo's distance, the index of association, and easily determine which multilocus genotypes are shared across populations.

1.2 Installation

This manual assumes that you have already installed R. If you have not, please refer to The CRAN home page at <http://cran.r-project.org/>. The author also recommends utilizing an R gui such as Rstudio (<http://www.rstudio.com/>) for a better R experience. To install *poppr* from CRAN is as simple as selecting "Package Installer" from the menu "Packages & Data" in the gui or by typing in your command line:

```
> install.packages("poppr", dependencies=TRUE)
```

If everything is working perfectly, all the dependencies (*adegenet*, *pegas*, *vegan*, *ggplot2*, *phangorn*, *ape* and *igraph*) should be installed. In the unfortunate case this does not work, consult <http://cran.r-project.org/doc/manuals/R-admin.html#Installing-packages>. If you choose to install *poppr* from a source file, you should first make sure to install all of the dependencies with the following command:

```
> install.packages(c("adegenet", "pegas", "vegan", "ggplot2", "phangorn", "ape", "igraph"))
```

After that, you can install the package with:

```
> install.packages("/path/to/poppr.tar.gz", type="source", repos=NULL)
```

1.3 Quick start

The author assumes that if you have reached this point in the manual, then you have successfully installed R and *poppr*. Before proceeding, you should be aware that R is case sensitive. This means that the words “Case” and “case” are different from R’s perspective. You should also know where your R package Library is located.

WHAT OR WHERE IS MY R PACKAGE LIBRARY?

R is as powerful as it is through a community of people who submit extra code called “Packages” to help it do specific things. These packages live in a certain place on your computer called an R library. You can find out where this library is by typing `.libPaths()`

Importing a file into R involves you knowing the path to your file and then typing that into R’s console. `getfile()` will help provide a point and click interface for selecting a file. There are two steps: Before you do anything, you’ll want to tell your computer to search R’s library to find the *poppr* and load the package:

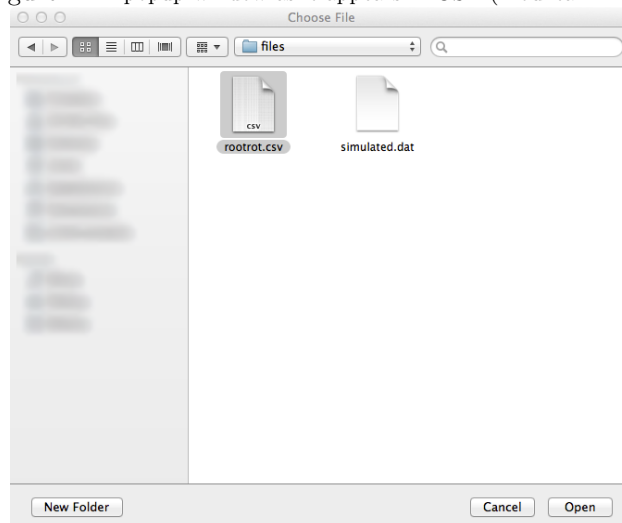
```
> library(poppr)
```

After that, you can use `getfile()`

```
> x <- getfile()
```

At this point, a pop up window will appear like this¹:

Figure 1: A popup window as it appears in OSX (Mountain Lion).



HEY! MY WINDOW DOESN'T LOOK LIKE THAT!

Now, this window will not match up to your window on your computer because you will probably not be in the right directory. Remember the first path in `.libPaths()`? Move to a folder called **poppr** in that path. In that folder, you will find another folder called **files**. Move there and your window will match the one displayed.

¹This window sometimes appears behind your current session of R, depending on the GUI and you will have to toggle to this window

We can navigate throughout your entire computer through this little window and tell R where to go. The example I'm using goes to your R library directory. If you don't know where that is, you can find it by typing `.libPaths()` into the R command line. Once we select a file, the file name and its path will be stored in the variable, `x`. We can confirm what we selected by simply typing `x` into R's command line.

```
> x
```

```
$files
[1] "/path/to/R/poppr/files/rootrot.csv"

$path
[1] "/path/to/R/poppr/files"
```

Here we can see that `x` is a list with two entries: `$files` giving you the files you selected and `$path` giving you the path to those files.

NOT SURE WHAT I MEAN BY PATH OR WORKING DIRECTORY?

For anyone who has never used a command line, this is a new concept. You can think of the path as an address. So instead of `"/path/to/R"`, you could have `"/USA/Oregon/Corvallis"`. Or on your computer, it could be `"C:/users/poppr-user/R/win-library/2.15"` on Windows (where "poppr-user" is your username) or `"/Library/Frameworks/R.framework/Versions/2.15/Resources/library"` on OSX. Each slash represents a folder that you would click through when you are using the mouse.

A working directory is simply the folder that R is working in. It is where you can access and write files. When you tell R to read a file, it will only look for that file in your working directory. Note that you will not endanger your files by reading them into R. R works by making a copy of the file into memory. This means that you can manipulate the data in any way that you want without ever losing the content.

To find out your current working directory, type `getwd()` into the R console. Usually, you will start off a session in your "home" directory, which will look like this: `"~/."`. The command `setwd()` will change your working directory to any place of your choice on your computer as indicated by the path that you provide. For more information, see Quick R at <http://www.statmethods.net>.

We will use `x$files` to access the file. The `poppr()` function provides a quick and convenient first analysis of your data directly from the file on the your disk (For information on importing your data into R, see section 1.4, *Get out of my dreams and into my R*).

```
> popdata <- poppr(x$files)
```

```
| Athena_1
| Athena_2
| Athena_3
| Athena_4
| Athena_5
| Athena_6
| Athena_7
| Athena_8
| Athena_9
| Athena_10
| Mt. Vernon_1
| Mt. Vernon_2
| Mt. Vernon_3
| Mt. Vernon_4
| Mt. Vernon_5
| Mt. Vernon_6
| Mt. Vernon_7
| Mt. Vernon_8
| Total
```

The output of `poppr()` was assigned to the variable `popdata`, so let's look at the data.

```
> popdata
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	rbarD	File
1	Athena_1	9	7	7.000	0.000	1.889	6.231	0.944	0.932	2.925	0.210	rootrot.csv
2	Athena_2	12	12	10.000	NaN	2.485	12.000	1.000	1.000	4.160	0.128	rootrot.csv
3	Athena_3	10	2	2.000	0.000	0.325	1.220	0.200	0.571	2.000	1.000	rootrot.csv
4	Athena_4	13	9	7.154	0.769	1.946	5.121	0.872	0.687	5.495	0.372	rootrot.csv
5	Athena_5	10	7	7.000	0.000	1.834	5.556	0.911	0.866	4.532	0.353	rootrot.csv
6	Athena_6	5	5	5.000	0.000	1.609	5.000	1.000	1.000	2.464	0.190	rootrot.csv
7	Athena_7	11	10	9.182	0.386	2.272	9.308	0.982	0.955	2.129	0.086	rootrot.csv
8	Athena_8	8	6	6.000	0.000	1.667	4.571	0.893	0.831	3.857	0.323	rootrot.csv
9	Athena_9	10	10	10.000	0.000	2.303	10.000	1.000	1.000	2.815	0.118	rootrot.csv
10	Athena_10	9	8	8.000	0.000	2.043	7.364	0.972	0.948	2.849	0.137	rootrot.csv
11	Mt. Vernon_1	10	9	9.000	0.000	2.164	8.333	0.978	0.952	7.132	0.276	rootrot.csv
12	Mt. Vernon_2	6	6	6.000	0.000	1.792	6.000	1.000	1.000	20.649	0.492	rootrot.csv
13	Mt. Vernon_3	8	6	6.000	0.000	1.667	4.571	0.893	0.831	2.117	0.106	rootrot.csv
14	Mt. Vernon_4	12	8	6.833	0.665	1.814	4.500	0.848	0.681	3.008	0.255	rootrot.csv
15	Mt. Vernon_5	17	7	5.541	0.828	1.758	5.070	0.853	0.848	2.677	0.340	rootrot.csv
16	Mt. Vernon_6	12	11	9.318	0.466	2.369	10.286	0.985	0.958	19.498	0.467	rootrot.csv
17	Mt. Vernon_7	12	9	7.818	0.649	2.095	7.200	0.939	0.870	1.208	0.153	rootrot.csv
18	Mt. Vernon_8	13	9	7.346	0.764	2.032	6.259	0.910	0.794	1.153	0.169	rootrot.csv
19	Total	187	119	9.612	0.612	4.558	68.972	0.991	0.720	14.371	0.271	rootrot.csv

The fields you see in the output include:

- Pop – Population name (Note that “Total” also means “Pooled”).
- N – Number of individuals observed.
- MLG – Number of multilocus genotypes (MLG) observed.
- eMLG – The number of expected MLG at the smallest sample size ≥ 10 based on rarefaction. [8]
- SE – Standard error based on eMLG [7]
- H – Shannon-Wiener Index of MLG diversity. [16]
- G – Stoddart and Taylor’s Index of MLG diversity. [18]
- Hexp – Nei’s 1978 genotypic diversity (corrected for sample size), or Expected Heterozygosity. [11]
- E.5 – Evenness, E_5 . [15][10][4]
- Ia – The index of association, I_A . [2] [17] [1]
- rbarD – The standardized index of association, \bar{r}_d . [1]

These fields are further described in section 5, *I know what you did last summary table* at the end of this vignette.

1.4 Get out of my dreams and into my R {importing data into poppr}

There are several ways of reading data into R.

1.4.1 Function: getfile

`getfile` gives the user an easy way to point R to the directory in which your data is stored. It is only meant for R GUIs such as Rstudio. Using this on the command line has very little advantage over setting the working directory manually.

Default Command:

```
getfile(multi = FALSE, pattern = NULL, combine = TRUE)
```

- multi – This is normally set to FALSE, meaning that it will only grab the file you selected. If it’s TRUE, it will grab all files within the directory, constrained only by what you type into the pattern field.

- **pattern** - A pattern that you want to filter the files you get. This accepts regular expressions, so you must be careful with anything that is not an alphanumeric character.
- **combine** - This tells **getfile** to combine the path and all the files. This is set to **TRUE** by default so that you can access your files no matter what working directory you are in.

This method works for a single file, but let's say you had a lot of data sets you wanted to import. You would have to do all of these one by one, right? Not so. **getfile** has a nice little flag called **multi** telling the computer that you want to grab multiple files in the folder. You would use this with **poppr.all** to produce a summary table for all of your files²:

```
> x <- getfile(multi=TRUE)
```

A window would pop up again, and you should navigate to the same directory as you had before, and select any of the files in that directory.

```
> x
```

```
$files
[1] "/path/to/R/poppr/files/rootrot2.csv"  "/path/to/R/poppr/files/rootrot.csv"
[3] "/path/to/R/poppr/files/simulated.dat"

$path
[1] "/path/to/R/poppr/files"
```

As you can see, now all of the files that existed in that directory are there! Now you can look at all those files at once!

```
> poppr.all(x$files)
```

```
| File: rootrot.csv
| Athena_1
| Athena_2
| Athena_3
| Athena_4
| Athena_5
| Athena_6
| Athena_7
| Athena_8
| Athena_9
| Athena_10
| Mt. Vernon_1
| Mt. Vernon_2
| Mt. Vernon_3
| Mt. Vernon_4
| Mt. Vernon_5
| Mt. Vernon_6
| Mt. Vernon_7
| Mt. Vernon_8
| Total
| File: rootrot2.csv
| 1
| 2
| 3
| 4
| 5
| 6
| 7
| 8
| 9
| 10
| Total
| File: simulated.dat
| Total
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	rbarD	File
1	Athena_1	9	7	7.000	0.000	1.889	6.231	0.944	0.932	2.925	0.210	rootrot.csv
2	Athena_2	12	12	10.000	NaN	2.485	12.000	1.000	1.000	4.160	0.128	rootrot.csv
3	Athena_3	10	2	2.000	0.000	0.325	1.220	0.200	0.571	2.000	1.000	rootrot.csv
4	Athena_4	13	9	7.154	0.769	1.946	5.121	0.872	0.687	5.495	0.372	rootrot.csv
5	Athena_5	10	7	7.000	0.000	1.834	5.556	0.911	0.866	4.532	0.353	rootrot.csv
6	Athena_6	5	5	5.000	0.000	1.609	5.000	1.000	1.000	2.464	0.190	rootrot.csv
7	Athena_7	11	10	9.182	0.386	2.272	9.308	0.982	0.955	2.129	0.086	rootrot.csv

²These files do not need to be similar in any way to do this analysis

```

8      Athena_8      8      6      6.000 0.000 1.667 4.571 0.893 0.831 3.857 0.323 rootrot.csv
9      Athena_9     10     10     10.000 0.000 2.303 10.000 1.000 1.000 2.815 0.118 rootrot.csv
10     Athena_10      9      8      8.000 0.000 2.043 7.364 0.972 0.948 2.849 0.137 rootrot.csv
11 Mt. Vernon_1     10      9      9.000 0.000 2.164 8.333 0.978 0.952 7.132 0.276 rootrot.csv
12 Mt. Vernon_2      6      6      6.000 0.000 1.792 6.000 1.000 1.000 20.649 0.492 rootrot.csv
13 Mt. Vernon_3      8      6      6.000 0.000 1.667 4.571 0.893 0.831 2.117 0.106 rootrot.csv
14 Mt. Vernon_4     12      8      6.833 0.665 1.814 4.500 0.848 0.681 3.008 0.255 rootrot.csv
15 Mt. Vernon_5     17      7      5.541 0.828 1.758 5.070 0.853 0.848 2.677 0.340 rootrot.csv
16 Mt. Vernon_6     12     11      9.318 0.466 2.369 10.286 0.985 0.958 19.498 0.467 rootrot.csv
17 Mt. Vernon_7     12      9      7.818 0.649 2.095 7.200 0.939 0.870 1.208 0.153 rootrot.csv
18 Mt. Vernon_8     13      9      7.346 0.764 2.032 6.259 0.910 0.794 1.153 0.169 rootrot.csv
19      Total     187    119      9.612 0.612 4.558 68.972 0.991 0.720 14.371 0.271 rootrot.csv
20          1     19     16      9.211 0.701 2.726 14.440 0.982 0.942 14.229 0.313 rootrot2.csv
21          2     18     18     10.000 0.000 2.890 18.000 1.000 1.000 9.143 0.194 rootrot2.csv
22          3     18      8      5.265 1.009 1.609 3.375 0.745 0.594 22.843 0.573 rootrot2.csv
23          4     25     17      7.887 1.124 2.575 9.615 0.933 0.710 18.488 0.415 rootrot2.csv
24          5     27     14      7.511 1.057 2.446 9.720 0.932 0.828 23.002 0.520 rootrot2.csv
25          6     17     15      9.338 0.633 2.670 13.762 0.985 0.949 17.778 0.410 rootrot2.csv
26          7     23     19      9.178 0.764 2.872 16.030 0.980 0.902 19.162 0.405 rootrot2.csv
27          8     21     15      8.273 0.981 2.558 10.756 0.952 0.820 24.313 0.543 rootrot2.csv
28          9     10     10     10.000 0.000 2.303 10.000 1.000 1.000 2.815 0.118 rootrot2.csv
29         10      9      8      8.000 0.000 2.043 7.364 0.972 0.948 2.849 0.137 rootrot2.csv
30      Total     187    119      9.612 0.612 4.558 68.972 0.991 0.720 14.371 0.271 rootrot2.csv
31      Total     100      6      6.000 0.000 1.235 2.790 0.648 0.735 0.050 0.061 simulated.dat

```

You've seen examples of how to use `getfile` to extract a single file and all the files in a directory, but what if you wanted many files, but only wanted ones that were of a certain type or had a certain name? This is what you would use the `pattern` argument for. A perfect use would be the example data contained in the *adegenet* package. Let's take a look at the names of these files.

For the rest of this section, remember that every time you invoke `getfile()`, a window will pop up and you should select a file before hitting enter.

```
> getfile(multi=TRUE)
```

Navigate to the *adegenet* folder in your R library.

```

$files
[1] "/path/to/R/adegenet/files/AFLP.txt"
[2] "/path/to/R/adegenet/files/exampleSnpDat.snp"
[3] "/path/to/R/adegenet/files/monddata1.rda"
[4] "/path/to/R/adegenet/files/monddata2.rda"
[5] "/path/to/R/adegenet/files/nancycats.dat"
[6] "/path/to/R/adegenet/files/nancycats.gen"
[7] "/path/to/R/adegenet/files/nancycats.gtx"
[8] "/path/to/R/adegenet/files/nancycats.str"
[9] "/path/to/R/adegenet/files/pdH1N1-data.csv"
[10] "/path/to/R/adegenet/files/pdH1N1-HA.fasta"
[11] "/path/to/R/adegenet/files/pdH1N1-NA.fasta"
[12] "/path/to/R/adegenet/files/usflu.fasta"

$path
[1] "/path/to/R/adegenet/files"

```

We can see that we have a mix of files with different formats. If we tried to run all of these files using `poppr`, we would have a problem because some of the file formats have no direct import into a `genind` object (*.fasta, or *.snp), or just simply are not supported (eg. *.rda files). We want to be able to filter these files out, and we will do so with the `pattern` argument. Let's say we only wanted the files that have the word "nancy" in them.

```
> getfile(multi=TRUE, pattern="nancy")
```

```

$files
[1] "/path/to/R/adegenet/files/nancycats.dat" "/path/to/R/adegenet/files/nancycats.gen"
[3] "/path/to/R/adegenet/files/nancycats.gtx" "/path/to/R/adegenet/files/nancycats.str"

$path
[1] "/path/to/R/adegenet/files"

```

Now, let's exclude everything but gentix files (*.gtx).


```
> getfile(multi=TRUE, pattern="gtx")

$files
[1] "/path/to/R/adegenet/files/nancycats.gtx"

$path
[1] "/path/to/R/adegenet/files"
```

Now, let's only get FSTAT files (*.dat)

```
> getfile(multi=TRUE, pattern="dat")

$files
[1] "/path/to/R/adegenet/files/monddata1.rda"
[2] "/path/to/R/adegenet/files/monddata2.rda"
[3] "/path/to/R/adegenet/files/nancycats.dat"
[4] "/path/to/R/adegenet/files/pdH1N1-data.csv"

$path
[1] "/path/to/R/adegenet/files"
```

Uh-oh. We've run into a problem. Three out of our four files are not FSTAT files. Why did this happen? It happened because they happen to have "dat" within their name. This problem can be solved, by using regular expressions. If you are unfamiliar with regular expressions, you can think of them as special characters that you can use to make your search pattern more strict or more flexible. Since the topic of regular expressions can take up several lectures, I will spare you the gory details. For this situation, the only one you need to know is "\$". The dollar sign indicates the end of a word or string. If we want specific file extensions all we have to do is add this to the end of the search term like so:

```
> getfile(multi=TRUE, pattern="dat$")

$files
[1] "/path/to/R/adegenet/files/nancycats.dat"

$path
[1] "/path/to/R/adegenet/files"
```

Now we have our FSTAT file!

1.4.2 Function: read.genalex

A very popular program for population genetics is GenAlEx (<http://biology.anu.edu.au/GenAlEx/Welcome.html>) [14, 13]. GenAlEx runs within the Excel environment and can be very powerful in its analyses. *Poppr* has added the ability to read *.CSV files³ produced in the GenAlEx format. It can handle data types containing regions and geographic coordinates, but currently it cannot import allelic frequency data from GenAlEx. All the user has to do is to export a single sheet of GenAlEx data from Excel into a *.CSV file, and the *poppr* function `read.genalex` will import it into *adegenet*'s `genind` object (more information on that below). For ways of formatting a GenAlEx file, see the manual here: http://biology.anu.edu.au/GenAlEx/Download_files/GenAlEx%206.5%20Guide.pdf

Default Command:

```
read.genalex(genalex, ploidy = 2, geo = FALSE, region = FALSE)
```

- `genalex` - a *.CSV file exported from GenAlEx on your disk (For example: "my_genalex_file.csv").
- `ploidy` - a number indicating the ploidy for the data set (eg 2 for diploids, 1 for haploids).

