

Data import and manipulation in poppr version 1.1.4

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

February 3, 2015

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	Resources	3
1.3	Getting Help	3
1.4	Acknowledgements	4
1.5	Citation	4
1.6	Installation	4
1.6.1	From CRAN	4
1.6.2	From Source	4
1.6.3	From github	5
1.7	Quick start	5
1.8	Importing data into poppr {Get out of my dreams and into my R}	8
1.8.1	Function: getfile	9
1.8.2	Function: read.genalex	12
1.8.3	Other ways of importing data	15
1.8.4	Function: genind2genalex	15
1.9	Getting to know <i>adegenet</i> 's genind object	17
1.9.1	The other slot	17
1.10	The genclone object {send in the clones}	19
1.10.1	Function: as.genclone	19
1.11	Accessing the population hierarchy	22
1.12	About polyploid data	22
2	Data Manipulation	24
2.1	Population hierarchy construction {Can you take me hier(archy)?}	24
2.1.1	Defining hierarchies	24
2.1.2	Viewing hierarchies	27
2.1.3	Manipulating hierarchical levels	28
2.1.4	Defining populations with hierarchies	29
2.2	Replace or remove missing data {Inside the golden days of missing data}	30
2.2.1	Function: missingno	30
2.3	Extract populations {Divide (populations) and conquer (your analysis)}	34
2.3.1	Function: popsub	34
2.4	Clone-censor data sets {Attack of the clone correction}	35
2.4.1	Function: clonecorrect	36
2.5	Permutations and bootstrap resampling {every day I'm shuffling (data sets)}	37
2.5.1	Function: shufflepop	38
2.6	Removing uninformative loci {Cut It Out!}	40
2.6.1	Function: informloci	40
3	Multilocus Genotype Analysis	41
3.1	How many multilocus genotypes are in our data set? {Just a peek}	41
3.2	MLGs across populations {clone-ing around}	42
3.2.1	Function: mlg.crosspop	42
3.3	Producing MLG tables and graphs {bringing something to the table}	44
3.3.1	Function: mlg.table	44
3.4	Combining MLG functions {getting into the mix}	47

4	Appendix	50
4.1	General hierarchy method use	50
4.2	Manipulating Graphics	52
4.3	Exporting Graphics	54
4.3.1	Basics	54
4.3.2	Image Editors	54
4.3.3	Exporting ggplot2 graphics	55
4.3.4	Exporting any graphics	55
4.4	Table of Functions	55

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 microsattelites, 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, Genetix, Fstat, and Structure) and newly introduces compatibility with GenA1Ex 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 Resources

This vignette will cover all of the material you need to know to efficiently analyze data in *poppr*. For information on methods of analysis (eg. index of association, distance measures, AMOVA, ...), please read the manual pages provided for each function.

As *poppr* expanded from version 1.0, the vignette also expanded to be 80+ pages. As a result, it became clear that over 22,000 was less of a manual and more of a novella with a terrible plot. To remedy this, this vignette will focus only on data manipulation and a separate vignette, "algo", has been written to give algorithmic details of analyses introduced with *poppr*.

As of spring 2014, Drs. Niklaus J. Grünwald, Sydney E. Everhart, and I have co-authored a primer on using R for population genetic analysis. It is located [here](#) and the source code can be found [on our github site](#).

1.3 Getting Help

If you have any questions or feedback, feel free to send a message to the *poppr* forum at <http://groups.google.com/group/poppr>. You can submit bug reports there or on our github site: <https://github.com/grunwaldlab/poppr>

1.4 Acknowledgements

Much thanks goes to Sydney E. Everhart for alpha testing, beta testing, feature requests, proofreading, data contribution, and moral support throught the writing of this package and manual. Thanks also to Brian Knaus, Ignazio Carbone, David Cooke, Corine Schoebel, Jane Stewart, and Zaid Abdo for beta testing and feedback.

The following data sets are included in *poppr*: *Pinf* [4], *monpop* (Sydney E. Everhart, unpublished), *Aeut* [5]

1.5 Citation

The formal publication for *poppr* was published in the journal PeerJ: <http://peerj.com/articles/281/>. To cite *poppr*, please type in your R console:

```
citation(package = "poppr")
```

1.6 Installation

This manual assumes you have installed R. If you have not, please refer to The CRAN home page at <http://cran.r-project.org/>. We also recommend the Rstudio IDE (<http://www.rstudio.com/>), which allows the user to view the R console, environment, scripts, and plots in a single window.

1.6.1 From CRAN

To install *poppr* from CRAN, select “Package Installer” from the menu “Packages & Data” in the gui or type:

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

All dependencies (*adegetnet*, *pegas*, *vegan*, *ggplot2*, *phangorn*, *ape*, *reshape2* and *igraph*) will also be installed. In the unfortunate case this does not work, consult <http://cran.r-project.org/doc/manuals/R-admin.html#Installing-packages>.

1.6.2 From Source

The tarball for *poppr* can be downloaded from CRAN: <http://cran.r-project.org/package=poppr>, under the RESOURCES tab in the Grünwald Lab website: <http://http://grunwaldlab.cgrb.oregonstate.edu/>, or via github at <https://github.com/grunwaldlab/poppr>.

Since *poppr* contains C code, it needs to be compiled, which means that you need a working C compiler. If you are on Linux, you should have that, but if you are on Windows or OSX, you might need to download some special tools:

Windows Download Rtools: <http://cran.r-project.org/bin/windows/Rtools/>

OSX Download Xcode: <https://developer.apple.com/xcode/>

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("adegetnet", "pegas", "vegan", "ggplot2", "phangorn", "ape", "igraph"))
```

If you want to install from github, skip to the next section.

After installing dependencies, download the package to your computer and then install it with:

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

1.6.3 From github

Github is a repository where you can find all stable and development versions of *poppr*. Installing from github requires a C compiler, so be sure to read the section above for instructions on how to obtain that if you aren't on a Linux system.

To install from github, you do not need to download the tarball since there is a package called *devtools* that will download and install the package for you directly from github. After you have installed all dependencies (see above section), you should download *devtools*:

```
install.packages("devtools")
```

Now you can execute the command `install_github` with the user and repository name:

```
library("devtools")
install_github(repo = "grunwaldlab/poppr")
```

If you are the adventurous type and are willing to test out unreleased versions of the package, you can also install the development version:

```
library("devtools")
install_github(repo = "grunwaldlab/poppr", ref = "devel")
```

Users who install this version do so at their own risk. Since it is a development version, documentation may be incomplete or nonexistent for new functions.

1.7 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. 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:

Tell your computer to search R's library to find the *poppr* and load the package:

```
library("poppr")
```

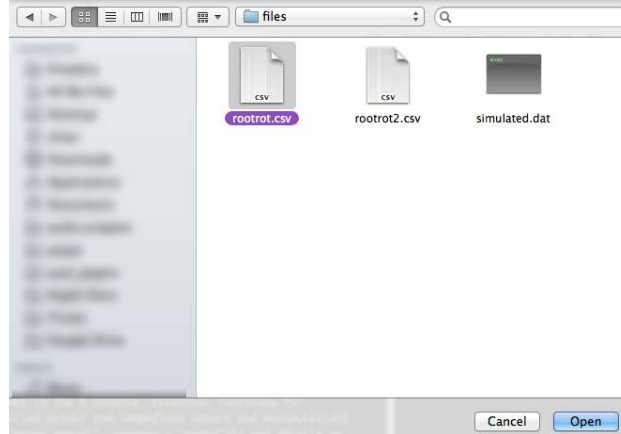
After that, you can use `getfile()`

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

A pop up window will appear like this¹:

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

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.

We can navigate throughout your entire computer through this window and tell R where to go. The example I'm using goes to your R library directory where *poppr* is stored. If you don't know where that is, you can find it by typing `find.package('poppr')` 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 that by 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: `$files` shows the files you selected and `$path` shows 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/3.1" on Windows (where "poppr-user" is your username) or "/Library/Frameworks/R.framework/Versions/3.1/Resources/library" on OSX. Each slash represents a folder that you would click on when you are using the mouse.

A working directory is 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 changing the original file.

To find what your current working directory is set to, 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 folder 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 simple first analysis of your data directly from the file on the your disk (For information on importing your data into R, see section [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	7.00	0.000	1.889	6.23	0.944	0.932	2.92	0.210	rootrot.csv
## 2	Athena_2	12	12	10.00	NaN	2.485	12.00	1.000	1.000	4.16	0.128	rootrot.csv
## 3	Athena_3	10	2	2.00	0.000	0.325	1.22	0.200	0.571	2.00	1.000	rootrot.csv
## 4	Athena_4	13	9	7.15	0.769	1.946	5.12	0.872	0.687	5.49	0.372	rootrot.csv