³*.CSV files are comma separated files that are easily machine readable.

- **geo** - GenAlEx allows you to have geographic data within your file. To do this for *poppr*, you will need to follow the first format outlined in the GenAlEx manual and place the geographic data AFTER all genetic and demographic data with one blank column separating it (See the GenAlEx Manual for details). If you have geographic information in your file, set this flag to **TRUE** and it will be included within the resulting *genind* object in the **@other** slot. (If you don't know what that is, don't worry. It will be explained later in section 1.4.4).
- **region** - To format your GenAlEx file to include regions along with your populations, You can choose to include a separate column for regional data, or, since regional data must be in contiguous blocks, you can simply format it in the same way you would any other data (see the GenAlEx manual for details). If you have your file organized in this manner, select this option and the regional information will be stored in the resulting *genind* object in the **@other** slot.

IF YOU ARE UNFAMILIAR WITH EXPORTING DATA FROM EXCEL

1. Click the Microsoft Office Button in the top left corner of Excel. (Or go to the File menu if you have an older version)
2. Click Save As...
3. In the "Save as type" drop down box, select CSV (comma delimited).

Note that regional data and geographic data are not mutually exclusive. You can have both in one file, just make sure that they are on the same sheet and that the geographic data is always placed after all genetic and demographic data.

We have a short example of *genalex* formatted data with no geographic or regional formatting. We will first see where the data is using the command `system.file()`

```
> system.file("files/rootrot.csv", package="poppr")
```

```
[1] "/path/to/R/library/poppr/files/rootrot.csv"
```

Now import the data into *poppr* like so:

```
> rootrot <- read.genalex(system.file("files/rootrot.csv", package="poppr"))
```

Executing *rootrot* shows that this file is now in *genind* format (ie. the format required by *poppr* and *adegenet*).

```
> rootrot

#####
### Genind object ###
#####
- genotypes of individuals -

S4 class: genind
@call: read.genalex(genalex = "rootrot.csv")

@tab: 187 x 56 matrix of genotypes

@ind.names: vector of 187 individual names
@loc.names: vector of 56 locus names
@loc.nall: NULL
@loc.fac: NULL
@all.names: NULL
@ploidy: 2
@type: PA

Optionnal contents:
@pop: factor giving the population of each individual
@pop.names: factor giving the population of each individual

@other: a list containing: population_hierarchy
```

1.4.3 Genalex formatting shortcuts

The GenAlEx format is a nice way to import data because it allows you to have geographic coordinates and two hierarchical levels of sampling (Region and population). If you have multiple levels of hierarchy, you will need to code them so that you combine multiple columns of hierarchy into one using a common separator (For an example, see section 2.2.1 of this manual). A problem arises when it becomes more work than it's worth to do that since, for the GenAlEx format, you must provide the sizes of each population in the header. Here, I'll show you a simple way to circumvent that. First, let's use the microbov data set from *adegenet* (for details, type `help("microbov")` into your R console). It contains three demographic factors: Country, Species and Breed contained within the `@other` slot (detailed in section 1.5.1). We will combine these and save the file to our desktop. We will cover these functions later in this manual. For now, just know they exist.

```
> library(poppr)
> data(microbov)
> microbov@other$population_hierarchy <- data.frame(list(Country = microbov@other$coun,
+ Species = microbov@other$spe, Breed = microbov@other$breed))
> microbov <- splitcombine(microbov, method=2, hier=c("Country", "Species", "Breed"))
> genind2genalex(microbov, file=~ /Desktop/microbov.csv")
```

Extracting the table ... Writing the table to ~/Desktop/microbov.csv ... Done.

After we do this, we can open the file in our favorite spreadsheet editor and see the following image.

Figure 2: The first 15 individuals and 4 loci of the microbov data set. The first column contains the individual names, the second column contains the population names, and each subsequent column represents microsatellite genetic data. Highlighted in red is a list of populations and their relative sizes.

	A	B	C	D	E	F	G	H	I	I
1	30	704	15	50	50	51	30	50	50	47
2	Unmodified Data		AF BI Borgou	AF BI Zebu	AF BT Lagunaire	AF BT NDama	AF BT Somba	FR BT Aubrac	FR BT Bazadais	
3	Ind	Pop	INRA63		INRA5		ETH225		ILSTS5	
4	AFBIBOR9503	AF BI Borgou	183	183	137	141	147	157	190	190
5	AFBIBOR9504	AF BI Borgou	181	183	141	141	139	157	186	186
6	AFBIBOR9505	AF BI Borgou	177	183	141	141	139	139	194	194
7	AFBIBOR9506	AF BI Borgou	183	183	141	141	141	147	184	190
8	AFBIBOR9507	AF BI Borgou	177	183	141	141	153	157	184	186
9	AFBIBOR9508	AF BI Borgou	177	183	137	143	149	157	184	186
10	AFBIBOR9509	AF BI Borgou	177	181	139	141	147	157	184	190
11	AFBIBOR9510	AF BI Borgou	183	183	139	141	155	157	184	186
12	AFBIBOR9511	AF BI Borgou	177	183	139	141	139	143	182	190
13	AFBIBOR9512	AF BI Borgou	183	183	141	141	157	159	186	186
14	AFBIBOR9513	AF BI Borgou	177	177	141	141	147	157	184	190
15	AFBIBOR9514	AF BI Borgou	183	183	143	143	139	157	186	186
16	AFBIBOR9515	AF BI Borgou	183	183	137	137	143	157	0	0
17	AFBIBOR9516	AF BI Borgou	177	183	137	143	139	157	182	184
18	AFBIBOR9517	AF BI Borgou	177	183	141	141	157	157	186	194

All that *poppr* needs from the first header row are the first three numbers (unless you are including regional data, but it's not terribly necessary with the hierarchical support *poppr* provides.), which represent the number of loci, individuals, and populations, respectively. After that, you have counts of individuals per population in each subsequent cell. For *poppr*, These cells don't matter because we already have that information in column 2.

If you have a large data set with many population levels, you can use the following shortcut by setting the number in the third cell to 1. The number in cell 4 is arbitrary (but must be there). In the following figure, it is set to the number of individuals in your data set, but can easily be replaced with any other number (perhaps your favorite number?).

Figure 3: The first 15 individuals and 4 loci of the microbov data set. This is the same figure as above, however the populations and counts have been removed from the header row and the third number in the header has been replaced by 1.

	A	B	C	D	E	F	G	H	I	J
1	30	704	1	704						
2	Example Modified Data			ALL						
3	Ind	Pop	INRA63	INRA5			ETH225		ILSTS5	
4	AFBIBOR9503	AF_BI_Borgou	183	183	137	141	147	157	190	190
5	AFBIBOR9504	AF_BI_Borgou	181	183	141	141	139	157	186	186
6	AFBIBOR9505	AF_BI_Borgou	177	183	141	141	139	139	194	194
7	AFBIBOR9506	AF_BI_Borgou	183	183	141	141	141	147	184	190
8	AFBIBOR9507	AF_BI_Borgou	177	183	141	141	153	157	184	186
9	AFBIBOR9508	AF_BI_Borgou	177	183	137	143	149	157	184	186
10	AFBIBOR9509	AF_BI_Borgou	177	181	139	141	147	157	184	190
11	AFBIBOR9510	AF_BI_Borgou	183	183	139	141	155	157	184	186
12	AFBIBOR9511	AF_BI_Borgou	177	183	139	141	139	143	182	190
13	AFBIBOR9512	AF_BI_Borgou	183	183	141	141	157	159	186	186
14	AFBIBOR9513	AF_BI_Borgou	177	177	141	141	147	157	184	190
15	AFBIBOR9514	AF_BI_Borgou	183	183	143	143	139	157	186	186
16	AFBIBOR9515	AF_BI_Borgou	183	183	137	137	143	157	0	0
17	AFBIBOR9516	AF_BI_Borgou	177	183	137	143	139	157	182	184
18	AFBIBOR9517	AF_BI_Borgou	177	183	141	141	157	157	186	194

1.4.4 Other ways of importing data

Adegenet already supports the import of FSTAT, Structure, Genpop, and Gentix formatted files, so if you have those formats, you can import them using the function `import2genind`. For sequence data, check if you can use `read.dna` from the *ape* package to import your data. If you can, then you can use the *adegenet* function `DNABin2genind`. If you don't have any of these formats handy, you can still import your data using R's `read.table` along with `df2genind` from *adegenet*. For more information, see *adegenet*'s "Getting Started" vignette.

1.4.5 Function: `genind2genalex`

Of course, being able to export data is just as useful as being able to import it, so we have this handy little function that will write a GenAlEx formatted file to wherever you desire.

WARNING: This will overwrite any file that exists with the same name.

Default Command:

```
genind2genalex(pop, filename = "genalex.csv", quiet = FALSE)
```

- `pop` - a `genind` object.
- `filename` - This is where you specify where you want the file to go. If you simply type the file name, it will deposit the file in the directory R is currently in. If you don't know what directory you are in, you can type `getwd()` to find out.

- **quiet** - If this is set to **FALSE**, a message will be printed to the screen.
- **geo** - This is set to **FALSE** by default. If it is set to **TRUE**, then that means you have a data frame or matrix in the **@other** slot of your **genind** object that contains geographic coordinates for all individuals or all populations. Setting this to **TRUE** means that you want the resulting file to have two extra columns at the end of your file with geographic coordinates.
- **geodf** - The name of the data frame or matrix containing the geographic coordinates. The default is **geodf = "xy"**.

First, a simple example for the **rootrot** data we demonstrated in section 1.4.2:

```
> genind2genalex(rootrot, "~/Desktop/rootrot.csv")
```

Extracting the table ... Writing the table to ~/Desktop/rootrot.csv ... Done.

Now here's an example of exporting the **nancycats** data set into **GenAlEx** format with geographic information. If we look at the **nancycats** geographic information, we can see it's coordinates for each population, but not each individual:

```
> data(nancycats)
> nancycats@other$xy
```

```

      x      y
P01 263.3498 171.10939
P02 183.5028 122.40790
P03 391.1050 254.70148
P04 458.6121  41.72336
P05 182.7769 219.08398
P06 335.2121 344.83557
P07 359.1662 375.36486
P08 271.3345  67.89132
P09 256.8169 150.02964
P10 270.6086  17.00917
P11 493.4544 237.25618
P12 305.4510  85.33663
P13 462.9674  86.79040
P14 429.5768 291.04587
P15 531.2003 115.13903
P16 407.8003  99.87438
P17 345.3745 251.79393
```

And we can export it easily:

```
> genind2genalex(nancycats, "~/Desktop/nancycats_pop_xy.csv")
```

Extracting the table ... Writing the table to ~/Desktop/nancycats_pop_xy.csv ... Done.

If we wanted to assign a geographic coordinate to each individual, we can simply use this little repetition trick knowing that there are 17 populations in the data set:

```
> nan2 <- nancycats
> nan2@other$xy <- nan2@other$xy[rep(1:17, table(pop(nan2))), ]
> head(nan2@other$xy)
```

```

      x      y
P01 263.3498 171.1094
P01 263.3498 171.1094
P01 263.3498 171.1094
P01 263.3498 171.1094
P01 263.3498 171.1094
P01 263.3498 171.1094
```

Now we can export it to a different file.

```
> genind2genalex(nan2, "~/Desktop/nancycats_inds_xy.csv")
```

Extracting the table ... Writing the table to ~/Desktop/nancycats_inds_xy.csv ... Done.

1.5 Getting to know adegenet's genind object

Since *poppr* was built around adegenet's framework, it is important to know how *adegenet* stores data in the *genind* object, as that is the object used by *poppr*. To create a *genind* object, *adegenet* takes a data frame of genotypes (rows) across multiple loci (columns) and converts them into a matrix of individual allelic frequencies at each locus [9].

For example, if you had a data frame with 3 diploid individuals each with 3 loci that had 3, 4, and 5 allelic states respectively, the resulting *genind* object would contain a matrix that has 3 rows and 12 columns. Here's the example data frame:

```
  locus1 locus2 locus3
1 101/101 201/201 301/302
2 102/103 202/203 301/303
3 102/102 203/204 304/305
```

And the resulting matrix after importing to *genind*.

```
  L1.1 L1.2 L1.3 L2.1 L2.2 L2.3 L2.4 L3.1 L3.2 L3.3 L3.4 L3.5
1    1    0    0    1    0    0    0    0.5  0.5  0.0  0.0  0.0
2    0    0.5  0.5    0    0.5  0.5  0.0  0.5  0.0  0.5  0.0  0.0
3    0    1.0  0.0    0    0.0  0.5  0.5  0.0  0.0  0.0  0.5  0.5
```

The first three columns represent the alleles of locus 1, the next four represent locus 2, and the last five represent locus 3.

Do you see what I mean when I say individual allele frequencies at each locus? For a diploid individual, you only have three possible allele frequencies at each locus: 1, 0.5, or 0. Now, this is not the entire *genind* object, but it is the main feature. The object also has various elements associated with it that give you information about the population membership, the names of loci, individuals, and alleles among other things that *poppr* uses to work [9]. If you wish to know more, see the *adegenet* "Getting Started" manual.

1.5.1 The other slot

The element that you as a *poppr* user needs to be concerned with is the "other" slot. No, I'm not trying to be cryptic. If you look at an *adegenet* object, you will see that it has several "slots" (starting with "@"). [9] Let's start by recreating that data frame I showed you earlier.

```
> df <- data.frame(list(locus1=c("101/101", "102/103", "102/102"),
+                        locus2=c("201/201", "202/203", "203/204"),
+                        locus3=c("301/302", "301/303", "304/305")
+                    )
+                  )
> dfg <- df2genind(df, sep="/")
```

Next we will display the contents of the *genind* object *dfg*

```
> dfg

#####
### Genind object ###
#####
- genotypes of individuals -

S4 class:   genind
@call: df2genind(X = df, sep = "/")

@tab:  3 x 12 matrix of genotypes

@ind.names: vector of  3 individual names
@loc.names: vector of  3 locus names
@loc.nall: number of alleles per locus
@loc.fac: locus factor for the 12 columns of @tab
@all.names: list of  3 components yielding allele names for each locus
@ploidy: 2
@type: codom

Optionnal contents:
@pop: - empty -
@pop.names: - empty -

@other: - empty -
```

The matrix containing our allelic frequencies is located in the `@tab` slot. All of the slots below that have very specific properties related to the matrix in `@tab`, but the `@other` slot is more or less a grab bag, where you can place anything you want, even if it doesn't make sense!

Here, I'll give you an example of placing the `genind` object inside itself. Notice first, that the `@other` slot is empty and pay attention to the commands I use, noting that you can use either "\$" or "@" to access the slots.

```
> # First off, how big is the object?
> print(object.size(dfg), units="auto")
```

8.5 Kb

```
> dfg$other$dfg <- dfg
> dfg # we can now see that the @other slot is now filled.
```

```
#####
### Genind object ###
#####
- genotypes of individuals -

S4 class:  genind
@call: df2genind(X = df, sep = "/")

@tab:  3 x 12 matrix of genotypes

@ind.names: vector of  3 individual names
@loc.names: vector of  3 locus names
@loc.nall: number of alleles per locus
@loc.fac: locus factor for the  12 columns of @tab
@all.names: list of  3 components yielding allele names for each locus
@ploidy:  2
@type:  codom

Optionnal contents:
@pop:  - empty -
@pop.names:  - empty -

@other: a list containing: dfg
```

```
> dfg$other$dfg
```

```
#####
### Genind object ###
#####
- genotypes of individuals -

S4 class:  genind
@call: df2genind(X = df, sep = "/")

@tab:  3 x 12 matrix of genotypes

@ind.names: vector of  3 individual names
@loc.names: vector of  3 locus names
@loc.nall: number of alleles per locus
@loc.fac: locus factor for the  12 columns of @tab
@all.names: list of  3 components yielding allele names for each locus
@ploidy:  2
@type:  codom

Optionnal contents:
@pop:  - empty -
@pop.names:  - empty -

@other: - empty -
```

```
> print(object.size(dfg), units="auto") # How big is it now?
```

17.2 Kb

WHAT IS THE # SIGN FOR?

This is called a comment. If you type something in R with the “#” sign in front of it, R will not interpret it.

And we can continue to do this until we reach the limit of our available memory. Why am I showing this silliness to you? For one thing I want to show you that you can stick anything you want into that slot and the object will not be hurt in any way. It’s also important when considering how you are going to deal with the population structure of your *genind* object. For the *poppr* functions *clonecorrect* (Section 2.4.1) and *splitcombine* (Section 2.2.1) to work, a data frame of the population hierarchy must be present in the *@other* slot and it must have the same number of rows as individuals in the data set. There are several ways to go about this. If you know how to create a data frame or import data into R, the command is no more difficult than `obj$other$population_hierarchy <- df`. If you do not know how to create a data frame or import data into R, you can visit Quick R at <http://www.statmethods.net/input/importingdata.html>.

1.5.2 Setting the population factor {adegenet’s function: pop}

A *genind* object can contain several populations, and, if you have differing population structures, you might want to switch among them for different analyses. The tools you as the user would need, are the slot *@pop.names* and the *adegenet* function *pop()*. I’ll use the H3N2 data set packaged with *adegenet* as an example.

```
> data(H3N2)
> H3N2

#####
### Genind object ###
#####
- genotypes of individuals -

S4 class:   genind
@call:    .local(x = x, i = i, j = j, drop = drop)

@tab:    1903 x 334 matrix of genotypes

@ind.names: vector of 1903 individual names
@loc.names: vector of 125 locus names
@loc.nall: number of alleles per locus
@loc.fac: locus factor for the 334 columns of @tab
@all.names: list of 125 components yielding allele names for each locus
@ploidy: 1
@type:   codom

Optionnal contents:
@pop:    - empty -
@pop.names: - empty -

@other: a list containing: x  xy  epid

> pop(H3N2)

NULL

> H3N2$pop.names

NULL
```

Notice how both the *pop* and *pop.names* are empty. This means that the population information needs to be set. Notice, however that there are 1903 individuals in the data set and that the *@other* slot is not empty. Let’s investigate an object in this slot.

```
> head(H3N2$other$x)
```


	accession	length	host	segment	subtype	country	year	Virus name
AB434107	AB434107	1701	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
AB434108	AB434108	1701	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
AB438242	AB438242	827	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
AB438243	AB438243	827	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
AB438244	AB438244	827	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
AB438245	AB438245	827	Human	4 (HA)	H3N2	Japan	2002	Influenza A virus
	Misc info	Age	Gender			date	usePreciseLoc	localisation
AB434107	(A/Morioka/34/2002(H3N2))	<NA>	<NA>			2002/02/25	FALSE	japan
AB434108	(A/Morioka/52/2002(H3N2))	<NA>	<NA>			2002/03/01	FALSE	japan
AB438242	(A/Niigata/F11/2002(H3N2))	<NA>	<NA>			2002/01/22	FALSE	japan
AB438243	(A/Niigata/F100/2002(H3N2))	<NA>	<NA>			2002/02/18	FALSE	japan
AB438244	(A/Niigata/F175/2002(H3N2))	<NA>	<NA>			2002/02/25	FALSE	japan
AB438245	(A/Niigata/F245/2002(H3N2))	<NA>	<NA>			2002/03/04	FALSE	japan
	lon	lat	month					
AB434107	137.2155	35.58418	2					
AB434108	137.2155	35.58418	3					
AB438242	137.2155	35.58418	1					
AB438243	137.2155	35.58418	2					
AB438244	137.2155	35.58418	2					
AB438245	137.2155	35.58418	3					

```
> nrow(H3N2$other$x)
```

```
[1] 1903
```

WHAT IS HEAD()?

`head()` is a command that will show you only the top portion of an R object. By default it will show you the first six elements (or rows of a data frame or matrix). This is so that you can quickly check the contents of an object.

We can see that it's a data frame containing a wealth of information that we could use to subset our data. So, let's start by setting the population structure by country. How do we do that? Well, the function `pop()` will allow us to set that structure using a vector that is the same length as the number of individuals in the data set. Since the number of rows in the data frame `x` meets that criteria, we can use any item in that data frame. Let's take a look.

```
> pop(H3N2) <- H3N2$other$x$country
> head(pop(H3N2))
```

```
[1] Japan Japan Japan Japan Japan Japan
37 Levels: Japan USA Finland China South Korea Norway Taiwan France ... Algeria
```

```
> H3N2$pop.names
```

[1] "Japan"	"USA"	"Finland"	"China"	"South Korea"
[6] "Norway"	"Taiwan"	"France"	"Latvia"	"Netherlands"
[11] "Bulgaria"	"Turkey"	"United Kingdom"	"Denmark"	"Austria"
[16] "Canada"	"Italy"	"Russia"	"Bangladesh"	"Egypt"
[21] "Germany"	"Romania"	"Ukraine"	"Czech Republic"	"Greece"
[26] "Iceland"	"Ireland"	"Sweden"	"Nepal"	"Saudi Arabia"
[31] "Switzerland"	"Iran"	"Mongolia"	"Spain"	"Slovenia"
[36] "Croatia"	"Algeria"			

Notice how useful the `@other` slot can be. We now have population structure in the data set and you now know how to set the population factor. The other slot will become useful later on when we are talking about multilocus genotypes.

2 Data Manipulation

One tedious aspect of population genetic analysis is the need for repeated data manipulation. *Adegenet* has some functions for manipulating data that are limited to replacing missing data and dividing data into populations, loci, or by sample size [9]. *Poppr* includes novel functions for clone-censoring your data sets or sub-setting a `genind` object by specific populations.

2.1 Inside the golden days of missing data {replace or remove missing data}

A data set without missing data is always ideal, but often not achievable. Many functions in *adegenet* cannot handle missing data and thus the function `na.replace` exists [9]. It will replace missing data with either “0” representing a mysterious extra allele in the data set resulting in more diversity or the mean of allelic frequencies at the locus. There is no set method, however, for simply removing missing data from analyses, which is why the *poppr* function `missingno` (see below) exists. If the name makes you uneasy it’s because it should. Missing data can mean different things based on your data type. For microsatellites, missing data might represent any source of error that could cause a PCR product to not amplify in gel electrophoresis, which may or may not be biologically relevant. For a DNA alignment, missing data could mean something as simple as an insertion or deletion, which is biologically relevant. The choice to exclude or estimate data has very different implications for the type of data you have.

2.1.1 Function: `missingno`

`missingno` is a function that serves partially as a wrapper for *adegenet*’s `na.replace` to replace missing data and as a way to exclude specific areas that contain systematic missing data.

Default Command:

```
missingno(pop, type = "loci", cutoff = 0.05, quiet = FALSE)
```

- `pop` - a *genind* object.
- `type` - This could be one of four options:
 - “**mean**” This replaces missing data with the mean allele frequencies in the entire data set.
 - “**zero**” or “**0**” This replaces missing data with zero, signifying a new allele.
 - “**loci**” This is to be used for a data set that has systematic problems with certain loci that contain null alleles or simply failed to amplify. This will remove loci with a defined threshold of missing data from the data set.
 - “**geno**” This is to be used for genotypes (individuals) in your data set where many null alleles are present. Individuals with a defined threshold missing data will be removed.
- `cutoff` - This is a numeric value from 0 to 1 indicating the percent allowable missing data for either loci or genotypes. If you have, for example, two loci containing missing 5% and 10% missing data, respectively and you set `cutoff = 0.05`, `missingno` will remove the second locus. Percent missing data for genotypes is considered the percent missing loci over number of total loci.
- `quiet` - When this is set to `FALSE`, the number of missing values replaced will be printed to screen if the method is “zero” or “mean”. It will print the number of loci or individuals removed if the method is “loci” or “geno”.

Of course, seeing is believing. Let’s take a look at what this does by focusing in on areas with missing data. Note that I will be using some sub-setting functions here that are described in *adegenet*’s *Getting Started* vignette. First, let’s take a look at what the missing data in R looks like as well as how many loci and individuals the data set *nancycats* contains. We need to first tell R to look in its library for the package *poppr*.

```
> library(poppr)
```

Next, we’ll initialize the *adegenet* data set *nancycats* and load it into memory.

```
> data(nancycats)
```

Now, we’ll take a quick look at the *nancycats* data set using *adegenet*’s `summary()` function:

```
> summary(nancycats)

# Total number of genotypes: 237

# Population sample sizes:
 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17
10 22 12 23 15 11 14 10  9 11 20 14 13 17 11 12 13

# Number of alleles per locus:
L1 L2 L3 L4 L5 L6 L7 L8 L9
16 11 10  9 12  8 12 12 18

# Number of alleles per population:
01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17
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:
      L1      L2      L3      L4      L5      L6      L7      L8      L9
0.6682028 0.6666667 0.6793249 0.7083333 0.6329114 0.5654008 0.6497890 0.6184211 0.4514768

# Expected heterozygosity:
      L1      L2      L3      L4      L5      L6      L7      L8      L9
0.8657224 0.7928751 0.7953319 0.7603095 0.8702576 0.6884669 0.8157881 0.7603493 0.6062686
```

We can see here a lot of summary statistics about nancycats. Here we can see that there are 17 populations, 237 individuals, and 9 loci. Nancycats also has a little over 2.3% missing data. Let's take a look at the names of the loci and the structure of the data. In order to save space, I will only show you the first five individuals (rows) and a portion of the alleles in the first locus (columns).

```
> nancycats$loc.names # Names of the loci

      L1      L2      L3      L4      L5      L6      L7      L8      L9
"fca8" "fca23" "fca43" "fca45" "fca77" "fca78" "fca90" "fca96" "fca37"

> nancycats$tab[1:5, 8:13]

      L1.08 L1.09 L1.10 L1.11 L1.12 L1.13
001      NA      NA      NA      NA      NA      NA
002      NA      NA      NA      NA      NA      NA
003      0.0      0.5      0      0      0      0.5
004      0.5      0.5      0      0      0      0.0
005      0.5      0.5      0      0      0      0.0
```

When looking at this data set, recall how a `genind` object is formatted. You have a matrix of 0's, 1's and 0.5's. For diploids, if you see 0.5, that means it is heterozygous at that allele, and a 1 means it's homozygous. Here we see three heterozygotes and two individuals with missing data (indicated by NA). Now, there are more places with missing data in the data set, but I'm only showing a little bit at one locus so it's easier to digest. Let's first replace it by zero and mean, respectively.

```
> nanzero <- missingno(nancycats, type = "zero")
```

```
Replaced 617 missing values
```

```
> nanmean <- missingno(nancycats, type = "mean")
```

```
Replaced 617 missing values
```

```
> nanzero$tab[1:5, 8:13]
```

	L1.08	L1.09	L1.10	L1.11	L1.12	L1.13
001	0.0	0.0	0	0	0	0.0
002	0.0	0.0	0	0	0	0.0
003	0.0	0.5	0	0	0	0.5
004	0.5	0.5	0	0	0	0.0
005	0.5	0.5	0	0	0	0.0

```
> nanmean$tab[1:5, 8:13]
```

	L1.08	L1.09	L1.10	L1.11	L1.12	L1.13
001	0.07603687	0.2419355	0.1912442	0.06221198	0.09447005	0.1013825
002	0.07603687	0.2419355	0.1912442	0.06221198	0.09447005	0.1013825
003	0.00000000	0.5000000	0.0000000	0.00000000	0.00000000	0.5000000
004	0.50000000	0.5000000	0.0000000	0.00000000	0.00000000	0.0000000
005	0.50000000	0.5000000	0.0000000	0.00000000	0.00000000	0.0000000

You notice how the values of NA changed, yet the basic structure stayed the same. These are the replacement options from adegenet. Let's look at the same example with the exclusion options (set to the default cutoff of 5%).

```
> nanloci <- missingno(nancycats, "loci")
```

```
Found 617 missing values.
2 loci contained missing values greater than 5%.
Removing 2 loci : fca8 fca45
```

```
> nangeno <- missingno(nancycats, "geno")
```

```
Found 617 missing values.
38 genotypes contained missing values greater than 5%.
Removing 38 genotypes : N215 N216 N188 N189 N190 N191 N192 N302 N304 N310
N195 N197 N198 N199 N200 N201 N206 N182 N184 N186 N298 N299 N300 N301 N303 N282
N283 N288 N291 N292 N293 N294 N295 N296 N297 N281 N289 N290
```

```
> nanloci$tab[1:5, 8:13]
```

	L1.08	L1.09	L1.10	L1.11	L2.01	L2.02
001	0	0.5	0	0	0	0
002	0	1.0	0	0	0	0
003	0	0.5	0	0	0	0
004	0	0.0	0	0	0	0
005	0	0.5	0	0	0	0

Notice how we now see columns named “L2.01” and “L2.02”. This is showing us another locus because we have removed the first. Recall from the summary table that the first locus had 16 alleles, and the second had 11. Now that we’ve removed loci containing missing data, all others have shifted over. Let’s look at the loci names and number of individuals.

```
> length(nanloci$ind.names) # Individuals
```

```
[1] 237
```

```
> nanloci$loc.names # Names of the loci
```

L1	L2	L3	L4	L5	L6	L7
"fca23"	"fca43"	"fca77"	"fca78"	"fca90"	"fca96"	"fca37"

You can see that the number of individuals stayed the same but the loci “fca8”, “fca45”, and “fca96” were removed.

Let’s look at what happened when we removed individuals.

```
> nangeno$tab[1:5, 8:13]
```

	L1.08	L1.09	L1.10	L1.11	L1.12	L1.13
001	0.0	0.5	0	0	0	0.5
002	0.5	0.5	0	0	0	0.0
003	0.5	0.5	0	0	0	0.0
004	0.0	0.5	0	0	0	0.5
005	0.0	1.0	0	0	0	0.0

```
> length(nangeno$ind.names) # Individuals
```

```
[1] 199
```

```
> nangeno$loc.names # Names of the loci
```

L1	L2	L3	L4	L5	L6	L7	L8	L9
"fca8"	"fca23"	"fca43"	"fca45"	"fca77"	"fca78"	"fca90"	"fca96"	"fca37"

We can see here that the number of individuals decreased, yet we have the same number of loci. Notice how the frequency matrix changes in both scenarios? In the scenario with “loci”, we removed several columns of the data set, and so with our sub-setting, we see alleles from the second locus. In the scenario with “geno”, we removed several rows of the data set so we see other individuals in our sub-setting.

2.2 Can you take me hier(archy)? {population hierarchy construction}

Remember all that fuss we made about the `@other` slot above in section 1.4.4? The way you can achieve hierarchical analysis in *poppr* is through a data frame in that slot. Many of the file formats that *adegenet* and *poppr* can import do not allow for more than two hierarchies. If you need more levels, you have a couple of choices:

1. Import them into R as a data frame with each column being a separate hierarchical element.
2. Collapse them into a single population factor so that you can trick these file formats into taking multiple population hierarchies (eg. instead of “Pop1”, “Subpop1”, “Subsubpop1”, you would have “Pop1_Subpop1_Subsubpop1”).

Whichever choice you make, The *poppr* function `splitcombine` can help you divide and combine those factors in any way you can think of.

2.2.1 Function: splitcombine

This function will allow you to combine your population hierarchies in ways meaningful to your data without needing to know R programming. It can either split a vector of combined population hierarchies or it can combine columns of a data frame containing population hierarchies (Note that it will only split the first column of the data frame if you choose `method = 1`).

Default Command:

```
splitcombine(pop, method = 1, dfname = "population_hierarchy", sep = "_", hier = c(1),
setpopulation = TRUE, fixed = TRUE)
```

- `pop` - a `genind` object with a data frame in the `@other` slot.
- `method` - An integer indicating what you want to do on your data frame:
 1. **split** Any populations combined using a common separator in your data frame. So, a population hierarchy of “Pop1_Subpop1_Subsubpop1” would be split into a data frame containing the columns “Pop1”, “Subpop1”, “Subsubpop1”. Since it will split the population factor, it needs only to be used once.

2. **combine** If you have your population hierarchy split into a data frame, you can do the exact opposite of method 1 and combine separate elements into one.

- **dfname** - This is the name of the data frame containing your population factor. Note that you are not limited to one data frame in your `genind` object. If you do not have that data frame in the `@other` slot, a warning will be returned and nothing will happen.
- **sep** - A separation factor you want to separate your populations with. Note, that you can choose whatever you want, but be careful because some characters have special meanings (regular expressions) in R and could give you incorrect results (“_” is the suggested default).
- **hier** - This can be a vector of words or numbers referring to what you want to name your population hierarchies in `method = 1`, or specific column names in your data frame in `method = 2`.
- **setpopulation** - if `TRUE` (default), this will automatically set the population factor to either the highest population factor (with `method = 1`, split) or the combined population hierarchy (with `method = 2`, combine). if this is set to `FALSE`, the population factor will not be set.
- **fixed** - This is an option to be passed onto the *base* function `strsplit`. For those not familiar with regular expressions, it will tell R whether or not the character in `split` should be treated as a special character or not. If you don’t know regular expressions, don’t touch it.

Let’s give an example using AFLP data of different populations of *A. euteiches* collected in Washington and Oregon. [5]

```
> Aeut <- read.genalex(system.file("files/rootrot.csv", package="poppr"))
> summary(Aeut)

# Total number of genotypes: 187

# Population sample sizes:
  Athena_1  Athena_2  Athena_3  Athena_4  Athena_5  Athena_6
    9        12        10        13        10         5
  Athena_7  Athena_8  Athena_9  Athena_10 Mt. Vernon_1 Mt. Vernon_2
   11         8        10         9         10         6
Mt. Vernon_3 Mt. Vernon_4 Mt. Vernon_5 Mt. Vernon_6 Mt. Vernon_7 Mt. Vernon_8
    8         12        17        12        12        13

# Percentage of missing data:
[1] 0
```

DOES THIS SUMMARY SEEM A LITTLE LACKING?

The data that we have here is presence absence data. This means that many of the functions that *adegenet* uses to calculate heterozygosity and number of alleles are slightly useless in this regard.

Notice that we have 18 different “populations” here, but they are really a hierarchy. Let’s say we want to analyze the diversity statistics of the two overall populations. Take a look at how the combined population factor is kept in the data frame.

```
> head(Aeut$other$population_hierarchy)
```

```
      Pop
1 Athena_1
2 Athena_1
3 Athena_1
4 Athena_1
5 Athena_1
6 Athena_1
```

We’ll use `splitcombine` to split that into a population and sub-population and set the population factor to the population.

IMPORTANT POINT ABOUT SPLITCOMBINE

Ideally, method split should only be used once after you read in your data. The reason for this is that when you select this method, it will look in the first column of your data frame to choose the combined population factor to split. If you do not name your hierarchy or if you attempt to give your hierarchy too many names, it will automatically name the columns “h1”, “h2”, etc.

```
> Aeut.pop <- splitcombine(Aeut, method=1, dfname="population_hierarchy", hier=c("Pop", "Subpop"), setpopulation=TRUE)
> head(Aeut.pop$other$population_hierarchy)
```

```
Pop_Subpop  Pop Subpop
1  Athena_1 Athena    1
2  Athena_1 Athena    1
3  Athena_1 Athena    1
4  Athena_1 Athena    1
5  Athena_1 Athena    1
6  Athena_1 Athena    1
```

```
> summary(Aeut.pop)
```

```
# Total number of genotypes: 187
# Population sample sizes:
  Athena Mt. Vernon
    97      90
# Percentage of missing data:
[1] 0
```

Now we can see that we have a data frame with all of our population factors separated, and we still have our original combined hierarchy, but it is now called “Pop_Subpop”. This allows you to keep track of what you named your population hierarchies. We can now run the function `poppr` to get a diversity analysis.

```
> poppr(Aeut.pop, quiet=TRUE)
```

```
Pop      N MLG  eMLG  SE    H      G  Hexp  E.5    Ia rbarD      File
1  Athena  97  70 65.981 1.246 4.063 42.193 0.986 0.721  2.906 0.072 rootrot.csv
2 Mt. Vernon 90  50 50.000 0.000 3.668 28.723 0.976 0.726 13.302 0.282 rootrot.csv
3      Total 187 119 68.453 2.989 4.558 68.972 0.991 0.720 14.371 0.271 rootrot.csv
```

It’s as simple as that. Now, let’s take a look at the same data set, except the input file is a GenAlEx file that has been formatted with Regional data (See section 1.4.2 for details). First, let’s see how the data set is laid out:

Figure 4: Part of the rootrot2.csv data set. Note the last two columns denoting the Regions and the number of individuals per region.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	56	187	10	19	18	18	25	27	17	23	21	10	9	2	97	90
2				1	2	3	4	5	6	7	8	9	10		Athena	Mt. Vernon
3	Ind	Pop	1	2	3	4	5	6	7	8	9	10	11	12	13	14
4	1	1	1	0	1	0	1	0	0	1	1	0	1	1	1	0
5	2	1	1	0	1	0	0	0	0	1	1	1	1	1	1	0
6	3	1	1	0	1	0	0	0	0	0	1	0	1	1	1	0
7	4	1	1	0	1	0	0	0	0	0	1	0	1	1	1	0
8	5	1	1	0	1	0	0	0	0	0	1	0	1	1	1	0

THE AMAZING DISAPPEARING OPTIONS!

Notice that I'm not writing in many of the options? This is because they have defaults. Since the data frame in my `@other` slot is called "population hierarchy", I don't have to specify that every time I do the function call, and that saves a lot of typing!

We'll import our data using `read.genalex` and take a look at the population hierarchy.

```
> Aeut2 <- read.genalex(system.file("files/rootrot2.csv", package="poppr"), region=TRUE)
> head(Aeut2@other$population_hierarchy)
```

```
Pop Region
1 1 Athena
2 1 Athena
3 1 Athena
4 1 Athena
5 1 Athena
6 1 Athena
```

```
> summary(Aeut2)
```

```
# Total number of genotypes: 187

# Population sample sizes:
 1  2  3  4  5  6  7  8  9 10
19 18 18 25 27 17 23 21 10  9

# Percentage of missing data:
[1] 0
```

What we see is that we have both of the population factors, but the names have changed and they are not combined. Note that since we specified "Athena" and "Mt. Vernon" as regions, the other level of the hierarchy was set as the population factor. We'll use `splitcombine` to combine both of these in the proper order. Note that we can use the indexes of the data frame columns to index these.

```
> Aeut2.combine <- splitcombine(Aeut2, method=2, hier=2:1)
> head(Aeut2.combine@other$population_hierarchy)
```

```
Pop Region Region_Pop
1 1 Athena Athena_1
2 1 Athena Athena_1
3 1 Athena Athena_1
4 1 Athena Athena_1
5 1 Athena Athena_1
6 1 Athena Athena_1
```

```
> summary(Aeut2.combine)
```

```
# Total number of genotypes: 187

# Population sample sizes:
  Athena_1  Athena_2  Athena_3  Athena_4  Athena_5  Athena_6
    9       12       10       13       10       5
  Athena_7  Athena_8  Athena_9  Athena_10 Mt. Vernon_1 Mt. Vernon_2
   11       8       10       9       10       6
Mt. Vernon_3 Mt. Vernon_4 Mt. Vernon_5 Mt. Vernon_6 Mt. Vernon_7 Mt. Vernon_8
    8       12       17       12       12       13

# Percentage of missing data:
[1] 0
```

Having these hierarchies in your data set is important when it comes to clone-censoring your data set (see section 2.4 *Attack of the Clone Correction*).

2.3 Divide (populations) and conquer (your analysis) {extract populations}

As I've mentioned before, *adegenet* has many ways of sub-setting the data, but you cannot easily subset a `genind` object by population in an efficient way. *Poppr* allows sub-setting a population from a `genind` object with one command.