```
## 5      Athena_5 10  7  7.00 0.000 1.834  5.56 0.911 0.866  4.53 0.353 rootrot.csv
## 6      Athena_6  5  5  5.00 0.000 1.609  5.00 1.000 1.000  2.46 0.190 rootrot.csv
## 7      Athena_7 11 10  9.18 0.386 2.272  9.31 0.982 0.955  2.13 0.086 rootrot.csv
## 8      Athena_8  8  6  6.00 0.000 1.667  4.57 0.893 0.831  3.86 0.323 rootrot.csv
## 9      Athena_9 10 10 10.00 0.000 2.303 10.00 1.000 1.000  2.82 0.118 rootrot.csv
## 10     Athena_10  9  8  8.00 0.000 2.043  7.36 0.972 0.948  2.85 0.137 rootrot.csv
## 11 Mt. Vernon_1 10  9  9.00 0.000 2.164  8.33 0.978 0.952  7.13 0.276 rootrot.csv
## 12 Mt. Vernon_2  6  6  6.00 0.000 1.792  6.00 1.000 1.000 20.65 0.492 rootrot.csv
## 13 Mt. Vernon_3  8  6  6.00 0.000 1.667  4.57 0.893 0.831  2.12 0.106 rootrot.csv
## 14 Mt. Vernon_4 12  8  6.83 0.665 1.814  4.50 0.848 0.681  3.01 0.255 rootrot.csv
## 15 Mt. Vernon_5 17  7  5.54 0.828 1.758  5.07 0.853 0.848  2.68 0.340 rootrot.csv
## 16 Mt. Vernon_6 12 11  9.32 0.466 2.369 10.29 0.985 0.958 19.50 0.467 rootrot.csv
## 17 Mt. Vernon_7 12  9  7.82 0.649 2.095  7.20 0.939 0.870  1.21 0.153 rootrot.csv
## 18 Mt. Vernon_8 13  9  7.35 0.764 2.032  6.26 0.910 0.794  1.15 0.169 rootrot.csv
## 19      Total 187 119  9.61 0.612 4.558 68.97 0.991 0.720 14.37 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][6]
- 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 the function `poppr`. You can access the help page for `poppr` by typing `?poppr` in your R console.

One thing to note about this output is the `NaN` in the column labeled `SE`. In R, `NaN` means “Not a number”. This is produced from calculation of a standard error based on rarefaction analysis. Occasionally, this calculation will encounter a situation in which it must attempt to take a square root of a negative number. Since the root of any negative number is not defined in the set of real numbers, it must therefore have an imaginary component, i . Unfortunately, R will not represent the imaginary components of numbers unless you specifically tell it to do so. By default, R represents these as `NaN`.

1.8 Importing data into poppr {Get out of my dreams and into my R}

There are several ways of reading data into R. One way is using the function `getfile`.

1.8.1 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 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. Instead of doing these one-by-one, `getfile` has a flag called **multi** telling the computer that you want to grab multiple files in the folder:

```
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/rootrot.csv"  "/path/to/R/poppr/files/rootrot2.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! We will use `poppr.all` to produce a summary table for all of your files². Let's set **digits** = 2 to only print 2 significant digits.

```
all_files <- poppr.all(x$files)
print(all_files, digits = 2)
```

```
## \
## | File: rootrot.csv
## /
## | Athena_1
## | Athena_2
## | Athena_3
## | Athena_4
## | Athena_5
```

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

```

## | 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  7.0 0.0e+00 1.89  6.2 0.94 0.93  2.92 0.210  rootrot.csv
## 2      Athena_2  12  12 10.0      NaN  2.48 12.0 1.00 1.00  4.16 0.128  rootrot.csv
## 3      Athena_3  10   2  2.0 0.0e+00 0.33  1.2 0.20 0.57  2.00 1.000  rootrot.csv
## 4      Athena_4  13   9  7.2 7.7e-01 1.95  5.1 0.87 0.69  5.49 0.372  rootrot.csv
## 5      Athena_5  10   7  7.0 0.0e+00 1.83  5.6 0.91 0.87  4.53 0.353  rootrot.csv
## 6      Athena_6   5   5  5.0 0.0e+00 1.61  5.0 1.00 1.00  2.46 0.190  rootrot.csv
## 7      Athena_7  11  10  9.2 3.9e-01 2.27  9.3 0.98 0.96  2.13 0.086  rootrot.csv
## 8      Athena_8   8   6  6.0 0.0e+00 1.67  4.6 0.89 0.83  3.86 0.323  rootrot.csv
## 9      Athena_9  10  10 10.0 0.0e+00 2.30 10.0 1.00 1.00  2.82 0.118  rootrot.csv
## 10     Athena_10  9   8  8.0 0.0e+00 2.04  7.4 0.97 0.95  2.85 0.137  rootrot.csv
## 11     Mt. Vernon_1 10  9  9.0 0.0e+00 2.16  8.3 0.98 0.95  7.13 0.276  rootrot.csv
## 12     Mt. Vernon_2  6   6  6.0 0.0e+00 1.79  6.0 1.00 1.00 20.65 0.492  rootrot.csv
## 13     Mt. Vernon_3  8   6  6.0 0.0e+00 1.67  4.6 0.89 0.83  2.12 0.106  rootrot.csv
## 14     Mt. Vernon_4 12   8  6.8 6.6e-01 1.81  4.5 0.85 0.68  3.01 0.255  rootrot.csv
## 15     Mt. Vernon_5 17   7  5.5 8.3e-01 1.76  5.1 0.85 0.85  2.68 0.340  rootrot.csv
## 16     Mt. Vernon_6 12  11  9.3 4.7e-01 2.37 10.3 0.98 0.96 19.50 0.467  rootrot.csv
## 17     Mt. Vernon_7 12   9  7.8 6.5e-01 2.09  7.2 0.94 0.87  1.21 0.153  rootrot.csv
## 18     Mt. Vernon_8 13   9  7.3 7.6e-01 2.03  6.3 0.91 0.79  1.15 0.169  rootrot.csv
## 19           Total 187 119  9.6 6.1e-01 4.56 69.0 0.99 0.72 14.37 0.271  rootrot.csv
## 20           1  19  16  9.2 7.0e-01 2.73 14.4 0.98 0.94 14.23 0.313  rootrot2.csv
## 21           2  18  18 10.0 5.4e-07 2.89 18.0 1.00 1.00  9.14 0.194  rootrot2.csv
## 22           3  18   8  5.3 1.0e+00 1.61  3.4 0.75 0.59 22.84 0.573  rootrot2.csv
## 23           4  25  17  7.9 1.1e+00 2.58  9.6 0.93 0.71 18.49 0.415  rootrot2.csv
## 24           5  27  14  7.5 1.1e+00 2.45  9.7 0.93 0.83 23.00 0.520  rootrot2.csv
## 25           6  17  15  9.3 6.3e-01 2.67 13.8 0.99 0.95 17.78 0.410  rootrot2.csv
## 26           7  23  19  9.2 7.6e-01 2.87 16.0 0.98 0.90 19.16 0.405  rootrot2.csv
## 27           8  21  15  8.3 9.8e-01 2.56 10.8 0.95 0.82 24.31 0.543  rootrot2.csv
## 28           9  10  10 10.0 0.0e+00 2.30 10.0 1.00 1.00  2.82 0.118  rootrot2.csv
## 29          10  9   8  8.0 0.0e+00 2.04  7.4 0.97 0.95  2.85 0.137  rootrot2.csv
## 30          Total 187 119  9.6 6.1e-01 4.56 69.0 0.99 0.72 14.37 0.271  rootrot2.csv
## 31          Total 100   6  6.0 0.0e+00 1.23  2.8 0.65 0.73  0.05 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 only of a certain type or with a certain name? This is what you would use the `pattern` argument for. For example, there are several data files with different formats in the *adegenet* folder in your R library. 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-HA.fasta"
## [10] "/path/to/R/adegenet/files/pdH1N1-NA.fasta"
## [11] "/path/to/R/adegenet/files/pdH1N1-data.csv"
## [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). To filter these files, use 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 genetix 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!