2.3.1 Function: popsub

The command `popsub` is powerful in that it allows you to choose exactly what populations you choose to include or exclude from your analyses. As with many R functions, you can easily use this within a function to avoid creating a new variable to keep track of.

Default Command:

```
popsub(pop, sublist = "ALL", blacklist = NULL, mat = NULL)
```

- `pop` - a `genind` object.
- `sublist` - The vector of populations or integers representing the populations in your data set you wish to retain. For example: `sublist = c("pop_z", "pop_y")` or `sublist = 1:2`.
- `blacklist` - The vector of populations or integers representing the populations in your data set you wish to exclude. This can take the same type of arguments as `sublist`, and can be used in conjunction with `sublist` for when you want a range of populations, but know that there is one in there that you do not want to analyze. For example: `sublist = 1:15, blacklist = "pop_x"`. One very useful thing about the `blacklist` is that it allows the user to be extremely paranoid about the data. You can set the `blacklist` to contain populations that are not even in your data set and it will still work!
- `mat` - (see section 3, *Multilocus Genotype Analysis* for more information) This is where you would put a matrix that's produced by `mlg.table` to be subsetted instead of the `genind` object. If you do this, the matrix will return with only the rows equal to your populations and only the multilocus genotypes (columns) pertaining to those populations.

To demonstrate this tool, let's revisit the H3N2 data set. Let's say we wanted to analyze only the data in North America. To make sure we are all on the same page, we will reset the population factor to "country". Remember that this is located in a data frame in the `@other` slot called "x".

```
> data(H3N2)
> pop(H3N2) <- H3N2$other$x$country
> H3N2$pop.names # Only two countries from North America.
```

[1] "Japan"	"USA"	"Finland"	"China"	"South Korea"
[6] "Norway"	"Taiwan"	"France"	"Latvia"	"Netherlands"
[11] "Bulgaria"	"Turkey"	"United Kingdom"	"Denmark"	"Austria"
[16] "Canada"	"Italy"	"Russia"	"Bangladesh"	"Egypt"
[21] "Germany"	"Romania"	"Ukraine"	"Czech Republic"	"Greece"
[26] "Iceland"	"Ireland"	"Sweden"	"Nepal"	"Saudi Arabia"
[31] "Switzerland"	"Iran"	"Mongolia"	"Spain"	"Slovenia"
[36] "Croatia"	"Algeria"			

```
> H.na <- popsub(H3N2, sublist=c("USA", "Canada"))
> H.na$pop.names
```

P1	P2
"USA"	"Canada"

Since this is a larger data set, running the `summary` function might take a few seconds longer than we want it to. If we want to see the population size, we can use the *adegenet* function `nInd()`:

```
> nInd(H.na)
```

```
[1] 665
```

```
> nInd(H3N2)
```

```
[1] 1903
```

You can see that the population factors are correct and that the size of the data set is considerably smaller. Let's see the data set without the North American countries.

```
> H.minus.na <- popsub(H3N2, blacklist=c("USA", "Canada"))
> H.minus.na$pop.names
```

P01	P02	P03	P04	P05
"Japan"	"Finland"	"China"	"South Korea"	"Norway"
P06	P07	P08	P09	P10
"Taiwan"	"France"	"Latvia"	"Netherlands"	"Bulgaria"
P11	P12	P13	P14	P15
"Turkey"	"United Kingdom"	"Denmark"	"Austria"	"Italy"
P16	P17	P18	P19	P20
"Russia"	"Bangladesh"	"Egypt"	"Germany"	"Romania"
P21	P22	P23	P24	P25
"Ukraine"	"Czech Republic"	"Greece"	"Iceland"	"Ireland"
P26	P27	P28	P29	P30
"Sweden"	"Nepal"	"Saudi Arabia"	"Switzerland"	"Iran"
P31	P32	P33	P34	P35
"Mongolia"	"Spain"	"Slovenia"	"Croatia"	"Algeria"

Let's make sure that the number of individuals in both data sets added up equals the number of individuals in our original data set:

```
> (nInd(H.minus.na) + nInd(H.na)) == nInd(H3N2)
```

```
[1] TRUE
```

Now we have data sets with and without North America. Let's try something a bit more challenging. Let's say that we want The first 10 populations in alphabetical order, but we know that we still don't want any countries in North America. We can easily do this by using the *base* function `sort`.

```
> Hsort <- sort(H3N2$pop.names)[1:10]
> Hsort
```

[1] "Algeria"	"Austria"	"Bangladesh"	"Bulgaria"	"Canada"
[6] "China"	"Croatia"	"Czech Republic"	"Denmark"	"Egypt"

```
> H.alpha <- popsub(H3N2, sublist=Hsort, blacklist=c("USA", "Canada"))
> H.alpha$pop.names
```

P1	P2	P3	P4	P5
"China"	"Bulgaria"	"Denmark"	"Austria"	"Bangladesh"
P6	P7	P8	P9	
"Egypt"	"Czech Republic"	"Croatia"	"Algeria"	

And that, is how you subset your data with *poppr*!

2.4 Attack of the clone correction {clone-censor data sets}

Clone correction refers to the ability of keeping one observation per clone in a given population (or sub-population). Clone correcting can be hazardous if its done by hand (even on small data sets) and it requires a defined population hierarchy to get relevant results. *Poppr* has a clone correcting function that is able to correct at the lowest level of any defined population hierarchy. Note that clone correction in *poppr* is sensitive to missing data, as it treats all missing data as a single extra allele.

2.4.1 Function: clonecorrect

This function will return a clone corrected data set corrected for the lowest population level. Population levels are specified with the `hier` flag. You can choose to combine the population hierarchy to analyze at the lowest population level by choosing `combine = TRUE`.

Default Command:

```
clonecorrect(pop, hier = c(1), dfname = "population_hierarchy", combine = FALSE, keep = 1)
```

- `pop` - a `genind` object that has a population hierarchy data frame in the `@other` slot. Note, the `genind` object does not necessarily require a population factor to begin with.
- `hier` - This can be a vector of words or numbers referring to specific column names in your data frame in the `@other` slot.
- `dfname` - The name of a data frame you have in the `@other` slot with the population factors.
- `combine` - Do you want to combine the population hierarchy? If it's set to `FALSE` (default), you will be returned a `genind` object with the top most hierarchical level as a population factor.
- `keep` - This flag is to be used if you set `combine = FALSE`. This will tell clone correct to return a specific combination of your hierarchy. For example, imagine a hierarchy that needs to be clone corrected at three levels: *Population by Year by Month*. If you wanted to only run an analysis on the *Population* level, you would set `keep = 1` since *Population* is the first level of the hierarchy. On the other hand, if you wanted to run analysis on *Year by Month*, you would set `keep = 2:3` since those are the second and third levels of the hierarchy.

Let's look at ways to clone-correct our data. We'll look at our *A. euteichies* data since that data set is known to include clonal populations [5]. Notice that I am not including the options `dfname` and `combine` because the default arguments suit my needs.

```
> data(Aeut)
> A.cc <- clonecorrect(Aeut, hier=c("Pop", "Subpop"), keep=1)
> poppr(A.cc, quiet=TRUE)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	rbarD	File
1	Athena	76	70	60.621	1.017	4.221	65.636	0.998	0.963	2.535	0.062	rootrot.csv
2	Mt. Vernon	65	50	50.000	0.000	3.796	36.739	0.988	0.821	14.310	0.298	rootrot.csv
3	Total	141	119	59.629	1.854	4.705	96.980	0.997	0.876	13.802	0.260	rootrot.csv

Now let's compare the clone corrected analysis to the uncorrected data set:

```
> poppr(Aeut, quiet=TRUE)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	rbarD	File
1	Athena	97	70	65.981	1.246	4.063	42.193	0.986	0.721	2.906	0.072	rootrot.csv
2	Mt. Vernon	90	50	50.000	0.000	3.668	28.723	0.976	0.726	13.302	0.282	rootrot.csv
3	Total	187	119	68.453	2.989	4.558	68.972	0.991	0.720	14.371	0.271	rootrot.csv

As you can see from the summary tables, everything all sub-populations have been clone censored to the sub population level with respect to the population hierarchy. Notice how the observed number of individuals (N) decreases in the clone corrected data set.

2.5 Every day I'm shuffling (data sets) {permutations and bootstrap resampling}

A common null hypothesis for populations with mixed reproductive modes is panmixia, or to put it simply: lots of sex. A handy way to test for that is permutation analysis to assess random linkage among loci whereupon you randomly shuffle your data. *Poppr* uses randomly shuffled data sets in order to calculate P-values for the index of association (I_A and \bar{r}_d) [1]. Since there might be other tests where a permutation analysis would be pertinent, a shuffler for `genind` objects was created with four shuffling schemes: two schemes shuffling without replacement and two shuffling with replacement. Details below.

2.5.1 Function: shufflepop

Default Command:

```
shufflepop(pop, method = 1)
```

- **pop** - a **genind** object.
- **method** - a number indicating the method of sampling you wish to use. The following methods are available for use:

1. **Multilocus permutation** This is called Multilocus permutation because it does the same thing as the permutation analysis in the program *multilocus* by Paul Agapow and Austin Burt [1]. This will shuffle the genotypes at each locus. For example, a single diploid locus with four alleles (1, 2, 3, 4) with the frequencies of 0.1, 0.2, 0.3, and 0.4, respectively:

	[,1]	[,2]
[1,]	4	4
[2,]	4	1
[3,]	4	3
[4,]	2	2
[5,]	3	3

might become:

	[,1]	[,2]
[1,]	3	3
[2,]	4	1
[3,]	2	2
[4,]	4	4
[5,]	4	3

Note that you have the same genotypes after shuffling, so at each locus, you will maintain the same allelic frequencies and heterozygosity. So, in this sample, you will only see a homozygote with allele 2. This also ensures that the P-values associated with I_A and \bar{r}_d are exactly the same (for an explanation, see the end of section 4.1.1 of this manual). Unfortunately, if you are trying to simulate a sexual population, this does not make much biological sense as it assumes that alleles are not independently assorting within individuals.

2. **Permute Alleles** This is a sampling scheme that will permute alleles within the locus. So, using our example above, our resampling might become:

	[,1]	[,2]
[1,]	3	2
[2,]	2	4
[3,]	1	3
[4,]	4	3
[5,]	4	4

As you can see, The heterozygosity has changed, yet the allelic frequencies remain the same. Overall this would show you, for example, what would happen if the sample you had underwent panmixis within this sample itself.

3. **Parametric Bootstrap** The previous two schemes reshuffled the observed sample, but the parametric bootstrap uses the allelic frequencies as estimates of what the true allelic frequencies are and uses those as probabilities for each allele when resampling the data with replacement. Here are two samples to show you what I mean.

First Sample

	[,1]	[,2]
[1,]	1	3
[2,]	3	3
[3,]	3	2
[4,]	4	4
[5,]	4	2

Second Sample

	[,1]	[,2]
[1,]	3	4
[2,]	2	3
[3,]	4	2
[4,]	4	4
[5,]	4	2

Notice how the heterozygosity has changed along with the allelic frequencies. The frequencies for alleles 3 and 4 have switched in the first data set, and we've lost allele 1 in the second data set purely by chance! This type of sampling scheme attempts to show you what the true population would look like if it were panmictic and your original sample gave you a basis for estimating expected allele frequencies. Since estimates are made from the observed allele frequencies, small samples will produce skewed results.

4. **Non-Parametric Bootstrap** The final method is sampling with replacement, but with no assumption about the distribution of the alleles.

```
[1,] [,1] [,2]
[2,] 1 3
[3,] 3 3
[4,] 3 1
[5,] 2 2
[6,] 3 1
```

Again, heterozygosity and allele frequencies are not maintained, but now all of the alleles have a 1 in 4 chance of being chosen.

These shuffling schemes have been implemented for the index of association, but there may be other summary statistics you can use **shufflepop** for. All you have to do is use the function **replicate**. Let's use I_A as an example:

```
> data(nancycats)
> nan1 <- popsub(nancycats, 1)
> ia(nan1)
```

```
      Ia      rbarD
0.16564272 0.02105965
```

```
> replicate(10, ia(shufflepop(nan1, method = 3), quiet=TRUE))
```

```
      Ia      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
rbarD -0.29112163 -0.053843659 -0.17996769 0.14010641 0.30063966 0.21221103 -0.22694841
      Ia      [,8]      [,9]      [,10]
rbarD 0.40814607 -0.24815362 -0.40796307
      Ia      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
rbarD -0.03681503 -0.006945724 -0.02273945 0.01789136 0.03811873 0.02683328 -0.02906914
```

You could use this method to replicate the resampling 999 times and then create a histogram to visualize a distribution of what would happen under different assumptions of panmixia.

3 Multilocus Genotype Analysis

In populations with mixed sexual and clonal reproduction, it is not uncommon to have multiple samples from the same population have the same genotype across multiple loci (multilocus genotype, MLG). Here, we introduce tools for tracking MLGs within and across populations in **genind** objects from the *adegenet* package. We will be using SNP data from isolates of the H3N2 virus from 2002 to 2006.

3.1 Just a peek {How many multilocus genotypes are in our data set?}

First, let's take a quick look at how many Multilocus Genotypes are present within the H3N2 data set using the **mlg** function. This will tell us if any MLG analysis is needed.

3.1.1 Function: `mlg`

The function `mlg` allows for the counting of the number of MLGs in a `genind` object. This is a very simple command for quick reference to determine if your data set needs further multilocus genotype analysis.

Default Command:

```
mlg(pop, quiet = FALSE)
```

- `pop` - a `genind` object.
- `quiet` - if TRUE, the number of individuals and multilocus genotypes will be printed to the screen, if FALSE, nothing will be printed to the screen and the number of multilocus genotypes will be reported.

```
> data(H3N2)
> mlg(H3N2, quiet=FALSE)
```

```
#####
# Number of Individuals: 1903
# Number of MLG: 752
#####
[1] 752
```

We can see that since the number of individuals exceeds the number of multilocus genotypes, this data set contains clones. Let's take a look at where those clones are with respect to populations.

3.2 Clone-ing around {MLGs across populations}

Since you have the ability to change the population structure of your data set freely, it is quite possible to see some of the same MLGs across different populations. Tracking them by hand can be a nightmare with large data sets. Luckily, `mlg.crosspop` has you covered in that regard.

3.2.1 Function: `mlg.crosspop`

Analyze the MLGs that cross populations within your data set. This has three output modes. The default one gives a list of MLGs, and for each MLG, it gives a named numeric vector indicating the abundance of that MLG in each population. Alternate outputs are described with `indexreturn` and `df`.

Default Command:

```
mlg.crosspop(pop, sublist = "ALL", blacklist = NULL, mlgsub = NULL, indexreturn =
FALSE, df = FALSE, quiet = FALSE)
```

- `pop` - a `genind` object.
- `sublist` - see `mlg.table`, Section 3.3.1. Analyze specified populations.
- `blacklist` - see `mlg.table`, Section 3.3.1. Do not include specified populations.
- `mlgsub` - see `mlg.table`, Section 3.3.1. Only analyze specified MLGs. The vector for this flag can be produced by this function as you will see later in this vignette.
- `indexreturn` - return a vector of indices of MLGs. (You can use these in the `mlgsub` flag, or you can use them to subset the columns of an MLG table).
- `df` - return a data frame containing the MLGs, the populations they cross, and the number of copies you find in each population. This is useful for making graphs in *ggplot2*.
- `quiet` - TRUE or FALSE. Should the populations be printed to screen as they are processed? (will print nothing if `indexreturn` is TRUE)

We can see what Multilocus Genotypes cross different populations and then give a vector that shows how many populations each multi-population MLG crosses.

```
> pop(H3N2) <- H3N2$other$x$country
> H.dup <- mlg.crosspop(H3N2, quiet=TRUE)
```

Here is a snippet of what the output looks like when quiet is FALSE. It will print out the MLG name, the total number of individuals that make up that MLG, and the populations where that MLG can be found.

```
MLG.3: (12 inds) USA Denmark
MLG.9: (16 inds) Japan USA Finland Denmark
MLG.31: (9 inds) Japan Canada
MLG.75: (23 inds) Japan USA Finland Norway Denmark Austria Russia Ireland
MLG.80: (2 inds) USA Denmark
MLG.86: (7 inds) Denmark Austria
MLG.95: (2 inds) USA Bangladesh
MLG.97: (8 inds) USA Austria Bangladesh Romania
MLG.104: (3 inds) USA France
MLG.110: (16 inds) Japan USA China
```

The output of this function is a list of MLGs, each containing a vector indicating the number of copies in each population. We'll count the number of populations each MLG crosses using the function `sapply` with `length`.

```
> head(H.dup)
```

```
$MLG.3
  USA Denmark
4      8

$MLG.9
  Japan  USA Finland Denmark
1      13      1      1

$MLG.31
  Japan Canada
2      7

$MLG.75
  Japan  USA Finland Norway Denmark Austria Russia Ireland
2      8      2      1      6      2      1      1

$MLG.80
  USA Denmark
1      1

$MLG.86
  Denmark Austria
3      4
```

```
> H.num <- sapply(H.dup, length) # count the number of populations each MLG crosses.
> H.num
```

MLG.3	MLG.9	MLG.31	MLG.75	MLG.80	MLG.86	MLG.95	MLG.97	MLG.104	MLG.110	MLG.119
2	4	2	8	2	2	2	4	2	3	2
MLG.149	MLG.158	MLG.163	MLG.205	MLG.206	MLG.207	MLG.210	MLG.213	MLG.221	MLG.224	MLG.227
2	6	2	2	3	2	2	4	2	3	3
MLG.234	MLG.241	MLG.244	MLG.246	MLG.252	MLG.253	MLG.258	MLG.274	MLG.277	MLG.283	MLG.285
6	3	2	10	2	9	2	5	3	3	2
MLG.290	MLG.291	MLG.314	MLG.315	MLG.317	MLG.321	MLG.325	MLG.326	MLG.334	MLG.344	MLG.350
3	2	2	3	3	2	2	2	2	2	2
MLG.368	MLG.370	MLG.381	MLG.401	MLG.405	MLG.417	MLG.439	MLG.453	MLG.461	MLG.471	MLG.508
3	2	3	3	3	5	2	2	3	2	3
MLG.529	MLG.530	MLG.540	MLG.548	MLG.552	MLG.556	MLG.570	MLG.578	MLG.580	MLG.582	MLG.589
5	3	3	2	2	2	2	2	2	2	2
MLG.597	MLG.605	MLG.611	MLG.615	MLG.619	MLG.620	MLG.621				
2	3	2	2	2	4	2				

3.3 Bringing something to the table {producing MLG tables and graphs}

We can also create a table of multilocus genotypes per population as well as bar graphs to give us a visual representation of the data. This is achieved through the function `mlg.table`

3.3.1 Function: mlg.table

Produce a matrix containing counts of MLGs (columns) per population (rows). If there is no population structure to your data set, a vector will be produced instead.

Default Command:

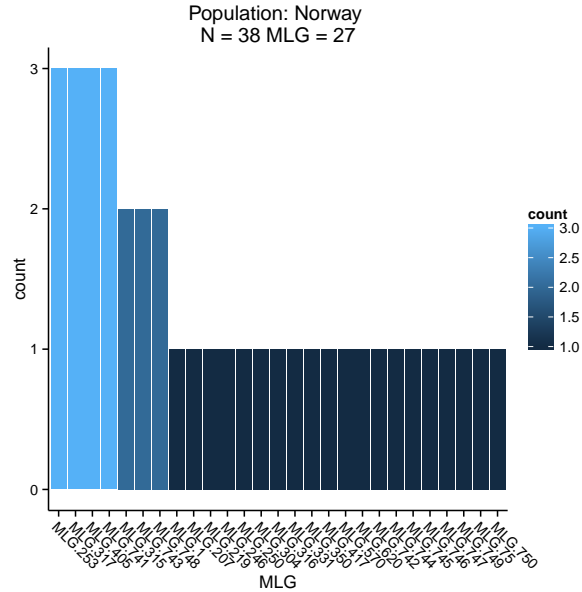
```
mlg.table(pop, sublist = "ALL", blacklist = NULL, mlgsub = NULL, bar = TRUE, total  
= FALSE, quiet = FALSE)
```

- **pop** - a **genind** object.
- **sublist** - a vector indicating which specific populations you want to produce a table for. This can be a numeric or character vector. See section 2.3.1 for details.
- **blacklist** - a vector indicating which specific populations you do not want to include in your table. This can be a numeric or character vector, and does not necessarily have to be the same type as **sublist**. eg. **sublist=1:10**, **blacklist="USA"**. See section 2.3.1 for details.
- **mlgsub** - a vector containing the indices of MLGs you wish to subset your table with.
- **bar** - TRUE or FALSE. If TRUE, a bar plot will be printed for each population with more than one individual.
- **total** - TRUE or FALSE. Should the entire data set be included in the table? This is equivalent to evoking **colSums** on the table.
- **quiet** - TRUE or FALSE. When **bar** is TRUE, should the populations be printed to screen as they are processed?

```
> H.tab <- mlg.table(H3N2, quiet=TRUE, bar=TRUE)  
> H.tab[1:10, 1:10] # Showing the first 10 columns and rows of the table.
```

	MLG.1	MLG.2	MLG.3	MLG.4	MLG.5	MLG.6	MLG.7	MLG.8	MLG.9	MLG.10
Japan	0	0	0	0	0	0	1	2	1	0
USA	0	2	4	1	1	0	0	0	13	0
Finland	0	0	0	0	0	0	0	0	1	0
China	0	0	0	0	0	0	0	0	0	0
South Korea	0	0	0	0	0	1	0	0	0	0
Norway	1	0	0	0	0	0	0	0	0	0
Taiwan	0	0	0	0	0	0	0	0	0	0
France	0	0	0	0	0	0	0	0	0	0
Latvia	0	0	0	0	0	0	0	0	0	0
Netherlands	0	0	0	0	0	0	0	0	0	0

Figure 5: An example of a bar-chart produced by `mlg.table`. Note that this data set would produce several such charts.



The MLG table is not restricted for use with just *Poppr*. One of the main advantages of the function `mlg.table` is that it allows easy access to diversity functions present in the package *vegan* [12]. One very simple example is to create a rarefaction curve for each population in your data set giving the number of expected MLGs for a given sample size. For more information, type `help("diversity", package="vegan")` in your R console.

For the sake of a simple example, instead of drawing a curve for each of the 37 countries represented in this sample, let's change the population structure to be the different years of the epidemics.

```
> H.year <- H3N2
> pop(H.year) <- H.year$other$x$year
> summary(H.year) # Check the data to make sure it's correct.
```

```
# Total number of genotypes: 1903
# Population sample sizes:
2002 2003 2004 2005 2006
158 415 399 469 462

# Number of alleles per locus:
L001 L002 L003 L004 L005 L006 L007 L008 L009 L010 L011 L012 L013 L014 L015 L016 L017 L018
3 3 4 2 4 2 3 2 4 3 4 2 4 3 2 2 3 3
L019 L020 L021 L022 L023 L024 L025 L026 L027 L028 L029 L030 L031 L032 L033 L034 L035 L036
2 2 3 3 3 2 2 2 2 2 2 2 2 2 2 4 4 3
L037 L038 L039 L040 L041 L042 L043 L044 L045 L046 L047 L048 L049 L050 L051 L052 L053 L054
3 3 4 2 2 2 4 3 2 3 4 2 3 2 3 2 2 2
L055 L056 L057 L058 L059 L060 L061 L062 L063 L064 L065 L066 L067 L068 L069 L070 L071 L072
4 2 2 2 2 2 2 2 4 4 4 3 3 2 3 4 3 2
L073 L074 L075 L076 L077 L078 L079 L080 L081 L082 L083 L084 L085 L086 L087 L088 L089 L090
3 3 3 3 2 3 2 4 2 3 2 2 3 3 3 3 2 2
L091 L092 L093 L094 L095 L096 L097 L098 L099 L100 L101 L102 L103 L104 L105 L106 L107 L108
2 2 3 2 3 2 3 2 3 2 3 2 2 2 3 2 2 2
L109 L110 L111 L112 L113 L114 L115 L116 L117 L118 L119 L120 L121 L122 L123 L124 L125
2 3 3 3 2 2 3 3 3 3 4 2 3 3 4 3 2

# Number of alleles per population:
1 2 3 4 5
203 255 232 262 240

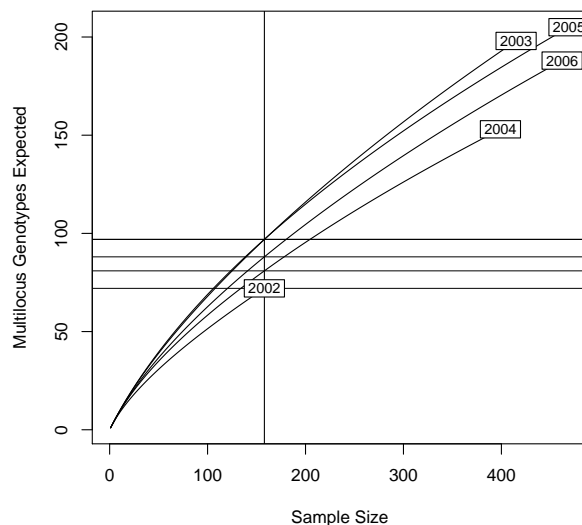
# Percentage of missing data:
[1] 2.363426

# Observed heterozygosity:
```

```
[1] 0
# Expected heterozygosity:
[1] 0

> H.year <- mlg.table(H.year, bar=FALSE)
> rarecurve(H.year, ylab="Multilocus genotypes expected", sample=min(rowSums(H.year)))
```

Figure 6: An example of a rarefaction curve produced using a MLG table.



The minimum value from the *base* function `rowSums()` of the table represents the minimum common sample size of all populations. Setting the “sample” flag draws the horizontal and vertical lines you see on the graph. The intersections of these lines correspond to the numbers you would find if you ran the function `poppr` on this data set (under the column “eMLG”).

3.4 Getting into the mix {combining MLG functions}

Alone, the different functionalities are neat. Combined, we can create interesting data sets. Let’s say we wanted to know which MLGs were duplicated across the regions of the United Kingdom, Germany, Netherlands, and Norway. All we have to do is use the `sublist` flag in the function:

```
> UGNN.list <- c("United Kingdom", "Germany", "Netherlands", "Norway")
> UGNN <- mlg.crosspop(H3N2, sublist=UGNN.list, indexreturn=TRUE)
```

OK, the output tells us that there are three MLGs that are crossing between these populations, but we do not know how many are in each. We can easily find that out if we subset our original table, `H.tab`.

```
> UGNN # Note that we have three numbers here. This will index the columns for us.
```

```
MLG.315 MLG.317 MLG.620
  315      317      620
```

```
> UGNN.list # And let's not forget that we have the population names.
```

```
[1] "United Kingdom" "Germany"          "Netherlands"    "Norway"
```

```
> H.tab[UGNN.list, UGNN]
```

	MLG.315	MLG.317	MLG.620
United Kingdom	1	0	0
Germany	0	1	1
Netherlands	0	0	0
Norway	2	3	1

Now we can see that Norway has a higher incidence of nearly all of these MLGs. We can go even further and subset the original data set to only give us those MLGs by utilizing the function `mlg.vector`:

3.4.1 Function: `mlg.vector`

This function is the backbone for `mlg.table` and `mlg.crosspop`, and is The function that determines what your MLGs are. This is quite useful for sub-setting the data set to only contain the MLGs of interest. The numbers in the vector correspond to the number of columns in a matrix produced by `mlg.table`. It is important to remember that this is also sensitive to missing data and will treat it as a single extra allele.

Default Command:

```
mlg.vector(pop)
```

- `pop` - a `genind` object.

```
> H.vec <- mlg.vector(H3N2)
> H.sub <- H3N2[H.vec %in% UGNN, ]
> mlg.table(H.sub, bar=FALSE)
```

	MLG.1	MLG.2	MLG.3
Austria	0	0	7
Germany	0	1	1
Greece	0	0	1
Norway	2	3	1
Japan	4	1	0
United Kingdom	1	0	0

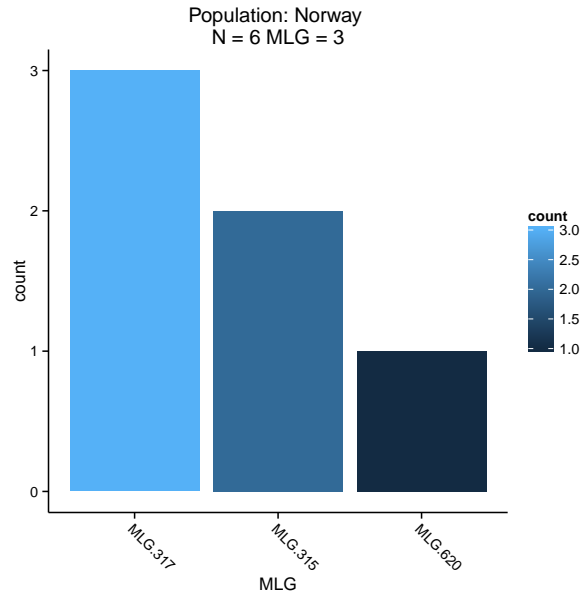
You can also do the same thing using the `mlgsub` flag.

```
> mlg.table(H3N2, mlgsub=UGNN, bar=TRUE)
```

	MLG.315	MLG.317	MLG.620
Japan	4	1	0
Norway	2	3	1
United Kingdom	1	0	0
Austria	0	0	7
Germany	0	1	1
Greece	0	0	1

And we can see where exactly these three MLGs fall within our data set.

Figure 7: An example of the same bar-chart as *Figure 1*, but focusing on three MLGs.



Now, you might notice that the MLG vector no longer matches up with our data after we subset it.

```
> H.vec[1:22]
```

```
[1] 605 605 672 675 674 673 670 671 670 678 678 678 678 582 615 580 581 570 615 582 582
[22] 592
```

```
> mlg.vector(H.sub)
```

```
[1] 3 3 3 3 3 3 3 3 3 3 2 1 1 2 2 1 2 1 1 1 1 2
```

Well, this is unfortunate because it means that we can't compare any subsetting data with non-subsetting data. Luckily, there's a little trick we can do using our old friend, the `@other` slot. If we place the MLG vector in the `@other` slot of our original data set, it will be subsetting along with the data.

```
> H3N2@other$MLG.vector <- H.vec
> H.sub <- H3N2[H.vec %in% UGNN, ]
> H.sub@other$MLG.vector
```

```
[1] 620 620 620 620 620 620 620 620 620 620 317 315 315 317 317 315 317 315 315 315 315
[22] 317
```

Magic!

So, we've gotten this far, yet we haven't actually seen what the genotypes look like! For analyses where the genotypic signature is important, this is a crucial identification step. Lucky for us, the `genind` object retains all of the genotypic information and can be accessed using the `genind2df` function. Let's take a look at the three genotypes we specified above utilizing the vector of MLGs we created above, `H.vec`.

```
> H.df <- genind2df(H3N2)
> H.df[H.vec %in% UGNN, 1:15] # Showing only 15 columns because it is a large dataset.
```

```
> H.df <- genind2df(H3N2[, loc=names(H3N2@loc.names)[1:15]])
> H.df[H.vec %in% UGNN, 1:15] # Showing only 15 columns because it is a large dataset.
```

Notice that there seems to be a clear separation between the SNPs of the first 10 isolates and the rest? This is no coincidence. Take a look at the output of our sub-setting.

```
> UGNN
```

```
MLG.315 MLG.317 MLG.620
      315      317      620
```

```
> H.vec[H.vec %in% UGNN]
```

```
[1] 620 620 620 620 620 620 620 620 620 620 317 315 315 317 317 315 317 315 315 315
[22] 317
```

We have the MLGs 315, 317, and 620, and the result of the sub-setting shows us that 620 occurs earlier in our data set, and that MLGs 315 and 317 are mixed in together. The reason why we do not see a mixture of three different sets of SNP calls in our little window is because `mlg.vector` creates the MLGs by first concatenating and then sorting the genotypes. This way, the closer two MLG indexes are to each other, the fewer differences they will have between one another.

3.5 Do you see what I see? {alternative data visualization}

The graphs that are output by *poppr* are simply aids for the user to make data analysis easier. We want to better visualize how these MLGs cross populations by MLG or population. We also want to see exactly what MLGs are in which populations, and how prevalent they are. As the package *ggplot2* is based on data frames, we have to give ourselves a data frame to work with. We can do this using the `df = TRUE` flag.

```
> df <- mlg.crosspop(H3N2, df=TRUE, quiet=TRUE)
> names(df)
```

```
[1] "MLG"          "Population" "Count"
```

Now that we have our data frame, we can do a couple of things. We can first see where the most omnipresent MLG occurs. After that, we will plot the top ten MLGs using *ggplot2*.

```
> H.max <- names(sort(H.num, decreasing=TRUE)[1:10])
> # Showing the data frame by the largest MLG complex.
> df[df$MLG %in% H.max[1], ]
```

	MLG	Population	Count
76	MLG.246	Japan	3
77	MLG.246	USA	8
78	MLG.246	China	4
79	MLG.246	Norway	1
80	MLG.246	Austria	6
81	MLG.246	Russia	1
82	MLG.246	Egypt	1
83	MLG.246	Iceland	1
84	MLG.246	Nepal	15
85	MLG.246	Switzerland	1

And now we can visualize the largest ten MLG complexes using *ggplot2*'s `qplot` function.

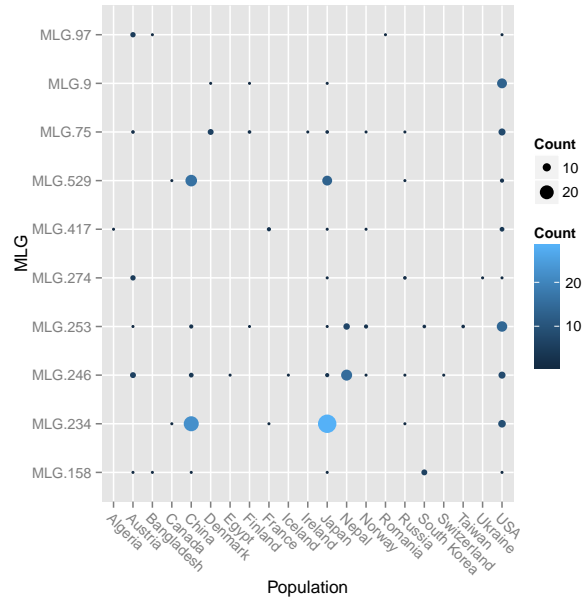
4 Index and Distance Calculations

4.1 The missing linkage disequilibrium {calculating the index of association, I_A and \bar{r}_d }

The index of association has previously been implemented in a couple of programs including *multilocus* [1] and *LIAN* [6], but these are programs that are no longer well supported, can be difficult to install and do not run in batch mode.

Figure 8: An example of the versatility of the MLG information.

```
> df2 <- df[df$MLG %in% H.max, ]
> library(ggplot2)
> qplot(y=MLG, x=Population, data=df2, color=Count, size=Count) +
+   theme(axis.text.x = element_text(size = 10, angle = -45, hjust = 0))
```



4.1.1 Function: ia

This function is a quick look at a single data set. It can do almost everything that `poppr` can do except for sorting through populations.

Default Command:

```
ia(pop, sample = 0, method = 1, quiet = "minimal", missing = "ignore", hist = TRUE)
```

- `pop` - a `genind` object.
- `sample` - You should use this flag whenever you want to reshuffle your data set. Indicate how many times you want to reshuffle your data set to obtain a P-value.
- `method` - a number from 1 to 4 indicating the sampling method:
 1. *multilocus* style permutation [1].
 2. permutation over alleles.
 3. parametric bootstrap.
 4. non-parametric bootstrap.

The methods are detailed in section 2.5.1 of this manual.

- `quiet` - This has three settings, `TRUE`, `FALSE`, and `"noisy"`. If set to `TRUE`, nothing will be printed to the screen as the sampling progresses. If `FALSE` and if there is sampling, a single dot for each sampling replicate will be printed to the screen to show the progress of the sampling. Choosing `"noisy"` is not recommended for the average user as it is meant for debugging. It will print the values of I_A and \bar{r}_d to the screen as they are produced.

- **missing** - This will preprocess your missing values. It is set to ignore missing data, so that they do not contribute to the distance measure. It can also be set to "loci", "geno", "zero", or "mean". For details, see section 2.1.1 of this manual.
- **hist** - This will produce a pair of histograms for each population showing the distribution of I_A and \bar{r}_d across the sampled data sets, and plot the observed value as a single vertical line.

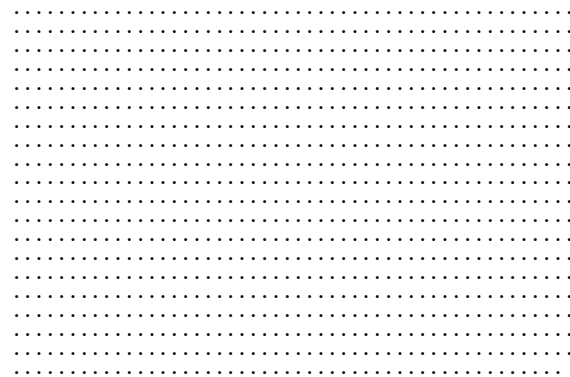
Running the analysis is as simple as this:

```
> ia(nancycats)

      Ia      rbarD
0.17207262 0.02178965
```

We can use `popsb` to subset for specific populations. Here, we'll also demonstrate the sampling flag and show you what the histogram looks like.

```
> set.seed(1001)
> ia(popsb(nancycats, 5), sample=999)
```



```
      Ia      p.Ia      rbarD      p.rD
-0.047539953 0.589000000 -0.006004254 0.589000000
```

This analysis produced the histograms you see below. What these histograms represent are 999 resamplings of the data under the null hypothesis (H_0) of sexual reproduction. The way that H_0 is created is determined by the sampling method chosen. In this case, the method was to shuffle genotypes at each locus to simulate unlinked loci. Since the $P = 0.589$, we would fail to reject H_0 and we therefore might conclude that this population is sexually reproducing [2] [17] [1].

Note that both of the P-values are exactly the same. This is what you would obtain with the default sampling scheme because the variances within loci do not change. When you change the sampling scheme, however, the P-values can end up being different. That's why there are two fields for P-values. Let's look at what happens to the P-values once we change the sampling scheme to a parametric bootstrap.

```
> set.seed(1001)
> ia(popsb(nancycats, 5), sample=999, method=3, quiet=TRUE, hist=FALSE)
```

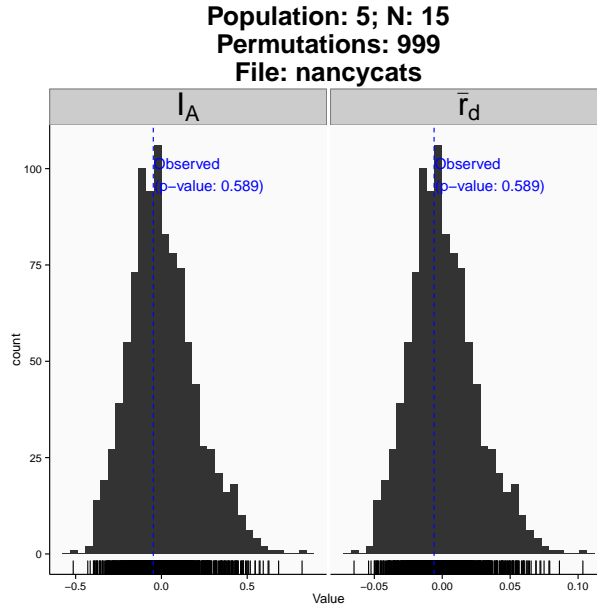
```
      Ia      p.Ia      rbarD      p.rD
-0.047539953 0.589000000 -0.006004254 0.596000000
```

There, is, of course one little caveat that needs to be mentioned. The P-values are calculated by comparing how many permuted values are greater than or equal to the observed value. This includes the observed value (which is why setting the randomizations to 999 will give you a round P-value) which means that the lowest P-value you will ever have is $1/(n+1)$ where n is the number of permutations you select. Take for example this population of a clonal root rot pathogen, *Aphanomyces euteiches*:

```
> data(Aeut)
> set.seed(1001)
> ia(popsb(Aeut, 1), sample=999, method=3, quiet=TRUE, hist=FALSE)
```

```
      Ia      p.Ia      rbarD      p.rD
2.90602921 0.00100000 0.07237008 0.00100000
```

Figure 9: Histograms of 999 values of I_A and \bar{r}_d calculated from 999 resamplings of population 5 from the data set “nancycats”. The observed values of I_A and \bar{r}_d are represented as vertical blue lines overlaid on the distributions. The ticks at the bottom of each histogram represent individual observations.