1.8.2 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 cannot import allelic frequency data from GenAlEx. Using the *poppr* function `read.genalex` will import your data into *adegenet*'s `genind` object or *poppr*'s `genclone` 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

Below is an example of the GenAlEx format. We will use the data set called `microbov` from the *adegenet* package to generate it. The data contains three demographic factors: Country, Species and Breed contained within the `@other` slot (detailed in [THE OTHER SLOT](#)). We will combine these and save the file to the desktop. Details of these functions are presented elsewhere in this manual.

```
library("poppr")
data(microbov)
microbov <- as.genclone(microbov)
sethierarchy(microbov) <- data.frame(other(microbov))
setpop(microbov) <- ~coun/breed/spe
genind2genalex(microbov, file=~"/Desktop/microbov.csv")
```

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

The GenAlEx format contains individuals in rows and loci in columns. Individual data begins at row 4. Column A always contains individual names and column B defines the population of each individual. Notice here that the three demographic factors from the data have been concatenated with a "_". This allows us to import more than one population factor to use as hierarchical levels in a [GENCLONE OBJECT](#).

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

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	J
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 Bazadals	
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

The First three rows contain information pertaining to the global data set. The only important information for *poppr* is the information contained in row 3 and the first three columns of row 1.

	A	B	C	D
1	# of Loci	# of Individuals	# of Populations	Pop1 Size ...
2	-	-	-	Pop1 Name ...
3	-	-	Locus 1	...

Highlighted in red in figure 2 are definitions of the number of populations and their respective sizes. As this is redundant information, we can remove it. Below is an example of a valid data set that can be imported into *poppr*.

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

All GenAlEx formatted data can be imported with the command `read.genalex`, detailed below:

Default Command:

```
read.genalex(genalex, ploidy = 2, geo = FALSE, region = FALSE,  
  genclone = TRUE, sep = ",")
```

- **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).
- **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 [THE OTHER SLOT](#).)
- **region** - To format your GenAlEx file to include regions, 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 **@other** slot of the resulting *genind* object or be incorporated into the hierarchy of the *genclone* object.
- **genclone** - This flag will convert your data into a **genclone** object (see [SEND IN THE CLONES](#) for more info).
- **sep** - The separator argument for columns in your data. It defaults to ",".

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 `genclone` format and can be used with any function in *poppr* and *adegenet*

```
rootrot
##
## This is a genclone object
## -----
## Genotype information:
##
##   119 multilocus genotypes
##   187 diploid individuals
##    56 dominant loci
##
## Population information:
##
##    1 hierarchical level - Pop
##   18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8
```

1.8.3 Other ways of importing data

Adegenet already supports the import of FSTAT, Structure, Genpop, and Genetix formatted files, so if you have data in 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.

1.8.4 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,
  geo = FALSE, geodf = "xy", sep = ",")
```

- `pop` - a `genind` object.
- `filename` - This is where you specify the path to the new file you wish to create. If you specify only a filename with no path, it will place the file in your current working directory.
- `quiet` - If this is set to `FALSE`, a status message will be printed to the console as the extraction progresses.
- `geo` - Set to `TRUE`, if 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 the resulting file will 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.

- `sep` - A separator to separate columns in the resulting file.

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.
```

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
```

To export it:

```
genind2genalex(nancycats, "~/Desktop/nancycats_pop_xy.csv", geo = TRUE)
```

```
## 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 use this 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", geo = TRUE)
```

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

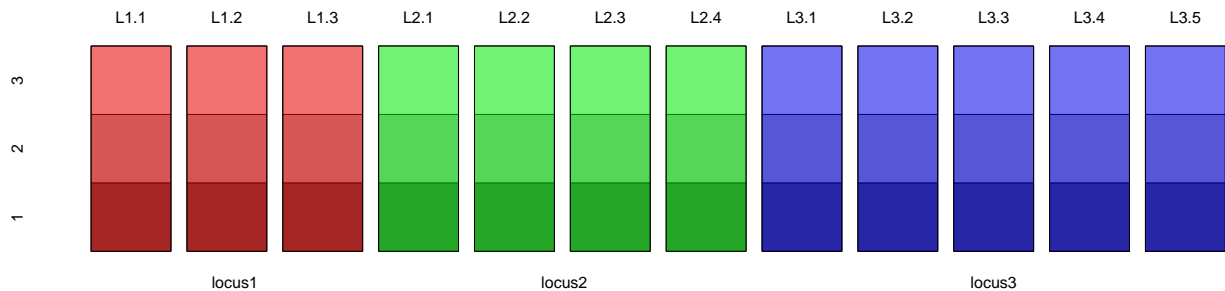
1.9 Getting to know *adegenet*'s genind object

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, Let's say we had data with 3 diploid individuals each with 3 loci that had 3, 4, and 5 allelic states respectively:

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

The resulting **genind** object would contain a matrix that has 3 rows and 12 columns. Below is a schematic of what that would look like. Each column represents a separate allele, each row represents an individual and each color represents a different locus.



When we look at the data derived from table 1.9, we see that we have a matrix of individual allele frequencies at each locus.

```
##  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  0.0   1  0.0  0.0  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
```

At each locus, the allele frequencies for each individual sum to one. Homozygotes are denoted as having an allele frequency of 1 at a particular allele while heterozygotes have their allele frequencies represented as $1/p$ where p = ploidy. Along with this matrix, are elements that define the names of the individuals, loci, alleles, and populations. If you wish to know more, see the *adegenet* "Getting Started" manual.

1.9.1 The other slot

The other slot is a place in the genind object that can be used to store useful information about the data. We saw earlier that it could store demographic information, now let's explore a different example. Bruvo's distance is based off of a stepwise mutation model for microsatellites. This requires us to know the length

of the repeat of each locus. We could store the repeat lengths in a separate variable in our R environment, but we are at risk of losing that. One way to prevent it from being lost would be to place it in the “other” slot. For the purpose of this example, we will use the “nancycats” data set from the *adegenet* package and assume that it has dinucleotide repeats at all of its loci.

```
data(nancycats) # Load the data
other(nancycats) # geographical coordinates

## $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

repeats <- rep(2, nLoc(nancycats)) #nLoc = number of loci
repeats

## [1] 2 2 2 2 2 2 2 2 2

other(nancycats)$repeat_lengths <- repeats
other(nancycats) # two items named xy and repeat_lengths

## $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
##
## $repeat_lengths
## [1] 2 2 2 2 2 2 2 2 2
```

1.10 The genclone object {send in the clones}

The **genclone** class was defined in order to make working with hierarchies more intuitive. It is built off of the **genind** object and has dedicated slots for the population hierarchy and defined multilocus genotypes. The name **genclone** refers to the fact that it has the ability to handle genotypes of clonal organisms (but it is also used for sexual populations).

The main difference between the **genclone** and **genind** objects is how they handle populations: With the **genind** object, the user must find the vector defining the population and set it using that vector. The **genclone** object already defines a data frame with different population factors or hierarchical levels in the object. The user simply supplies a formula defining the desired hierarchy with which to set the population. This formula driven method is also used for clone correction, combining hierarchical levels and conducting AMOVA. These will all be explained in later chapters. For examples and details, type `help("genclone")` in your R console.

The function `as.genclone` allows the user to convert a **genind** object to a **genclone** object. The following example will demonstrate that the **genclone** object is an extension of the **genind** object as well as the advantages of having populations pre-defined in your data set.

1.10.1 Function: as.genclone

Default Command:

```
as.genclone(x, hierarchy = NULL)
```

- **x** - a **genind** object to be converted.
 - **hierarchy** - an optional data frame where each column represents a hierarchical level of the population hierarchy in the data set.
-

Let's show an example of a **genclone** object. First, we will take an existing **genind** object and convert it using the function `as.genclone` (We can also use the function `read.genalex` to import as **genclone** or **genind** objects). We will use the **Aeut** data set because it is a clonal data set that has a simple population hierarchy [5]. The data set is here: <http://dx.doi.org/10.6084/m9.figshare.877104> and it is AFLP data of the root rot pathogen *Aphanomyces euteiches* collected from two different fields in NW Oregon and W Washington, USA. These fields were divided up into subplots from which samples were collected. The fields represent the population and the subplots represent the subpopulation. Let's take a look at what the **genind** object looks like:

```
library("poppr")
data(Aeut)
Aeut

##
## #####
## ### 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
##
## Optional contents:
## @pop: factor giving the population of each individual
## @pop.names: factor giving the population of each individual
##
## @other: a list containing: population_hierarchy
```

This gives us a lot of information about the object, and is useful once you become more comfortable with programming. Once we convert it to a `genclone` object, the multilocus genotypes will be defined and the population hierarchy (if a data frame is defined in the `@other` slot called “population_hierarchy”) will be set.

```
agc <- as.genclone(Aeut)
agc

##
## This is a genclone object
## -----
## Genotype information:
##
## 119 multilocus genotypes
## 187 diploid individuals
## 56 dominant loci
##
## Population information:
##
## 3 hierarchical levels - Pop Subpop Pop Subpop
## 2 populations defined - Athena Mt. Vernon

# We can also manually set the hierarchy.
as.genclone(Aeut, hierarchy = other(Aeut)$population_hierarchy[-1])

##
## This is a genclone object
## -----
## Genotype information:
##
## 119 multilocus genotypes
## 187 diploid individuals
## 56 dominant loci
##
## Population information:
##
## 2 hierarchical levels - Pop Subpop
## 2 populations defined - Athena Mt. Vernon
```

We can see here that it shows less information, but it gives us a very simple overview of our data. Don’t be fooled, however, because it contains the same information as a `genind` object and the advantage of the `genclone` object is that setting the population from the hierarchy becomes **much** easier.

```
c(is.genind(Aeut), is.genclone(Aeut), is.genind(agc), is.genclone(agc))

## [1] TRUE FALSE TRUE TRUE

# Adegenet functions work the same, too
c(nInd(Aeut), nInd(agc))
```

```
## [1] 187 187
```

```
# We'll look at the population names
Aeut$pop.names

##          P1          P2
##    "Athena" "Mt. Vernon"

# genind way:
# Extract the combined hierarchical levels.
pophier <- other(Aeut)$population_hierarchy$Pop_Subpop
pop(Aeut) <- pophier
Aeut$pop.names

## [1] "Athena_1"      "Athena_2"      "Athena_3"      "Athena_4"      "Athena_5"
## [6] "Athena_6"      "Athena_7"      "Athena_8"      "Athena_9"      "Athena_10"
## [11] "Mt. Vernon_1" "Mt. Vernon_2" "Mt. Vernon_3" "Mt. Vernon_4" "Mt. Vernon_5"
## [16] "Mt. Vernon_6" "Mt. Vernon_7" "Mt. Vernon_8"

# genclone way:
agc

##
## This is a genclone object
## -----
## Genotype information:
##
##    119 multilocus genotypes
##    187 diploid individuals
##    56 dominant loci
##
## Population information:
##
##    3 hierarchical levels - Pop_Subpop Pop Subpop
##    2 populations defined - Athena Mt. Vernon

setpop(agc) <- ~Pop/Subpop
agc

##
## This is a genclone object
## -----
## Genotype information:
##
##    119 multilocus genotypes
##    187 diploid individuals
##    56 dominant loci
##
## Population information:
##
##    3 hierarchical levels - Pop_Subpop Pop Subpop
##    18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8

# Notice that you only see the first and last three population names.
# Use the print function to display all population names.

print(agc)


##
## This is a genclone object
## -----
## Genotype information:
```

```
##
## 119 multilocus genotypes
## 187 diploid individuals
## 56 dominant loci
##
## Population information:
##
## 3 hierarchical levels - Pop_Subpop Pop Subpop
## 18 populations defined - 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
```

1.11 Accessing the population hierarchy

The hierarchy slot in the `genclone` object allows us to access and manipulate different levels of a population hierarchy without the complication of creating new data sets. We implemented the use of hierarchical nested formulae:

```
hier = ~Population/Subpopulation/Year
```

 **NOTE:** The `~` symbol is absolutely required for formulas, even if you only are specifying one level. The `"/"` symbolizes a hierarchical nesting. The above formula represents years nested within subpopulations, nested within populations. This allows the user to easily restructure the population hierarchy. Manipulation of the hierarchy is necessary for [CLONE CORRECTION](#) and AMOVA analysis. See the section [CAN YOU TAKE ME HIER\(ARCHY\)](#) for more details on manipulation of hierarchies.

1.12 About polyploid data



WARNING

Treat polyploid data with care. Please read this section carefully and consult the help pages for all functions mentioned here.

With diploid or haploid data, genotypes are unambiguous. It is often clear when it is homozygous or heterozygous. With polyploid data, genotypes can be ambiguous. For example, a tetraploid individual with the apparent genotype of **A/B** could actually have one of three genotypes: **A/A/A/B**, **A/A/B/B**, or **A/B/B/B**. This ambiguity prevents a researcher from accurately calling all alleles present. In *adegenet*, it was previously difficult to import polyploid data because of this ambiguity as data was required to be unambiguous or missing.

A solution to this problem is to code missing alleles as “0”. An example of this is found within the `Pinf` data set in *poppr* [4]. We look at the last six samples over two loci, `Pi63` (3 alleles, triploid) and `Pi70` (3 alleles, diploid) to examine how the data is represented.

```
data(Pinf)
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
## 72 multilocus genotypes
## 86 tetraploid individuals
## 11 codominant loci
##
```

```
## Population information:
##
##      2 hierarchical levels - Continent Country
##      2 populations defined - South America North America

tail(truenames(Pinf[loc = c("L09", "L10")])$tab)

##      Pi63.000 Pi63.148 Pi63.151 Pi63.157 Pi70.000 Pi70.189 Pi70.192 Pi70.195
## PiPE22      0.25      0.25      0.25      0.25      0.5      0.25      0.25      0.00
## PiPE23      0.25      0.25      0.25      0.25      0.5      0.25      0.25      0.00
## PiPE24      0.25      0.25      0.25      0.25      0.5      0.25      0.25      0.00
## PiPE25      0.25      0.25      0.25      0.25      0.5      0.25      0.25      0.00
## PiPE26      0.50      0.00      0.00      0.50      0.5      0.00      0.25      0.25
## PiPE27      0.25      0.25      0.25      0.25      0.5      0.25      0.25      0.00
```

Each column in this data represents a different allele at a particular locus. Pi63.148 is the allele 148 at locus Pi63. Each row is an individual. The numbers represent the fraction of a given allele that makes up the individual genotype at a particular locus. What we can see here is that the number of columns is 8 when we expect only 6 (2 loci \times 3 alleles). The first allele at each locus is 000. Let's take a look at the data in a human-readable format.

```
Pinfdf <- genind2df(Pinf, sep = "/")
tail(Pinfdf[10:11])

##      Pi63      Pi70
## PiPE22 000/148/151/157 000/000/189/192
## PiPE23 000/148/151/157 000/000/189/192
## PiPE24 000/148/151/157 000/000/189/192
## PiPE25 000/148/151/157 000/000/189/192
## PiPE26 000/000/157/157 000/000/192/195
## PiPE27 000/148/151/157 000/000/189/192
```

It's more clear now that we have a data set of tetraploid individuals where some genotypes appear diploid (000/000/157/157) and some appear triploid (000/148/151/157). The tetraploid genotype is padded with zeroes to make up the difference in ploidy.

This method allows BRUVO's DISTANCE [3] and the INDEX OF ASSOCIATION [2, 17, 1] to work with polyploids as they specifically recognize the zeroes as being missing data. A side effect, unfortunately is that the extra zeroes appear as extra alleles. As this affects all frequency-based statistics (except for the ones noted above), the user should reformat their data set with the function `recode_polyploids`, which will remove the zeroes and recode the allele frequencies to the observed frequencies.