4.2 Going the distance {dissimilarity distance}

Since *poppr* is still in its infancy, the number of distance measures it can offer are few. Bruvo’s distance is well supported and allows you to quickly visualize your data, but it only allows for microsatellites. The index of association, above, utilizes a discrete dissimilarity distance matrix. It is with this matrix that we have constructed a relative dissimilarity distance where the distance is the ratio of the number of dissimilarities to the number of dissimilarities possible. The number of dissimilarities possible is the number of loci multiplied by the ploidy, so if you have 10 loci from a diploid population, then there are 20 dissimilarities possible. For details, see equations (2) and (3) in section 6.1.1.

4.2.1 Function: `diss.dist`

Use this function to calculate relative dissimilarity between individuals and return a distance matrix for use in creating cladograms or minimum spanning networks.

Default Command:

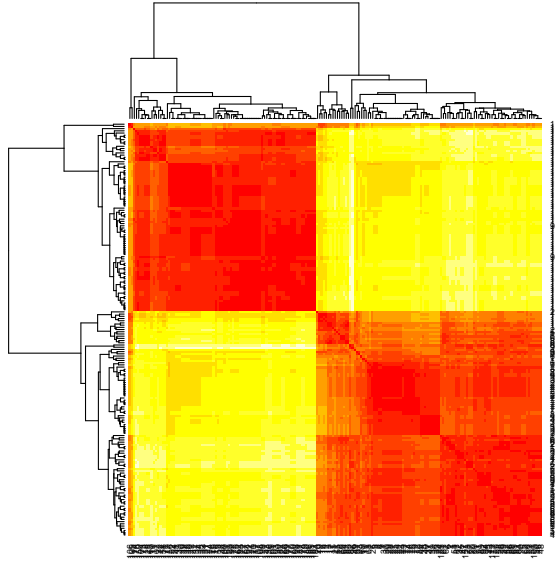
```
diss.dist(pop)
```

- `pop` – a `genind` object.

Since we have a data set that we know is very clonal, let’s analyze the *A. euteiches* data set [5] and create a heatmap to visualize the degree of difference between populations.

```
> data(Aeut)
> A.dist <- diss.dist(Aeut)
> heatmap(as.matrix(A.dist), symm=TRUE)
```


Figure 10: Heatmap representation of a dissimilarity distance for the data set “Aeut”



4.3 Step by stepwise mutation {Bruvo’s distance}

Bruvo’s distance is a genetic distance measure for microsatellite markers utilizing a stepwise mutation model that allows for differing ploidy levels [3]. As *adeigenet*’s *genind* object has an all or none approach to missing data, any genotypes not exhibiting full ploidy will be treated as missing. This means that only non-special cases will be considered for the calculation and missing data will be ignored [3]. It is important to note that this is a distance between individuals, not populations, unlike Nei’s 1978 distance [11]. For distances between populations, see the *adeigenet* function `dist.genpop`

4.3.1 Function: `bruvo.dist`

Bruvo’s distance requires knowledge of the repeat lengths of each locus, so take care to read the description below.

Default Command:

```
bruvo.dist(pop, replen = 1, add = TRUE, loss = TRUE)
```

- `pop` - a *genind* object.
- `replen` - This is a vector of numbers indicating the repeat length for each locus in your sample. If you have two dinucleotide repeats and five tetranucleotide repeats, you would put `c(2,2,4,4,4,4,4)` in this field. If you have imported data where that represents the raw number of steps, all you would have to type is `rep(1, n)`, replacing *n* with the number of loci in your sample. It is important that you place something in this field because this function will attempt to estimate the repeat length based on the minimum difference of the alleles represented; with variability of position calls, relying on this estimation is NOT recommended.
- `add` - For missing data: use the genome addition model (see section 6.1.2)
- `loss` - For missing data: use the genome loss model (see section 6.1.2)

This function will return a distance matrix (displaying the smallest population in the data set “nancycats”):

```
> dist9 <- bruvo.dist(popsb(nancycats, 9), replen=rep(1,9))
> dist9
```

```

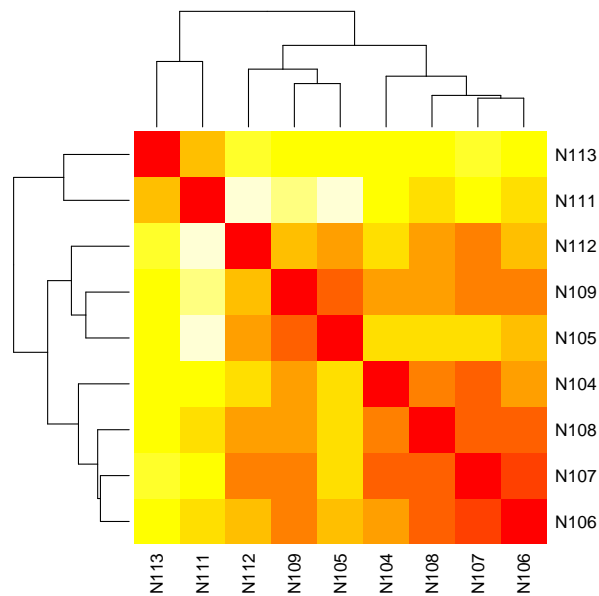
      N104      N105      N106      N107      N108      N109      N111      N112
N105 0.5778800
N106 0.4008213 0.4563124
N107 0.2202691 0.5093036 0.1805522
N108 0.3270365 0.5533719 0.2352431 0.2178786
N109 0.4016376 0.2760247 0.3192139 0.3330612 0.4294671
N111 0.6150004 0.8707648 0.5533278 0.6167331 0.5219014 0.7330560
N112 0.5492079 0.4086363 0.4391785 0.3289388 0.3886142 0.4528936 0.8037448
N113 0.5925835 0.6227587 0.6203071 0.6585558 0.6186252 0.6099514 0.5060450 0.7026198

```

You can visualize this better with a simple heatmap:

Figure 11: Heatmap representation of Bruvo's distance for population 9 of the data set "nancycats"

```
> heatmap(as.matrix(dist9), symm=TRUE)
```



Let's take a closer look at the two individuals, N113 and N111. They seem to have large distances between everyone else and themselves. The names and columns of the matrix contain the names of individuals, but not the population information. We can make a comparison of Bruvo's distance across populations easier by editing the "Labels" attribute of the distance object. Let's take a look at the labels attribute using the `attr()` command.

```
> attr(dist9, "Labels")
```

```

      1      2      3      4      5      6      7      8      9
"N104" "N105" "N106" "N107" "N108" "N109" "N111" "N112" "N113"

```

Remember that they all came from population 9, so let's append that to each label using the `paste()` command.

```

> dist9.attr <- attr(dist9, "Labels")
> attr(dist9, "Labels") <- paste(rep("P09", 9), dist9.attr)
> dist9

```

```

      P09 N104 P09 N105 P09 N106 P09 N107 P09 N108 P09 N109 P09 N111 P09 N112
P09 N105 0.5778800
P09 N106 0.4008213 0.4563124
P09 N107 0.2202691 0.5093036 0.1805522
P09 N108 0.3270365 0.5533719 0.2352431 0.2178786
P09 N109 0.4016376 0.2760247 0.3192139 0.3330612 0.4294671
P09 N111 0.6150004 0.8707648 0.5533278 0.6167331 0.5219014 0.7330560
P09 N112 0.5492079 0.4086363 0.4391785 0.3289388 0.3886142 0.4528936 0.8037448
P09 N113 0.5925835 0.6227587 0.6203071 0.6585558 0.6186252 0.6099514 0.5060450 0.7026198

```

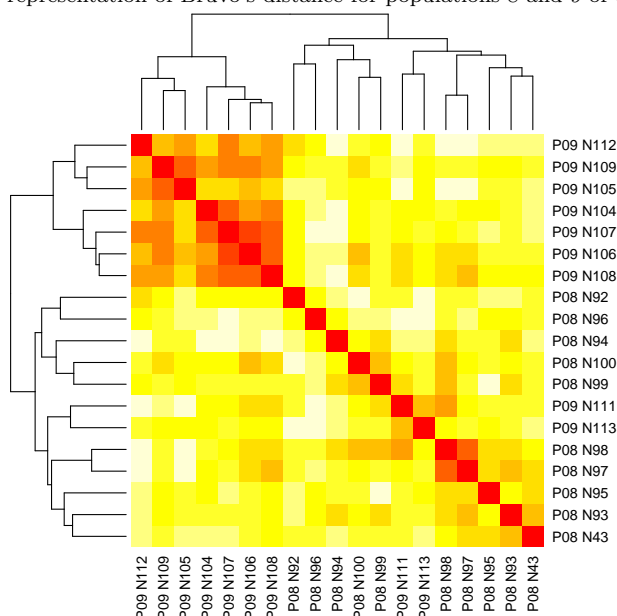
Now we can see that all of the labels are corresponding to population 9. Let's calculate Bruvo's distance between populations 8 and 9.

```

> dist9to8 <- bruvo.dist(popsb(nancycats, 8:9), replen=rep(1,9))
> dist9to8.attr <- attr(dist9to8, "Labels")
> nan9to8pop <- nancycats@pop[nancycats@pop %in% c("P08", "P09")]
> attr(dist9to8, "Labels") <- paste(nan9to8pop, dist9to8.attr)
> heatmap(as.matrix(dist9to8), symm=TRUE)

```

Figure 12: Heatmap representation of Bruvo's distance for populations 8 and 9 of the data set "nancycats"



Remember N113 and N111? Take a look at where they fall on the heatmap. They don't cluster together with population 9 anymore, but somewhere in population 8.

4.4 See the forest for the trees {visualizing distances with dendrograms and networks}

Staring at a raw distance matrix might be able to tell you something about your data, but it also might be able to ruin your eyesight. In this section, we present functions to display this data in trees and networks.

4.4.1 Function: bruvo.boot

This function provides the ability to draw a dendrogram based on Bruvo's distance including bootstrap support.

Default Command:

```

bruvo.boot(pop, replen = 1, add = TRUE, loss = TRUE, sample = 100, tree = "upgma",
showtree = TRUE, cutoff = NULL, quiet = FALSE)

```

- `pop` - a `genind` object.
- `replen` - see `bruvo.dist`, above.
- `add` - For missing data: use the genome addition model (see section 6.1.2)
- `loss` - For missing data: use the genome loss model (see section 6.1.2)
- `sample` - How many bootstraps do you want to perform?
- `tree` - Two trees are available, Neighbor-Joining "`nj`" or UPGMA "`upgma`".
- `showtree` - if `TRUE`, a tree will be plotted automatically.
- `cutoff` - This is a number between 0 and 100 indicating the cutoff value for the bootstrap node labels. If you only wanted to see the the bootstrap values for nodes that were present more than 75% of the time, you would use `cutoff = 75`. If you don't put anything for this parameter, all values will be shown.
- `quiet` - if `quiet = TRUE`, no standard messages will be printed to screen. If `quiet = FALSE` (default), then a progress bar and standard message will be printed to the screen.

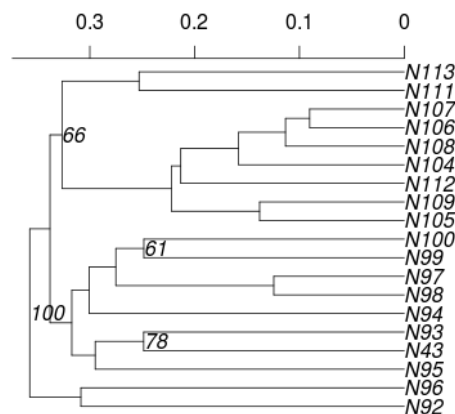
For this example, let's set the cutoff to 50%.

```
> set.seed(1001)
> nan9tree <- bruvo.boot(popsb(nancycats, 8:9), replen=rep(1,9), sample=1000, cutoff=50)
```

Bootstrapping... (note: calculation of node labels can take a while even after the progress bar is full)

```
|=====| 100%
```

Figure 13: UPGMA Tree of Bruvo's distance for population 9 of the data set "nancycats" with 1000 Bootstrap Replicates. Node labels represent percentage of bootstrap replicates that contained that node.



4.4.2 Function: `greycurve`

Use this function to display a gradient of grey values based on user-defined parameters. The following functions will display a minimum spanning network that utilize a grey scale to display the weight of the lines (referred to as “edges”) that connect two or more individuals. The darker the line the closer the distance. Since this is based off of a linear grey scale, what happens when you have a distance matrix comprised of values all below 0.2 or all above 0.8?

With linear grey scaling, it becomes very difficult to detect the differences in these ranges. The following function allows you to visualize and manipulate a gradient from black to white so that you can use it in *poppr*’s *msn* functions below to maximize the visual differences in your data.

Default Command:

```
greycurve(glim = c(0, 0.8), gadj = 3, gweight = 1)
```

This function does not return any values. It will print a visual gradient from black to white horizontally. On this gradient, it will plot the adjustment curve (in opposing grey values), yellow horizontal lines bounding the maximum and minimum values, and the equation used to calculate the correction in red. Keep in mind that this is plotting values from zero to one.

First, we'll see what happens when we change the weight parameter.

Figure 14: Default for `greywidth()`, weighted for small values.

```
> greycurve()
```

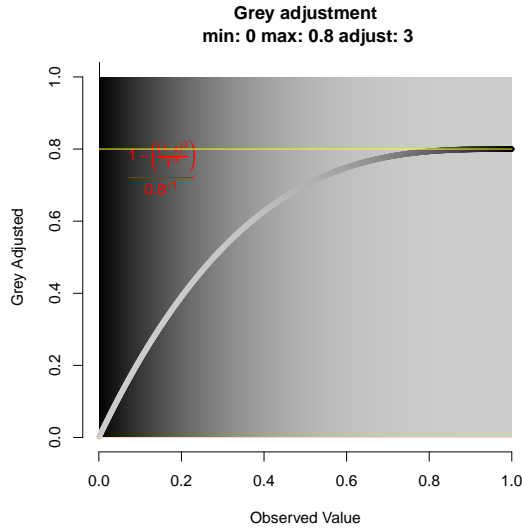
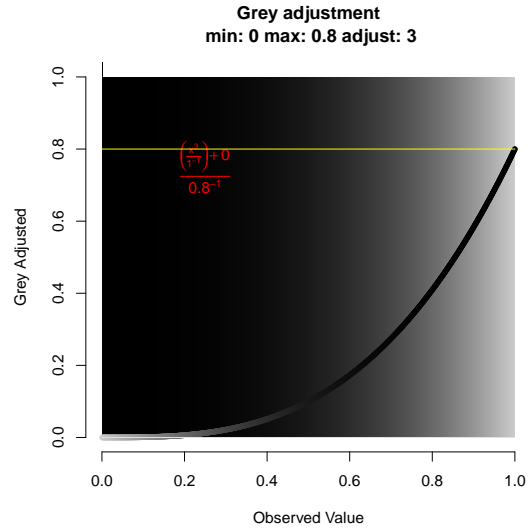


Figure 15: weighting for large values.

```
> greycurve(gweight = 2)
```



Now, we'll see what happens when we change the adjustment parameter (affects the shape of the curve) and the upper and lower limits of the grey scale.

Figure 16: Setting the lower and upper limits and weighting the curve heavily toward smaller values.

```
> greycurve(glim = c(0.2, 0.9), gadj=15)
```

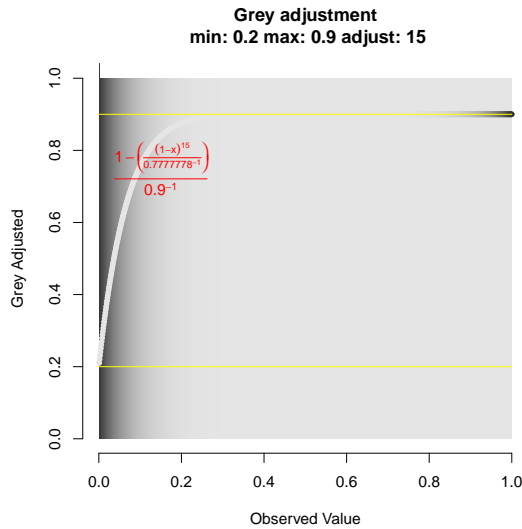
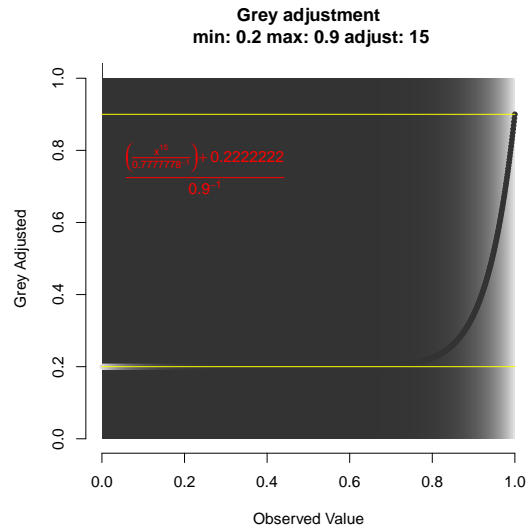


Figure 17: Same as the figure on the left, but weighting heavily toward larger values.

```
> greycurve(glim = c(0.2, 0.9), gadj=15, gweight=2)
```



4.4.3 Function: `bruvo.msn`

This function will automatically draw a minimum spanning network of MLGs based on Bruvo's distance. It's important to note that this will recalculate Bruvo's distance each time it is run, but the amount of time it takes to run is on the order of seconds. It will return a list containing the network, the populations and the related colors in the network so you can export or redraw it with the legend if you wanted to.