```
Pinf_rc <- recode_polyploids(Pinf, newploidy = 2)
tail(truenames(Pinf_rc[loc = c("L09", "L10")])$tab)

##      Pi63.148 Pi63.151 Pi63.157 Pi70.189 Pi70.192 Pi70.195
## PiPE22 0.3333333 0.3333333 0.3333333      0.5      0.5      0.0
## PiPE23 0.3333333 0.3333333 0.3333333      0.5      0.5      0.0
## PiPE24 0.3333333 0.3333333 0.3333333      0.5      0.5      0.0
## PiPE25 0.3333333 0.3333333 0.3333333      0.5      0.5      0.0
## PiPE26 0.0000000 0.0000000 1.0000000      0.0      0.5      0.5
## PiPE27 0.3333333 0.3333333 0.3333333      0.5      0.5      0.0
```

Notice that the triploid locus now has frequencies that are multiples of $\frac{1}{3}$ and the diploid locus has multiples of $\frac{1}{2}$. Below, we show the observed genotypes:

```
Pinfrcdf <- genind2df(Pinf_rc, sep = "/")
tail(Pinfrcdf[10:11])

##      Pi63      Pi70
## PiPE22 148/151/157 189/192
## PiPE23 148/151/157 189/192
```

```
## PiPE24 148/151/157 189/192
## PiPE25 148/151/157 189/192
## PiPE26      157/157 192/195
## PiPE27 148/151/157 189/192
```

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

in *poppr*, the **GENCLONE** object contains a slot called **hierarchy**. This slot contains a data frame used to define hierarchical levels of population factors describing your data. The preferred way of defining these hierarchical levels is to concatenate them using ‘.’ and use them to define a single population in your data before you import it into *poppr*. Examples of this format can be found in figures 2, 3, and at <http://dx.doi.org/10.6084/m9.figshare.877104>.

In this section, we will show you how to **DEFINE** hierarchical levels, **VIEW** those levels to ensure that they are correctly defined, **MANIPULATE** your hierarchical levels by adding and renaming them, and use these levels to **SET THE POPULATION** in your **GENCLONE** object using the following methods:

Method	Function	Input	Result
split	splithierarchy	formula	defined hierarchical levels
set	sethierarchy	data frame	new hierarchical levels
get	gethierarchy	formula	data frame
name	namehierarchy	formula	new hierarchical level names
add	addhierarchy	vector or data frame	new hierarchical level



A NOTE ABOUT FORMULAS

The formulas used by *genclone* objects always start with a ~ and are hierarchical levels are always separated by a /. Some examples are:

~Country/City/District

~Field/Year

Refer to [ACCESSING HIERARCHIES](#) for more details on how to access hierarchies.

In the next section, we’ll explore two ways of defining hierarchical levels.

2.1.1 Defining hierarchies

As explained above, the best way to define hierarchical levels is to concatenate them using ‘.’ and set that as your population factor. We will use the example data set from <http://dx.doi.org/10.6084/m9.figshare.877104>. It is an AFLP data set of the root rot pathogen *Aphanomyces euteiches* from two fields and multiple soil cores per field. First, we will follow the link and copy the download link from figshare.

```
aphan <- read.genalex("http://files.figshare.com/1314228/rootrot.csv")
aphan
```



```
##
## This is a genclone object
## -----
## Genotype information:
##
##      119 multilocus genotypes
##      187 diploid individuals
##      56 dominant loci
##
## Population information:
##
##      1 hierarchical level - Pop
##      18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8
```

The supplemental information in the data defined two hierarchical levels, yet we only see one here labeled 'Pop'. This is how populations are automatically defined when importing to a [GENCLONE](#) object since it does not know how many hierarchical levels you have defined. To define these levels present in the data set, we will need to split them up using the function `splithierarchy`:

```
splithierarchy(aphan) <- ~field/sample
aphan

##
## This is a genclone object
## -----
## Genotype information:
##
##      119 multilocus genotypes
##      187 diploid individuals
##      56 dominant loci
##
## Population information:
##
##      2 hierarchical levels - field sample
##      18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8
```

Now we have successfully defined our hierarchies. If you have imported your data in this manner, you may skip to the [SETTING POPULATION HIERARCHIES](#), [VIEWING HIERARCHIES](#), or [MANIPULATING HIERARCHICAL LEVELS](#).

If you have imported your data with a single population and want to add hierarchical levels separately, you can use the function `sethierarchy` with a data frame containing your hierarchical levels. For this example, we will use the data set H3N2, which contains SNP data from the H3N2 virus. This data set holds a data frame in the [OTHER SLOT](#) that contains many variables including country, year, and month of collection. We will first load the data and write that data frame to a file on the desktop.

```
data(H3N2)
write.table(other(H3N2)$x, file = "~/Desktop/virus_info.csv", row.names = FALSE)
```

Now we have our data and we have a separate table in a file on our desktop defining our hierarchical levels. Let's import those levels into R with `read.table` and see what they are:

```
virus_info <- read.table("~/Desktop/virus_info.csv", header = TRUE)
names(virus_info)
```

```
## [1] "accession"      "length"         "host"           "segment"        "subtype"
## [6] "country"        "year"           "Virus name"     "Misc info"      "Age"
## [11] "Gender"         "date"           "usePreciseLoc"  "localisation"   "lon"
## [16] "lat"           "month"
```

From here we will convert our `genind` object to a `genclone` object and use `sethierarchy` to define the hierarchical levels with the table we just imported.

```
virus <- as.genclone(H3N2) # Converting it to a genclone object.
sethierarchy(virus) <- virus_info # Setting the hierarchy
virus

##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##     1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##      17 hierarchical levels - accession length host ... lon lat month
##       0 populations defined.
```

In this data, levels such as `host` and `segment` are unimportant levels because they are all the same. Let's say that we are only interested in `year` and `country`. To make things easier to view, we will set the hierarchical levels to these two columns:

```
sethierarchy(virus) <- virus_info[c("country", "year")]
virus

##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##     1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##       2 hierarchical levels - country year
##       0 populations defined.
```

Notice that there are no populations defined. Now that we have the hierarchical levels in place, we can use it to define the population hierarchy. We will use the function `setpop` to define the population as year with respect to country:

```
setpop(virus) <- ~year/country
virus

##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##     1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##       2 hierarchical levels - country year
##      102 populations defined - 2002_Japan 2002_USA 2002_Finland ... 2006_Algeria
## 2006_Sweden 2006_Greece
```

2.1.2 Viewing hierarchies

If you wanted to view your hierarchies to make sure that you made no spelling errors in your population definitions, you can extract the data frame from your `genclone` object by using the function `gethierarchy`:

Default Command:

```
gethierarchy(x, formula = NULL, combine = TRUE)
```

Where **x** is the `genclone` object, **formula** defines the hierarchical levels, and **combine** indicates whether or not you want the lower levels of the hierarchy combined with the higher levels. For example, in the root rot data above, the hierarchical levels are explicitly hierarchical and should be combined. Note, if you don't supply a formula argument, the data frame as it exists will be returned.

```
head(gethierarchy(aphan))

##    field sample
## 1 Athena      1
## 2 Athena      1
## 3 Athena      1
## 4 Athena      1
## 5 Athena      1
## 6 Athena      1

head(gethierarchy(aphan, ~field/sample))

##    field  sample
## 1 Athena Athena_1
## 2 Athena Athena_1
## 3 Athena Athena_1
## 4 Athena Athena_1
## 5 Athena Athena_1
## 6 Athena Athena_1
```

If the hierarchical levels are not nested, or you simply want to use this hierarchy for another data set, you might want to set the **combine** flag to **FALSE**. Let's use the virus data as an example:

```
head(gethierarchy(virus, ~year/country))

##   year  country
## 1 2002 2002_Japan
## 2 2002 2002_Japan
## 3 2002 2002_Japan
## 4 2002 2002_Japan
## 5 2002 2002_Japan
## 6 2002 2002_Japan

head(gethierarchy(virus, ~year/country, combine = FALSE))

##   year country
## 1 2002   Japan
## 2 2002   Japan
## 3 2002   Japan
## 4 2002   Japan
## 5 2002   Japan
## 6 2002   Japan
```

It will return only the levels you ask it to return:

```
head(gethierarchy(virus, ~country))
```

```
## country
## 1 Japan
## 2 Japan
## 3 Japan
## 4 Japan
## 5 Japan
## 6 Japan
```

2.1.3 Manipulating hierarchical levels

Once we have our hierarchies set in place, we want to be able to rename and add to them. For this example, we will revisit the [VIRUS EXAMPLE](#) from above. We have set the population hierarchy to be based on year and country, but we've noticed that we left out month. And let's say that we accidentally overwrote the data object like this:

```
virus_info <- virus_info[["month"]]
names(virus_info)

## NULL
```

If we were saving our script the whole time, we could just go back and retrieve the data frame, but that defeats the purpose of this section where we imagine that we've recieved new information and wanted to add it to our hierarchy. If we want to add this to our hierarchy, we just use the function `addhierarchy` defined as thus:

Default Command:

```
addhierarchy(x, value, name = "NEW")
```

We can use this function to add on a new column to the data frame.

```
addhierarchy(virus) <- virus_info
virus

##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##     1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##      3 hierarchical levels - country year NEW
##     102 populations defined - 2002_Japan 2002_USA 2002_Finland ... 2006_Algeria
## 2006_Sweden 2006_Greece
```

Notice that the new hierarchical level is simply labeled as `NEW`. We will customize the name of the hierarchical levels with the function `namehierarchy`.

```
namehierarchy(virus) <- ~country/year/month
virus
```

```
##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##      1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##      3 hierarchical levels - country year month
##      102 populations defined - 2002_Japan 2002_USA 2002_Finland ... 2006_Algeria
## 2006_Sweden 2006_Greece
```

Of course, perhaps a better way still would be to use a data frame:

```
addhierarchy(virus) <- data.frame(month = virus_info)
```

2.1.4 Defining populations with hierarchies

Now that we have defined the hierarchical levels in the data set, setting the population hierarchy allows us to group our data according to the hierarchical level of your choice. This is a necessary step. For this example, we will use a data set of *Phytophthora infestans* collected from North America and South America.

```
data(Pinf)
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
##
##      2 hierarchical levels - Continent Country
##      2 populations defined - South America North America
```

Above we have two hierarchies for Continent and Country, but the populations only show Continent level populations. If we wanted to investigate each country separately, we would need to reset the population to be represented by Country. This can be done with the function `setpop`. This function utilizes the defined population hierarchies to set the population. We'll use our data set above to illustrate this:

```
setpop(Pinf) <- ~Country
Pinf # Now set by country

##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
```

```
##
##      2 hierarchical levels - Continent Country
##      4 populations defined - Colombia Ecuador Mexico Peru
```

The beauty about it is the fact that it will also combine all the hierarchical levels you want to use. Let's see when we ask it to set the population of Country with respect to Continent.

```
setpop(Pinf) <- ~Continent/Country
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
##
##      2 hierarchical levels - Continent Country
##      4 populations defined - South America_Colombia South America_Ecuador
## North America_Mexico South America_Peru
```

Nice!

2.2 Replace or remove missing data {Inside the golden days of 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.



WARNING

Treatment of Missing data is a non-trivial task. You should understand the nature of missing data in your data set before treatment.

2.2.1 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.

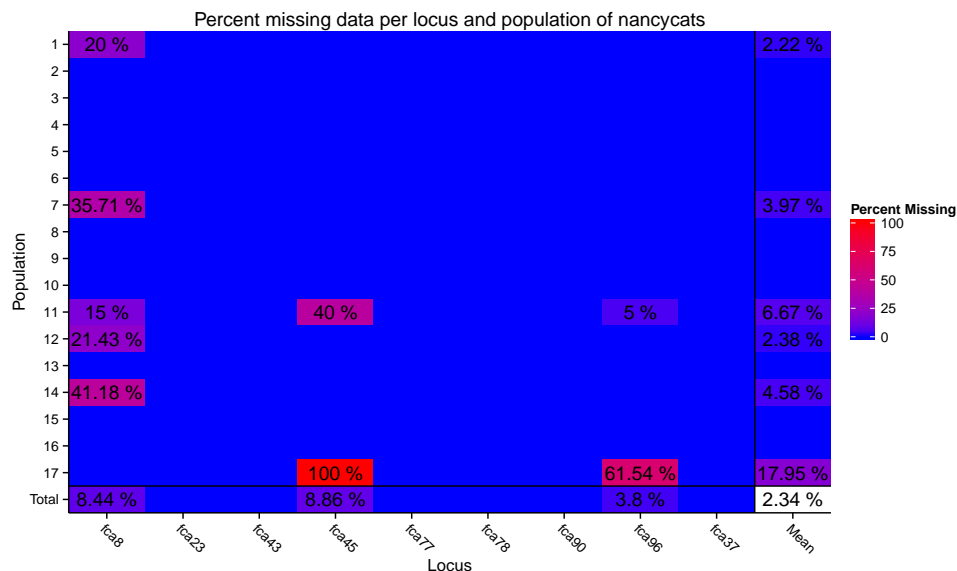
Default Command:

```
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”.

Let’s take a look at what this does by focusing in on areas with missing data. We’ll use the data set `nancycats` as an example. Using the `poppr` function `info_table`, we can assess missing data within populations.

```
library("poppr")
data(nancycats)
info_table(nancycats, plot = TRUE)
```



```
##      Locus
## Population fca8 fca23 fca43 fca45 fca77 fca78 fca90 fca96 fca37 Mean
##      1      0.200 .      .      .      .      .      .      .      0.022
```

```
##      2      .      .      .      .      .      .      .      .
##      3      .      .      .      .      .      .      .      .
##      4      .      .      .      .      .      .      .      .
##      5      .      .      .      .      .      .      .      .
##      6      .      .      .      .      .      .      .      .
##      7      0.357      .      .      .      .      .      .      0.040
##      8      .      .      .      .      .      .      .      .
##      9      .      .      .      .      .      .      .      .
##     10      .      .      .      .      .      .      .      .
##     11      0.150      .      0.400      .      .      0.050      0.067
##     12      0.214      .      .      .      .      .      .      0.024
##     13      .      .      .      .      .      .      .      .
##     14      0.412      .      .      .      .      .      .      0.046
##     15      .      .      .      .      .      .      .      .
##     16      .      .      .      .      .      .      .      .
##     17      .      .      1.000      .      .      0.615      0.179
##     Total 0.084      .      0.089      .      .      0.038      0.023
```

We can see that locus fca8 has a lot of missing data. To demonstrate the function `missingno`, we will zoom into the first five individuals at the first locus.

```
nancycats$stab[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). 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$stab[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$stab[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 N298 N299 N300
## N301 N302 N303 N304 N310 N195 N197 N198 N199 N200 N201 N206 N182 N184 N186 N282
## N283 N288 N291 N292 N293 N294 N295 N296 N297 N281 N289 N290

nanloci$stab[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.

```
nInd(nanloci) # Individuals

## [1] 237

locNames(nanloci) # 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$stab[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

nInd(nangeno) # Individuals

## [1] 199

locNames(nangeno) # 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.

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

Adegenet provides methods for subsetting data by individual or splitting all of the data into a list of populations. If you only want one or two populations, these methods become tedious. The *poppr* function `popsub` makes this easier:

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 also use this within a function to avoid creating a new variable to keep track of.

Default Command:

```
popsub(gid, sublist = "ALL", blacklist = NULL, mat = NULL, drop = TRUE)
```

- `pop` - a `genind` object.
 - `sublist` - 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` - 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 [MULTILOCUS GENOTYPE ANALYSIS](#) for more information) A matrix produced from the `mlg.table` function. This overrides the `pop` argument and subsets this table instead.
-

To demonstrate this tool, let’s revisit the [VIRUS](#) data set that we saw in [DEFINING HIERARCHIES](#). 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”.

```
setpop(virus) <- ~country
virus$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"
```

```
vna <- popsub(virus, sublist=c("USA", "Canada"))
vna$pop.names

##      P1      P2
##    "USA" "Canada"
```

If we want to see the population size, we can use the *adeigenet* function `nInd()`:

```
c(NorthAmerica = nInd(vna), Total = nInd(virus))

## NorthAmerica      Total
##           665      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.

```
vnaminus <- popsub(virus, blacklist=c("USA", "Canada"))
vnaminus$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 is equal to the number of individuals in our original data set:

```
(nInd(vnaminus) + nInd(vna)) == nInd(virus)

## [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`.

```
vsort <- sort(virus$pop.names)[1:10]
vsort

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

valph <- popsub(virus, sublist=vsort, blacklist=c("USA", "Canada"))
valph$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 Clone-censor data sets {Attack of the clone correction}

Clone correction refers to the ability of keeping one observation of each MLG 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 that will correct down to 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.

This function will create new data sets, but it is also utilized by the functions `poppr` and `poppr.amova` natively.

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 = 1, dfname = "population_hierarchy",
             combine = FALSE, keep = 1)
```

- `pop` - a `genclone` object with a defined hierarchy or 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` - A hierarchical formula (eg. `~Pop/Subpop`), representing the hierarchical levels in your data.
- `dfname` - **Only for use in `genind` objects, otherwise, deprecated** 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 an object with the top most hierarchical level as a population factor unless the `keep` argument is defined. If set to `TRUE`, the hierarchy will be returned combined.
- `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 defined as integers. 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 that we loaded earlier in the section [CAN YOU TAKE ME HIER\(ARCHY\)](#) since that data set is known to include clonal populations [5]. Try playing around with the data and see what different combinations of the `hier`, and `keep` flags produce. Below, I will give three examples of clone corrections at the sample level with respect to field, at the field level, and finally, at the level of the entire data set.

First, we will examine the original data set.