Default Command:

```
bruvo.msn(pop, replen = 1, add = TRUE, loss = TRUE, palette = topo.colors,  
          sublist = "All", blacklist = NULL, vertex.label = "MLG", gscale = TRUE,  
          glim = c(0, 0.8), gadj = 3, gweight = 1, wscale = TRUE, ...)
```

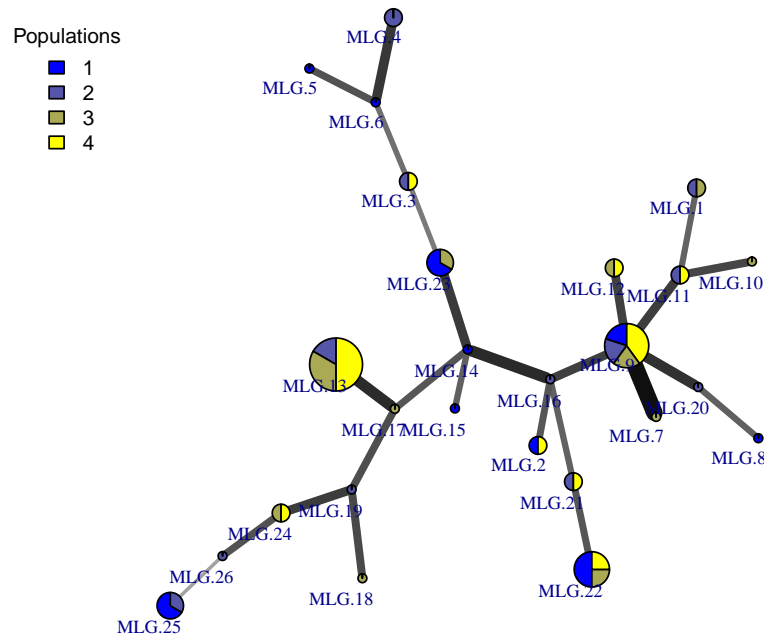
- `pop` - a `genind` object.
- `replen` - see `bruvo.dist`, above.
- `add` - For missing data: use the genome addition model (see section 6.1.2)
- `loss` - For missing data: use the genome loss model (see section 6.1.2)
- `palette` - this is a **function** defining a color palette to use. The default is `topo.colors`. There are different palettes, which you can search by typing `?rainbow`. If you want a custom color palette, an easy way is to use the function `colorRampPalette`.
- `sublist` - The populations you wish to analyze. This defaults to "All". See section 2.3.1 for details.
- `blacklist` - Populations you do not want to include in the graph. See section 2.3.1 for details.
- `vertex.label` - This is an option that is passed on to *igraph*'s `plot` function. *Poppr* has added two arguments specific to *poppr*. If you want to label the graph with the multilocus genotypes from the whole data set, use the argument `vertex.label = "mlg"`. If you want to display the representative individual names, you can use the argument `vertex.label = "inds"`. I say representative individual names because, only one representative from each MLG will be present in the clone corrected data set used to calculate the distance. For no labels, you can choose `vertex.label = NA`.
- `gscale` - If this is set to `TRUE`, the edge color will be converted to greyscale based on Bruvo's distance. If two nodes are closely related, the edge will appear darker. The limits of the scale can be set by the argument `glim`. If this is set to `FALSE`, all edge colors will be black.
- `glim` - This is a vector of numbers between 0 and 1. This lets you set the limits of the grey scaling based on R's internal `grey` function. For example, if you wanted a maximum of 50% white saturation (for use if you have distantly related nodes) and a minimum of 1%, you would use `glim = c(0.01, 0.5)`.
- `gadj` - This is an integer greater than zero used to adjust the scaling factor for the grey curve. Since very small changes in the grey scale are not easily perceived, it's useful to be able to adjust the grey scale to be able to show you the weights of each edge. For example, a population with most weights less than 0.3, you might want to set `gadj = 10` to exaggerate the grey scale.
- `gweight` - If `gweight = 1`, the grey scale adjustment will be weighted towards separating out smaller values of Bruvo's distance. If `gweight = 2`, the grey scale adjustment will be weighted towards separating out larger values of Bruvo's distance.
- `wscale` - If this is set to `TRUE`, edge widths will be displayed corresponding to Bruvo's distance in that thicker edges will represent a smaller distance between nodes. If this is set to `FALSE`, all edges will be set to a width of 2.

- ... - This is a placeholder for any other arguments that you want to supply to *igraph*. Useful arguments are `vertex.label.cex` to adjust the size of the labels, `vertex.label.dist` to adjust the position of the labels, and `vertex.label.color` to adjust the color of the labels.

Often, minimum spanning networks are the preferred way to visualize Bruvo's distance. *Poppr* offers an easy way to plot these. For a demonstration, let's analyze a simulated data set of 50 individuals from populations that reproduce at a 99.9% rate of clonal reproduction.

```
> data(partial_clone)
> set.seed(9005)
> pc.msn <- bruvo.msn(partial_clone, replen=rep(1, 10), vertex.label.cex=0.7,
+   vertex.label.dist=-0.5, palette=colorRampPalette(c("blue", "yellow")))
```

Figure 18: Minimum Spanning Network representing 4 simulated populations. Each node represents a different multi locus genotype (MLG). Node sizes and colors correspond to the number of individuals and population membership, respectively. Edge thickness and color are proportional to Bruvo's distance. Edge lengths are arbitrary.



Note that the thickness of the edges (the lines that are connecting the dots) is representative of relatedness between individuals, but the lengths do not necessarily mean anything due to the fact that with a larger data sets, displaying lengths proportional to relatedness would be impossible to draw on a 2D surface. Interpreting these data would show that MLG 9 has 5 individuals from all four populations and that it is most closely related to MLG 7, whereas the most distantly related connection exists between MLG 25 and MLG 26.

4.4.4 Function: poppr.msn

Use this function to draw a minimum spanning network from your data set and a distance matrix derived from your data set. Since there are hundreds of distances that can be calculated for genetic data, and since I want to be able to graduate at some point in this decade, functions to automatically calculate distances and draw the minimum spanning networks will be few and far between. This function is an attempt to meet the user halfway and draw a minimum spanning network provided that the user has supplied two things:

1. A distance matrix over all individuals.
2. The original data set containing demographic information.

That's it. For the most part, this function is functionally the same as `bruvo.msn`, except that instead of being exclusive to microsatellite markers, you can now visualize distances in any marker type provided that you have the two items listed above.

Default Command:

```
poppr.msn(pop, distmat, palette = topo.colors, sublist = "All",
  blacklist = NULL, vertex.label = "MLG", gscale = TRUE, glim = c(0, 0.8),
  gadj = 3, gweight = 1, wscale = TRUE, ...)
```

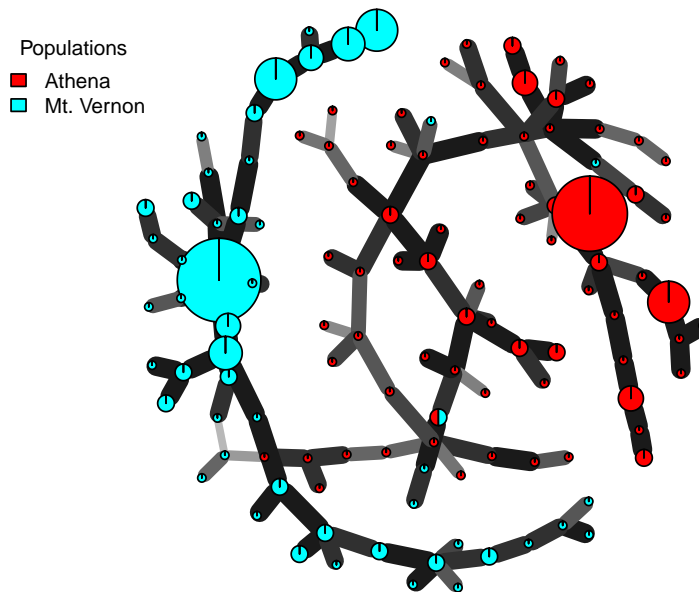
- `pop` - a `genind` object.
- `distmat` - a dissimilarity distance matrix derived from your data with distances between zero and one.
- `palette` - this is a **function** defining a color palette to use. The default is `topo.colors`. There are different palettes, which you can search by typing `?rainbow`. If you want a custom color palette, an easy way is to use the function `colorRampPalette`.
- `sublist` - The populations you wish to analyze. This defaults to "All".
- `blacklist` - Populations you do not want to include in the graph.
- `vertex.label` - This is an option that is passed on to *igraph*'s `plot` function. *Poppr* has added two arguments specific to *poppr*. If you want to label the graph with the multilocus genotypes from the whole data set, use the argument `vertex.label = "mlg"`. If you want to display the representative individual names, you can use the argument `vertex.label = "inds"`. I say representative individual names because, only one representative from each MLG will be present in the clone corrected data set used to calculate the distance. For no labels, you can choose `vertex.label = NA`.
- `gscale` - If this is set to `TRUE`, the edge color will be converted to greyscale based on the distance. If two nodes are closely related, the edge will appear darker. The limits of the scale can be set by the argument `glim`. If this is set to `FALSE`, all edge colors will be black.
- `glim` - This is a vector of numbers between 0 and 1. This lets you set the limits of the grey scaling based on R's internal `grey` function. For example, if you wanted a maximum of 50% white saturation (for use if you have distantly related nodes) and a minimum of 1%, you would use `glim = c(0.01, 0.5)`.
- `gadj` - This is an integer greater than zero used to adjust the scaling factor for the grey curve. Since very small changes in the grey scale are not easily perceived, it's useful to be able to adjust the grey scale to be able to show you the weights of each edge. For example, a population with most weights less than 0.3, you might want to set `gadj = 10` to exaggerate the grey scale.
- `gweight` - If `gweight = 1`, the grey scale adjustment will be weighted towards separating out smaller values of the distance. If `gweight = 2`, the grey scale adjustment will be weighted towards separating out larger values of Bruvo's distance.

- **wscale** - If this is set to **TRUE**, edge widths will be displayed corresponding to Bruvo's distance in that thicker edges will represent a smaller distance between nodes. If this is set to **FALSE**, all edges will be set to a width of 2.
- ... - This is a placeholder for any other arguments that you want to supply to *igraph*. Useful arguments are **vertex.label.cex** to adjust the size of the labels, **vertex.label.dist** to adjust the position of the labels, and **vertex.label.color** to adjust the color of the labels.

Since we have the ability, let's visualize the *A. euteiches* data set [5].

```
> data(Aeut)
> A.dist <- diss.dist(Aeut)
> set.seed(9005)
> A.msn <- poppr.msn(Aeut, A.dist, vertex.label=NA, palette=rainbow, gadj=15)
```

Figure 19: Minimum Spanning Network representing 4 simulated populations. Each node represents a different multi locus genotype (MLG). Node sizes and colors correspond to the number of individuals and population membership, respectively. Edge thickness and color are proportional to Bruvo's distance. Edge lengths are arbitrary.



5 I know what you did last summary table {diversity table}

Remember the summary function that you used to get all the diversity statistics in section 1.3? In this section, we will flesh out all that you can do with this function. This was the very first function that was written for *poppr* to make it easy for the user to manipulate and summarize the data in one function.

5.1 Function: *poppr*

This function is quite daunting with all its possibilities. You have the option to subset your data for specific populations, correct for missing data, and clone correct. With each of these possibilities, comes the need to provide all the arguments for their various functions.

Default Command:

```
poppr(pop, total = TRUE, sublist = c("ALL"), blacklist = c(NULL), sample = 0,
      method = 1, missing = "ignore", cutoff = 0.05, quiet = "minimal",
      clonecorrect = FALSE, hier = c(1), dfname = "population_hierarchy",
      hist = TRUE, minsamp = 10)
```

- **pop** - A *genind* object.
- **total** - This is also a synonym for “pooled”. This will calculate all diversity statistics on the entire data set if set to **TRUE** or if there is no population structure.
- *pops* **sub** *functions*: See section 2.3
 - sublist** - A list of populations you want to include in your analysis.
 - blacklist** - A list of populations you want to exclude from your analysis.
- *shuffle* *pop* *functions*: See section 2.5
 - Note that this only affects the calculation for I_A and \bar{r}_d .
 - sample** - The number of samples you desire (eg. 999)
 - method** - Which sampling method? 1: multilocus, 2: permute, 3: parametric bootstrap, 4: non-parametric bootstrap.
- *missing* *no* *functions*: See Section 2.1
 - Note that all analyses in this function ignore/impute missing data by default.
 - missing** - How to deal with missing data. This feeds into the **type** flag of *missingno*.
 - cutoff** - Allowable percentage of missing data per genotype or locus.
- **quiet** - This has three settings, **TRUE**, **FALSE**, and **"noisy"**. If set to **TRUE**, nothing will be printed to the screen as the sampling progresses. If **FALSE** and if there is sampling, a single dot for each sampling replicate will be printed to the screen to show the progress of the sampling. Choosing **"noisy"** is not recommended for the average user as it is meant for debugging. It will print the values of I_A and \bar{r}_d to the screen as they are produced.
- *clone* *correct* *functions*: See section 2.4
 - clonecorrect** - if this is set to **TRUE**, then you will need to set the next two parameters.
 - hier** - A list of the population hierarchy, or names of columns in the data frame noted below.
 - dfname** - A data frame in the **@other** slot of the *genind* object containing all of the population factors in different columns. For an example, see sections 2.2 and 2.4.

keep - A vector of integers as indexes for the **hier** flag indicating which levels of the hierarchy you want to analyze. See section 2.4 for details.

- **hist** - if TRUE, a histogram of distributions of I_A and \bar{r}_d will be displayed with each population if there is sampling.
- **minsamp** - The minimum number of individuals you want to use to calculate the expected number of MLGs. The default is set to 10.

This function produces a table that contains the population name, number of individuals observed, number of MLGs observed, number of MLGs expected at the lowest common sampling size within the data set [8] [7], the Shannon-Wiener index [16], Stoddart and Taylor's index for expected MLGs [18], Nei's 1987 genotypic diversity [11], evenness [15][10][4], the index of association [2][17], the standardized index of association [1], and the file name. Most of these indices are calculated by converting the population into an MLG table with **mlg.table** (see section 3.3) and using the *vegan* package's **diversity** function (To see details, type `?vegan::diversity` into the R console).

To begin, let's revisit our example data set of *Aphanomyces euteiches* [5].

```
> data(Aeut)
> poppr(Aeut)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	rbarD	File
1	Athena	97	70	65.981	1.246	4.063	42.193	0.986	0.721	2.906	0.072	rootrot.csv
2	Mt. Vernon	90	50	50.000	0.000	3.668	28.723	0.976	0.726	13.302	0.282	rootrot.csv
3	Total	187	119	68.453	2.989	4.558	68.972	0.991	0.720	14.371	0.271	rootrot.csv

OK, so we were able to get a table out of this. Now let's see what happens when we do some sampling to see if this is reproducing clonally or not (for simplicity's sake, we will stick with the default multilocus sampling method). We will turn quiet on and the histogram off to save space.

```
> poppr(Aeut, sample=999, hist=FALSE, quiet=TRUE)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	p.Ia	rbarD	p.rD
1	Athena	97	70	65.981	1.246	4.063	42.193	0.986	0.721	2.906	0.001	0.072	0.001
2	Mt. Vernon	90	50	50.000	0.000	3.668	28.723	0.976	0.726	13.302	0.001	0.282	0.001
3	Total	187	119	68.453	2.989	4.558	68.972	0.991	0.720	14.371	0.001	0.271	0.001

File
1 rootrot.csv
2 rootrot.csv
3 rootrot.csv

From now on, we'll set **quiet** = TRUE to save space on our vignette. Let's clone correct at different levels to see if that affects the index of association. First, we'll clone correct at the sub population level.

```
> poppr(Aeut, sample=999, clonecorrect=TRUE, hier=c("Pop","Subpop"),
+       dfname="population_hierarchy", quiet=TRUE, hist=FALSE)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	p.Ia	rbarD	p.rD
1	Athena	76	70	60.621	1.017	4.221	65.636	0.998	0.963	2.535	0.001	0.062	0.001
2	Mt. Vernon	65	50	50.000	0.000	3.796	36.739	0.988	0.821	14.310	0.001	0.298	0.001
3	Total	141	119	59.629	1.854	4.705	96.980	0.997	0.876	13.802	0.001	0.260	0.001

File
1 rootrot.csv
2 rootrot.csv
3 rootrot.csv

And at the population level.

```
> poppr(Aeut, sample=999, clonecorrect=TRUE, hier="Pop",
+       dfname="population_hierarchy", quiet=TRUE, hist=FALSE)
```

	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	p.Ia	rbarD	p.rD
1	Athena	70	70	50.000	0.000	4.248	70.000	1	1.000	2.438	0.001	0.060	0.001
2	Mt. Vernon	50	50	50.000	0.000	3.912	50.000	1	1.000	13.856	0.001	0.285	0.001
3	Total	120	119	49.828	0.377	4.776	118.033	1	0.995	12.497	0.001	0.234	0.001

File
1 rootrot.csv
2 rootrot.csv
3 rootrot.csv

As you can see, clone correction doesn't always have to involve creation of new data sets!

You might notice that the P-values for both I_A and \bar{r}_d are often equal to each other. This is due to the default sampling method [1]. Here, we show examples where they are not equal.

```
> set.seed(2001)
> poppr(nancycats, sublist=5:6, total=FALSE, sample=999, method=3, quiet=TRUE, hist=FALSE)
```

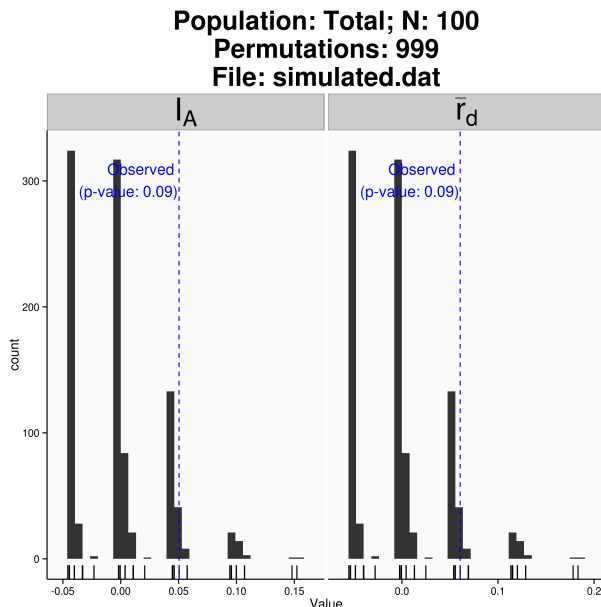
	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	p.Ia	rbarD	p.rD	File
1	5	15	15	11	0	2.708	15	1	1	-0.048	0.599	-0.006	0.599	truenames(nancycats)\$tab
2	6	11	11	11	0	2.398	11	1	1	0.334	0.064	0.043	0.065	truenames(nancycats)\$tab

The reason why the P-values would be different is described at the end of section 4.1.1. The differences in P-values are normally not very far off. It's important to note this because of what can happen in extremely clonal populations. You can end up with a large enough sample size consisting of very few MLGs. Upon shuffling, you find that there are very few values of I_A and \bar{r}_d that can be obtained. Observe with this simulated data set:

```
> set.seed(2004)
> poppr(system.file("files/simulated.dat", package="poppr"), sample=999, method=1, quiet=TRUE)
```

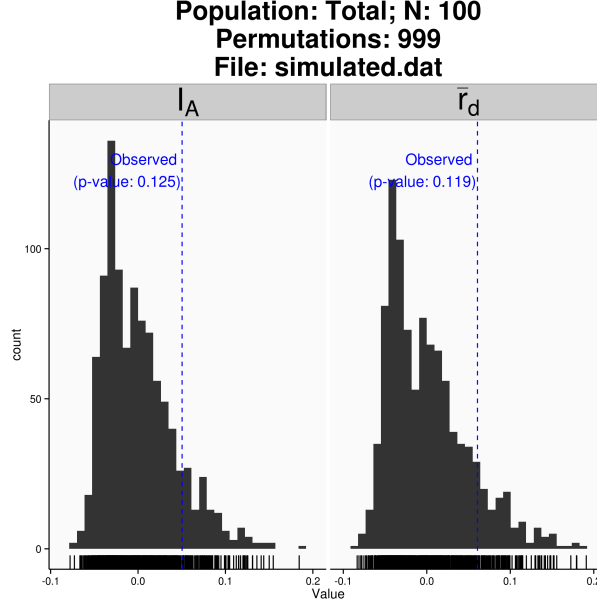
	Pop	N	MLG	eMLG	SE	H	G	Hexp	E.5	Ia	p.Ia	rbarD	p.rD	File
1	Total	100	6	6	0	1.235	2.79	0.648	0.735	0.05	0.09	0.061	0.09	simulated.dat

Figure 20: Output of multilocus-style sampling. Note the multi-modal distribution.



Take a look at these two histograms. The number of ways you can recombine the data with the default sampling method is very small. Other sampling methods could give a more theoretical distribution. Let's try the parametric bootstrap (For details, see section 2.5).

Figure 21: Output for parametric bootstrap sampling.



As you can see, the distribution is much closer to a distribution we would expect if this were a small sample of a larger population.

6 Appendix

6.1 Algorithmic Details

6.1.1 I_A and r_d

The index of association was originally developed by A.H.D. Brown analyzing population structure of wheat [2]. It has been widely used as a tool to detect clonal reproduction within populations [17]. Populations whose members are undergoing sexual reproduction, whether it be selfing or out-crossing, will produce gametes via meiosis, and thus have a chance to shuffle alleles in the next generation. Populations whose members are undergoing clonal reproduction, however, generally do so via mitosis. This means that the most likely mechanism for a change in genotype is via mutation. The rate of mutation varies from species to species, but it is rarely sufficiently high to approximate a random shuffling of alleles. The index of association is a calculation based on the ratio of the variance of the raw number of differences between individuals and the sum of those variances over each locus [17]. You can also think of it as the observed variance over the expected variance. If they are the same, then the index is zero after subtracting one (from Maynard-Smith, 1993 [17]):

$$I_A = \frac{V_O}{V_E} - 1 \quad (1)$$

Since the distance is more or less a binary distance, any sort of marker can be used for this analysis. In the calculation, phase is not considered, and any difference increases the distance between two individuals. Consider the genotypes of the dummy data frame we created earlier:

```

      locus1 locus2 locus3
1 101/101 201/201 301/302
2 102/103 202/203 301/303
3 102/102 203/204 304/305

```

Now, consider the first locus represented in the `genind` object:

```

      L1.1 L1.2 L1.3
1      1  0.0  0.0
2      0  0.5  0.5
3      0  1.0  0.0

```

Remember that each column represents a different allele and that each entry in the table represents the fraction of the genotype made up by that allele at that locus. Notice also that the sum of the rows all equal one. *Poppr* uses this to calculate distances by simply taking the sum of the absolute values of the differences between rows.

The calculation for the distance between two individuals at a single locus with a allelic states and a ploidy of k is as follows⁴:

$$d = \frac{k}{2} \sum_{i=1}^a |ind_{Ai} - ind_{Bi}| \quad (2)$$

```
> abs(dfg@tab[1, 1:3] - dfg@tab[2, 1:3])
```

```

L1.1 L1.2 L1.3
1.0  0.5  0.5

```

```
> abs(dfg@tab[1, 1:3] - dfg@tab[3, 1:3])
```

```

L1.1 L1.2 L1.3
1      1      0

```

```
> abs(dfg@tab[2, 1:3] - dfg@tab[3, 1:3])
```

```

L1.1 L1.2 L1.3
0.0  0.5  0.5

```

As you can see, these values of d at locus one add up to 2, 2, and 1, respectively.

To find the total number of differences between two individuals over all loci, you just take d over m loci, a value we'll call D :

$$D = \sum_{i=1}^m d_i \quad (3)$$

These values are calculated over all possible combinations of individuals in the data set, $\binom{n}{2}$ after which you end up with $\binom{n}{2} \cdot m$ values of d and $\binom{n}{2}$ values of D . Calculating the observed variances is fairly straightforward (modified from Agapow and Burt, 2001) [1]:

$$V_O = \frac{\sum_{i=1}^{\binom{n}{2}} D_i^2 - \frac{(\sum_{i=1}^{\binom{n}{2}} D_i)^2}{\binom{n}{2}}}{\binom{n}{2}} \quad (4)$$

⁴Individuals with Presence / Absence data will have the $k/2$ term dropped.

Calculating the expected variance is the sum of each of the variances of the individual loci. The calculation at a single locus, j is the same as the previous equation, substituting values of D for d [1]:

$$var_j = \frac{\sum_{i=1}^{\binom{n}{2}} d_i^2 - \frac{(\sum_{i=1}^{\binom{n}{2}} d_i)^2}{\binom{n}{2}}}{\binom{n}{2}} \quad (5)$$

The expected variance is then the sum of all the variances over all m loci [1]:

$$V_E = \sum_{j=1}^m var_j \quad (6)$$

Now you can plug the sums of equations (4) and (6) into equation (1) to get the index of association. Of course, Agapow and Burt showed that this index increases steadily with the number of loci, so they came up with an approximation that is widely used, \bar{r}_d [1]. For the derivation, see the manual for *multilocus*. The equation is as follows, utilizing equations (4), (5), and (6) [1]:

$$\bar{r}_d = \frac{V_O - V_E}{2 \sum_{j=1}^m \sum_{\substack{k=1 \\ j < k}}^m \sqrt{var_j \cdot var_k}} \quad (7)$$

6.1.2 Bruvo's distance

Bruvo's distance between two individuals calculates the minimum distance across all combinations of possible pairs of alleles at a single locus and then averaging that distance across all loci [3]. The distance between each pair of alleles is calculated as [3]:

$$m_x = 2^{-|x|} \quad (8)$$

$$d_a = 1 - m_x \quad (9)$$

Where x is the number of steps between each allele. So, let's say we were comparing two haploid ($k = 1$) individuals with alleles 228 and 244 at a locus that had a tetranucleotide repeat pattern (CATG) ^{n} . The number of steps for each of these alleles would be $228/4 = 57$ and $244/4 = 61$, respectively. The number of steps between them is then $|57 - 61| = 4$. Bruvo's distance at this locus between these two individuals is then $1 - 2^{-4} = 0.9375$. For samples with higher ploidy (k), there would be k such distances of which we would need to take the sum [3].

$$s_i = \sum_{a=1}^k d_a \quad (10)$$

Unfortunately, it's not as simple as that since we do not assume to know phase. Because of this, we need to take all possible combinations of alleles into account. This means that we will have k^2 values of d_a , when we only want k . How do we know which k distances we want? We will have to invoke parsimony for this and attempt to take the minimum sum of the alleles, of which there are $k!$ possibilities [3]:

$$d_l = \frac{\left(\min_{i \dots k!} s_i \right)}{k} \quad (11)$$

Finally, after all of this, we can get the average distance over all loci [3].

$$D = \frac{\sum_{i=1}^l d_i}{l} \quad (12)$$

This is calculated over all possible combinations of individuals and results in a lower triangle distance matrix over all individuals.

6.1.3 Special Cases of Bruvo's distance

As shown in the above section, ploidy is irrelevant with respect to calculation of Bruvo's distance. However, since it makes a comparison between all alleles at a locus, it only makes sense that the two loci need to have the same ploidy level. Unfortunately for polyploids, it's often difficult to fully separate distinct alleles at each locus, so you end up with genotypes that appear to have a lower ploidy level than the organism [3].

To help deal with these situations, Bruvo has suggested three methods for dealing with these differences in ploidy levels [3]:

- Infinite Model - The simplest way to deal with it is to count all missing alleles as infinitely large so that the distance between it and anything else is 1. Aside from this being computationally simple, it will tend to inflate distances between individuals.
- Genome Addition Model - If it is suspected that the organism has gone through a recent genome expansion, the missing alleles will be replaced with all possible combinations of the observed alleles in the shorter genotype. For example, if there is a genotype of [69, 70, 0, 0] where 0 is a missing allele, the possible combinations are: [69, 70, 69, 69], [69, 70, 69, 70], and [69, 70, 70, 70]. The resulting distances are then averaged over the number of comparisons.
- Genome Loss Model - This is similar to the genome addition model, except that it assumes that there was a recent genome reduction event and uses the observed values in the full genotype to fill the missing values in the short genotype. As with the Genome Addition Model, the resulting distances are averaged over the number of comparisons.
- Combination Model - Combine and average the genome addition and loss models.

As mentioned above, the infinite model is biased, but it is not nearly as computationally intensive as either of the other models. The reason for this is that both of the addition and loss models requires replacement of alleles and recalculation of Bruvo's distance. The number of replacements required is equal to the multiset coefficient: $\binom{n}{k} = \binom{n-k+1}{k}$ where n is the number of potential replacements and k is the number of alleles to be replaced. So, for the example given above, The genome addition model would require $\binom{2}{2} = 3$ calculations of Bruvo's distance, whereas the genome loss model would require $\binom{4}{2} = 10$ calculations.

To reduce the number of calculations and assumptions otherwise, Bruvo's distance will be calculated using the largest observed ploidy. This means that when comparing [69,70,71,0] and [59,60,0,0], they will be treated as triploids.

6.2 Exporting Graphics

R has the ability to produce nice graphics from most any type of data, but to get these graphics into a report, presentation, or manuscript can be a bit challenging. It's no secret that the R Documentation pages are a little difficult to interpret, so I will give the reader here a short example on how to export graphics from R. Note that any code here that will produce images will also be present in other places in this vignette. The default installation of the R GUI is quite minimal, and for an easy way to manage your plots and code, I strongly encourage the user to use Rstudio <http://www.rstudio.com/>.

6.2.1 Basics

Before you export graphics, you have to ask yourself what they will be used for. If you want to use the graphic for a website, you might want to opt for a low-resolution image so that it can load quickly. With printing, you'll want to make sure that you have a scalable or at least a very high resolution image. Here, I will give some general guidelines for graphics (note that these are merely suggestions, not defined rules).

- **What you see is not always what you get** I have often seen presentations where the colors were too light or posters with painfully pixellated graphs. Think about what you are going to be using a graphic for and how it will appear to the intended audience given the media type.
- **≥ 300 dpi unless its for a web page** For any sort of printed material that requires a raster based image, 300dpi (dots per inch) is the absolute minimum resolution you should use. For simple black and white line images, 1200dpi is better. This will leave you with crisp, professional looking images.
- **If possible, save to SVG, then rasterize** Raster images (bmp, png, jpg, etc...) are based off of the number of pixels or dots per inch it takes to render the image. This means that the raster image is more or less a very fine mosaic. Vector images (SVG) are built upon several interconnected polygons, arcs, and lines that scale relative to one another to create your graphic. With vector graphics, you can produce a plot and scale it to the size of a building if you wanted to. When you save to an SVG file first, you can also manipulate it in programs such as Adobe Illustrator or Inkscape.
- **Before saving, make sure the units and dimensions are correct** Unless you really wanted to save a graph that's over 6 feet wide.

6.2.2 Image Editors

Often times, fine details such as labels on networks need to be tweaked by hand. Luckily, there are a wide variety of programs that can help you do that. Here is a short list of image editors (both free and for a price) that you can use to edit your graphics.

- Bitmap based editors (for jpeg, bmp, png, etc...)

THE GIMP Free, cross-platform. <http://www.gimp.org>

PAINT.NET Free, Windows only. <http://www.getpaint.net>

ADOBE PHOTOSHOP Pricey, Windows and Mac. <http://www.adobe.com/products/photoshop.html>

- Scalable Vector Graphics based editors (for svg, pdf)

INKSCAPE Free, cross-platform <http://inkscape.org>

ADOBE ILLUSTRATOR Pricey, Windows and Mac. <http://www.adobe.com/products/illustrator.html>

6.2.3 Exporting ggplot2 graphics

ggplot2 is a fantastic package that *poppr* uses to produce graphs for the `mlg.table`, `poppr`, and `ia` functions. Saving a plot with *ggplot2* is performed with one command after your plot has rendered:

```
> data(nancycats) # Load the data set.  
> poppr(nancycats, sublist=5, sample=999) # Produce a single plot.  
> ggsave("nancy5.pdf")
```

Note that you can name the file anything, and `ggsave` will save it in that format for you. The details are in the documentation and you can access it by typing `help("ggsave")` in your R console. The important things to note are that you can set a `width`, `height`, and `unit`. The only downside to this function is that you can only save one plot at a time. If you want to be able to save multiple plots, read on to the next section.

6.2.4 Exporting any graphics

Some of the functions that *poppr* offers will give you multiple plots, and if you want to save them all, using *ggsave* will require a lot of tedious typing and clicking. Luckily, R has Functions that will save any plot you generate in nearly any image format you want. You can save in raster images such as png, bpm, and jpeg. You can also save in vector based images such as svg, pdf, and postscript. The important thing to remember is that when you are saving in a raster format, the default units of measurement are “pixels”, but you can change that by specifying your unit of choice and a resolution.

For raster images and svg files, you can only save your plots in multiple files, but pdf and postscript plots can be saved in one file as multiple pages. All of these functions have the same basic form. You call the function to specify the file type you want (eg. `pdf("myfile.pdf")`), create any graphs that you want to create, and then make sure to close the session with the function `dev.off()`. Let’s give an example saving to pdf and png files.

```
> data(H3N2)
> pop(H3N2) <- H3N2$other$х$country
> ####
> png("H3N2_barchart%02d.png", width = 14, height = 14, units = "in", res = 300)
> H.tab <- mlg.table(H3N2)
> dev.off()
> ####
```

Since this data set is made up of 30 populations with more than 1 individual, this will save 30 files to your working directory named “H3N2_barchart01.png...H3N2_barchart30.png”. The way R knows how to number these files is because of the %02d part of the command. That’s telling R to use a number that is two digits long in place of that expression. All of these files will be 14x14” and will have a resolution of 300 dots per inch. If you wanted to do the same thing, but place them all in one file, you should use the pdf option.

```
> pdf("H3N2_barcharts.png", width = 14, height = 14, compress = FALSE)
> H.tab <- mlg.table(H3N2)
> dev.off()
```

Remember, it is important not to forget to type `dev.off()` when you are done making graphs. Note that I did not have to specify a resolution for this image since it is based off of vector graphics.

6.3 Function calls

Here is a list of all the default function calls for *poppr*. Details can be found in the above sections.

- `getfile(multi = FALSE, pattern = NULL, combine = TRUE)` (Section 1.4.1)
- `read.genalex(genalex, ploidy = 2, geo = FALSE, region = FALSE)` (Section 1.4.2)
- `genind2genalex(pop, filename = "genalex.csv", quiet = FALSE, geo = FALSE, geodf = "xy")` (Section 1.4.5)
- `missingno(pop, type = "loci", cutoff = 0.05, quiet = FALSE)` (Section 2.1.1)
- `splitcombine(pop, method = 1, dfname = "population_hierarchy", sep = "_", hier = c(1), setpopulation = TRUE, fixed = TRUE)` (Section 2.2.1)
- `popsub(pop, sublist = "ALL", blacklist = NULL, mat = NULL)` (Section 2.3.1)
- `clonecorrect(pop, hier = c(1), dfname = "population_hierarchy", combine = FALSE, keep = 1)` (Section 2.4.1)
- `shufflepop(pop, method = 1)` (Section 2.5.1)
- `mlg(pop, quiet = FALSE)` (Section 3.1.1)
- `mlg.crosspop(pop, sublist = "ALL", blacklist = NULL, mlgsub = NULL, indexreturn = FALSE, df = FALSE, quiet = FALSE)` (Section 3.2.1)
- `mlg.table(pop, sublist = "ALL", blacklist = NULL, mlgsub = NULL, bar = TRUE, total = FALSE, quiet = FALSE)` (Section 3.3.1)

- `mlg.vector(pop)` (Section 3.4.1)
- `ia(pop, sample = 0, method = 1, quiet = "minimal", missing = "ignore", hist = TRUE)` (Section 4.1.1)
- `diss.dist(pop)` (Section 4.2.1)
- `bruvo.dist(pop, replen = 1, add = TRUE, loss = TRUE)` (Section 4.3.1)
- `bruvo.boot(pop, replen = 1, add = TRUE, loss = TRUE, B = 100, tree = "upgma", showtree = TRUE, cutoff = NULL, quiet = FALSE)` (Section 4.4.1)
- `greycurve(glim = c(0, 0.8), gadj = 3, gweight = 1)` (Section 4.4.2)
- `bruvo.msn(pop, replen = 1, add = TRUE, loss = TRUE, palette = topo.colors, sublist = "All", blacklist = NULL, vertex.label = "MLG", gscale = TRUE, glim = c(0, 0.8), gadj = 3, gweight = 1, wscale = TRUE, ...)` (Section 4.4.3)
- `poppr.msn(pop, distmat, palette = topo.colors, sublist = "All", blacklist = NULL, vertex.label = "MLG", gscale = TRUE, glim = c(0, 0.8), gadj = 3, gweight = 1, wscale = TRUE, ...)` (Section 4.4.4)
- `poppr(pop, total = TRUE, sublist = c("ALL"), blacklist = c(NULL), sample = 0, method = 1, missing = "ignore", cutoff = 0.05, quiet = "minimal", clonecorrect = FALSE, hier = c(1), keep = 1, dfname = "population_hierarchy", hist = TRUE, minsamp = 10)` (Section 5.1)
- `poppr.all(filelist, ...)` (Sections 1.4.1 and 5.1)

References

- [1] Paul-Michael Agapow and Austin Burt. Indices of multilocus linkage disequilibrium. *Molecular Ecology Notes*, 1(1-2):101-102, 2001.
- [2] A.H.D. Brown, M.W. Feldman, and E. Nevo. Multilocus structure of natural populations of *hordeum spontaneum*. *Genetics*, 96(2):523-536, 1980.
- [3] Ruzica Bruvo, Nicolaas K. Michiels, Thomas G. D'Souza, and Hinrich Schulenburg. A simple method for the calculation of microsatellite genotype distances irrespective of ploidy level. *Molecular Ecology*, 13(7):2101-2106, 2004.
- [4] Niklaus J. Grünwald, Stephen B. Goodwin, Michael G. Milgroom, and William E. Fry. Analysis of genotypic diversity data for populations of microorganisms. *Phytopathology*, 93(6):738-46, 2003.
- [5] N.J. Grünwald and G. Hoheisel. Hierarchical analysis of diversity, selfing, and genetic differentiation in populations of the oomycete *aphanomyces euteiches*. *Phytopathology*, 96(10):1134-1141, 2006.
- [6] Bernhard Haubold and Richard R. Hudson. Lian 3.0: detecting linkage disequilibrium in multilocus data. *Bioinformatics*, 16(9):847-849, 2000.
- [7] Kenneth L.Jr. Heck, Gerald van Belle, and Daniel Simberloff. Explicit calculation of the rarefaction diversity measurement and the determination of sufficient sample size. *Ecology*, 56(6):pp. 1459-1461, 1975.
- [8] S H Hurlbert. The nonconcept of species diversity: a critique and alternative parameters. *Ecology*, 52(4):577-586, 1971.
- [9] Thibaut Jombart. adegenet: a r package for the multivariate analysis of genetic markers. *Bioinformatics*, 24(11):1403-1405, 2008.
- [10] J.A. Ludwig and J.F. Reynolds. *Statistical Ecology. A Primer on Methods and Computing*. New York USA: John Wiley and Sons, 1988.

- [11] Masatoshi Nei. Estimation of average heterozygosity and genetic distance from a small number of individuals. *Genetics*, 89(3):583-590, 1978.
- [12] Jari Oksanen, F. Guillaume Blanchet, Roeland Kindt, Pierre Legendre, Peter R. Minchin, R. B. O'Hara, Gavin L. Simpson, Peter Solymos, M. Henry H. Stevens, and Helene Wagner. *vegan: Community Ecology Package*, 2012. R package version 2.0-5.
- [13] R. Peakall and P. E. Smouse. GENALEX 6: genetic analysis in excel. population genetic software for teaching and research. *Molecular Ecology Notes*, 6(1):288-295+, 2006.
- [14] Rod Peakall and Peter E. Smouse. Genalex 6.5: genetic analysis in excel. population genetic software for teaching and research--an update. *Bioinformatics*, 28(19):2537-2539, 2012.
- [15] E.C. Pielou. *Ecological Diversity*. Wiley, 1975.
- [16] Claude Elwood Shannon. A mathematical theory of communication. *Bell Systems Technical Journal*, 27:379-423,623-656, 1948.
- [17] J M Smith, N H Smith, M O'Rourke, and B G Spratt. How clonal are bacteria? *Proceedings of the National Academy of Sciences*, 90(10):4384-4388, 1993.
- [18] J.A. Stoddart and J.F. Taylor. Genotypic diversity: estimation and prediction in samples. *Genetics*, 118(4):705-11, 1988.