```
aphan # Original object

##
## This is a genclone object
## -----
## Genotype information:
##
##   119 multilocus genotypes
##   187 diploid individuals
##    56 dominant loci
##
## Population information:
##
##    2 hierarchical levels - field sample
##   18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8
```

Now we correct by sample with respect to field and keep the field as the population.

```
clonecorrect(aphan, hier = ~field/sample)

##
## This is a genclone object
## -----
## Genotype information:
##
##      119 multilocus genotypes
##      141 diploid individuals
##      56 dominant loci
##
## Population information:
##
##      2 hierarchical levels - field sample
##      2 populations defined - Athena Mt. Vernon

# Your turn: Use the same hierarchy and use combine = TRUE and then
# keep = 1:2. Is there any difference?
```

Correcting by field. Notice how the number of MLG is much closer to our census.

```
clonecorrect(aphan, hier = ~field)

##
## This is a genclone object
## -----
## Genotype information:
##
##      119 multilocus genotypes
##      120 diploid individuals
##      56 dominant loci
##
## Population information:
##
##      2 hierarchical levels - field sample
##      2 populations defined - Athena Mt. Vernon
```

Correcting over whole data set. Our MLG is equal to our census.

```
clonecorrect(aphan, hier = NA)

##
## This is a genclone object
## -----
## Genotype information:
##
##      119 multilocus genotypes
##      119 diploid individuals
##      56 dominant loci
##
## Population information:
##
##      2 hierarchical levels - field sample
##      18 populations defined - Athena_1 Athena_2 Athena_3 ... Mt. Vernon_6 Mt. Vernon_7
## Mt. Vernon_8
```

2.5 Permutations and bootstrap resampling {every day I'm shuffling (data sets)}

A common null hypothesis for populations with mixed reproductive modes is panmixia, or to put it simply: lots of sex. *Popprr* randomly shuffles data sets in order to calculate P-values for the index of

association (I_A and \bar{r}_d) [1] using 4 different methods:

method	strategy	units sampled
1	permutation	alleles
2	simulation	alleles
3	simulation	alleles
4	permutation	genotypes

These methods are detailed below. We will create a dummy data set to be shuffled by each example below. Let's assume 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:

A1/A2	
1	4/4
2	4/1
3	4/3
4	2/2
5	3/3

Table 1: Original

The 4 methods are detailed 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. **Permute Alleles (default)** This is a sampling scheme that will **permute alleles within the locus**.

The example above might become tables 2 and 3.

A1/A2	
1	3/4
2	2/3
3	4/4
4	2/1
5	3/4

Table 2: Permute 1

A1/A2	
1	1/3
2	2/4
3	3/4
4	4/3
5	4/2

Table 3: Permute 2

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

	A1/A2
1	1/3
2	3/3
3	3/2
4	4/4
5	4/2

Table 4: Parametric 1

	A1/A2
1	3/4
2	2/3
3	4/2
4	4/4
5	4/2

Table 5: Parametric 2

2. **Parametric Bootstrap** The previous scheme reshuffled the observed sample, but the parametric bootstrap **draws samples from a multinomial distribution using the observed allele frequencies as weights**. Tables 4 and 5 are examples of what I mean.

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.

3. **Non-Parametric Bootstrap** The third method is sampling with replacement, again **drawing from a multinomial distribution, but with no assumption about the allele frequencies** (tables 6 and 7).

	A1/A2
1	1/3
2	3/3
3	3/1
4	2/2
5	3/1

Table 6: Non-parametric 1

	A1/A2
1	1/3
2	3/1
3	2/3
4	2/1
5	4/3

Table 7: Non-parametric 2

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

4. **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. Using our example above, tables 8 and 9 are shuffled with method 4.

	A1/A2
1	3/3
2	4/1
3	2/2
4	4/4
5	4/3

Table 8: ML 1

	A1/A2
1	4/4
2	2/2
3	3/3
4	4/3
5	4/1

Table 9: ML 2

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. 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**.

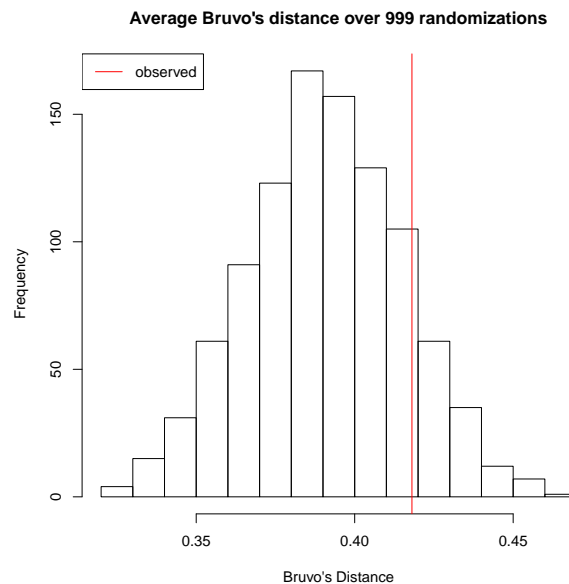
These shuffling schemes have been implemented for the index of association, but there may be other summary statistics you can use `shufflepop` for. All you have to do is use the function `replicate`. Let's use average Bruvo's distance with the first population of the data set `nancycats` as an example:

```
data(nancycats)
nan1 <- popsub(nancycats, 1)
reps <- rep(2, 9) # Assuming dinucleotide repeats.
observed <- mean(bruvo.dist(nan1, replen = reps))
observed
## [1] 0.4180619
```

```
set.seed(9999)
bd.test <- replicate(999, mean(bruvo.dist(shufflepop(nan1, method = 2), replen = reps)))
```

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.

```
hist(bd.test, xlab = "Bruvo's Distance", main = "Average Bruvo's distance over 999 randomizations")
abline(v = observed, col = "red")
legend('topleft', legend="observed", col="red", lty = 1)
```



2.6 Removing uninformative loci {Cut It Out!}

Phylogenetically uninformative loci are those that have only one sample differentiating from the rest. This can lead to biased results when using multilocus analyses such as the index of association [2, 17]. These nuisance loci can be removed with the following function.

2.6.1 Function: `informloci`

Default Command:

```
informloci(pop, cutoff = 2/nInd(pop), quiet = FALSE)
```

- **pop** - a **genind** object.
- **cutoff** - this represents the minimum fraction of individuals needed for a locus to be considered informative. The default is set to $2/n$ with n being the number of individuals in the data set (represented by the *adegenet* function **nInd**). Essentially, this means that any locus with fewer than 2 observations differing will be removed. The user can also specify a fraction of observations for the cutoff (eg. 0.05).
- **quiet** - if **TRUE**, nothing will be printed to the screen, if **FALSE** (default), the cutoff value in percentage and number of individuals will be printed as well as the names of the uninformative loci found.

Here's a quick example using the H3N2 virus SNP data set. We will only retain loci that have a minor allele frequency of $\geq 5\%$

```
data(H3N2)
H.five <- informloci(H3N2, cutoff = 0.05)
```

```
## cutoff value: 5 percent ( 95 individuals ).
## 47 uninformative loci found: 157
## 177 233 243 262 267 280 303 313 327 357 382 384 399 412 418 424 425 429 433 451
## 470 529 546 555 557 564 576 592 595 597 602 612 627 642 647 648 654 658 663 667
## 681 717 806 824 837 882
```

Now what happens when you have all informative loci. We'll use the *nancycats* data set, which has microsatellite loci. It is important to note that this is searching for loci with a specified genotype frequency as fixed heterozygous sites are also uninformative:

```
data(nancycats)
naninform <- informloci(nancycats, cutoff = 0.05)

## cutoff value: 5 percent ( 12 individuals ).
## No sites found with fewer than 12 different individuals.
```

3 Multilocus Genotype Analysis

In populations with mixed sexual and clonal reproduction, it common to have multiple samples from the same population that have the same set of alleles at all loci. Here, we introduce tools for tracking MLGs within and across populations in **GENIND** objects from the *adegenet* package. We will be using the **VIRUS** data set containing SNP data from isolates of the H3N2 virus from 2002 to 2006. Note that **genclone** objects are optimal for these analyses.

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

Counting the number of MLGs in a population is the first step for these analyses as they allow us to see how many clones exist. With the **GENCLONE** object, This information is already displayed when we view the object.

```
virus
```

```
##
## This is a genclone object
## -----
## Genotype information:
##
##      752 multilocus genotypes
##     1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##      3 hierarchical levels - country year month
##     37 populations defined - Japan USA Finland ... Slovenia Croatia Algeria
```

If we need to store the number of MLGs as a variable, we can simply run the `mlg` command.

```
virus_mlg <- mlg(virus)

## #####
## # Number of Individuals: 1903
## # Number of MLG: 752
## #####

virus_mlg

## [1] 752
```

Since the number of individuals exceeds the number of multilocus genotypes, we conclude that this data set contains clones. Let's examine what populations these clones belong to.

3.2 MLGs across populations {clone-ing around}

Since you have the ability to define hierarchical levels 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` - Populations to include (Defaults to "ALL"). see [popsup](#).
- `blacklist` - Populations to exclude. see [popsup](#).
- `mlgsub` - see [mlg.table](#). 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 MLGs cross different populations and then give a vector that shows how many populations each one of those MLGs crosses.

```
setpop(virus) <- ~country
v.dup <- mlg.crosspop(virus, 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` to loop over the data with the function `length`.

```
head(v.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

v.num <- sapply(v.dup, length) # count the number of populations each MLG crosses.
head(v.num)

## MLG.3 MLG.9 MLG.31 MLG.75 MLG.80 MLG.86
##      2      4      2      8      2      2
```

3.3 Producing MLG tables and graphs {bringing something to the table}

We can also create a table of MLGs 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`

This function will produce a matrix containing counts of MLGs (columns) per population (rows). If there are not populations defined in your data set, a vector will be produced instead.

Default Command:

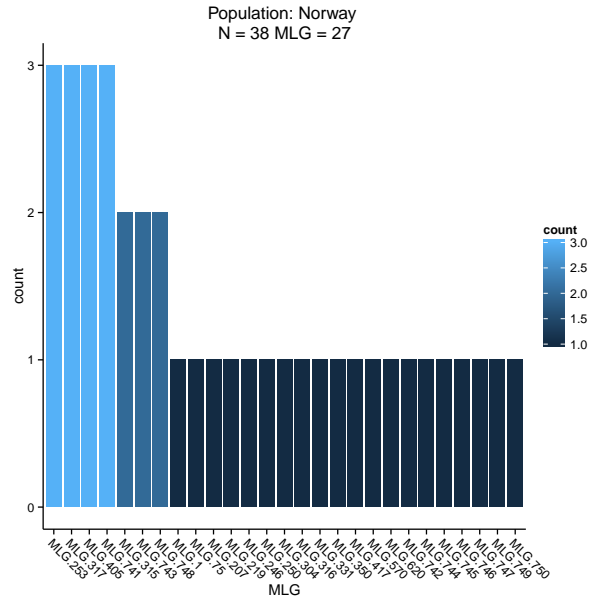
```
mlg.table(pop, sublist = "ALL", blacklist = NULL, mlgsub = NULL,  
          bar = TRUE, total = FALSE, quiet = FALSE)
```

- `pop` - a `genind` object.
 - `sublist` - Populations to include (Defaults to "ALL"). see [popsb](#).
 - `blacklist` - Populations to exclude. see [popsb](#).
 - `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` - When set to TRUE, the pooled data set will be added to the table. Defaults to FALSE.
 - `quiet` - Defaults to FALSE: population names will be printed to the console as they are processed.
-

```
v.tab <- mlg.table(virus, quiet=TRUE, bar=TRUE)  
v.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 4: An example of a bar-chart produced by `mlg.table`. Note that this data set would produce several such charts but only the chart for Norway is shown here.



The MLG table is not limited to use with *poppr*. In fact, one of the main advantages of `mlg.table` is that it allows easy access to diversity functions present in the package *vegan* [12]. One 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 this example, instead of drawing a curve for each of the 37 countries represented in this sample, let's set the hierarchical level to year.

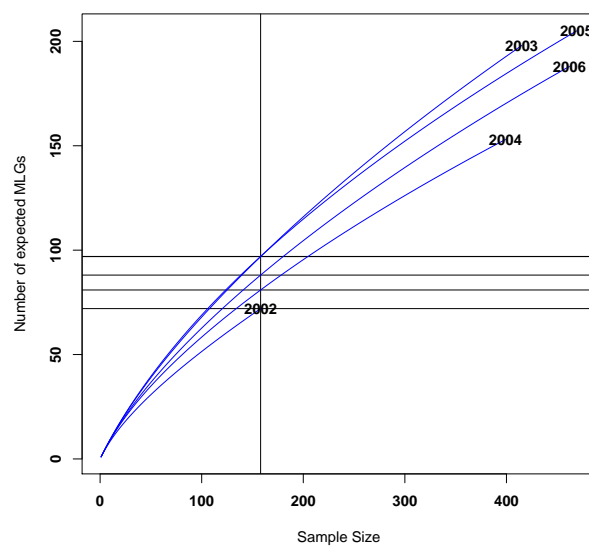
```
setpop(virus) <- ~year
summary(virus) # Check the data to make sure it's correct.
```

```
##
## # 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 2 3 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
```

```
library("vegan")
H.year <- mlg.table(virus, bar=FALSE)
rarecurve(H.year, ylab="Number of expected MLGs", sample=min(rowSums(H.year)),
          border = NA, fill = NA, font = 2, cex = 1, col = "blue")
```

Figure 5: 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 defined in the table. 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 Combining MLG functions {getting into the mix}

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:

```
setpop(virus) <- ~country
UGNN.list <- c("United Kingdom", "Germany", "Netherlands", "Norway")
UGNN <- mlg.crosspop(virus, 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, `v.tab`.

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

## [1] 315 317 620

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

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

v.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 investigate the incidence of these MLGs throughout our data set. One thing that the `GENCLONE` object keeps track of is a single vector defining the unique multilocus genotypes within the data. These are represented as integers and can be accessed with `mlg.vector`. This is useful for finding MLGs that correspond to certain individuals or populations. Let’s use `mlg.vector` to find individuals corresponding to the MLGs. First we’ll investigate what the output of this function looks like.

```
v.vec <- mlg.vector(virus)
str(v.vec) # Analyze the structure.

## int [1:1903] 605 605 672 675 674 673 670 671 670 678 ...
```

The integers produced are the MLG assignment of each individual in the same order as the data set. This means that the first two individuals have the exact same set of alleles at each locus, so they have the same MLG: 605. If we look at the number of unique integers in the vector, it corresponds to the number of observed multilocus genotypes:

```
length(unique(v.vec)) # count the number of MLGs

## [1] 752

virus # equal to the first number in this output.

##
## This is a genclone object
## -----
```

```
## Genotype information:
##
##      752 multilocus genotypes
##      1903 haploid individuals
##      125 codominant loci
##
## Population information:
##
##      3 hierarchical levels - country year month
##      37 populations defined - Japan USA Finland ... Slovenia Croatia Algeria
```

We will take UGNN (MLGs crossing UK, Germany, Netherlands, and Norway) and compare its elements to the MLG vector (`v.vec`) to see where else they occur.

```
UGNN # Show what we are looking for

## [1] 315 317 620

UGNN_match <- v.vec %in% UGNN
table(UGNN_match) # How many individuals matched to those three MLGs?

## UGNN_match
## FALSE TRUE
## 1881 22
```

22 individuals matched to those three MLGs. We can use this vector to show us the 22 individuals.

```
indNames(virus)[UGNN_match]

##      0329      0330      0331      0332      0341      0342      0345      0556
## "CY026119" "CY026120" "CY026121" "CY026122" "CY026131" "CY026132" "CY026135" "EU502462"
##      0557      0558      0870      0974      1112      1113      1114      1122
## "EU502463" "EU502464" "EU501513" "AB243868" "DQ883618" "DQ883619" "DQ883620" "DQ883628"
##      1193      1209      1210      1281      1288      1426
## "EU501609" "EU501642" "EU501643" "EU501735" "EU501742" "EU502513"
```

Note that there is an alternative way to list individuals matching specific MLGs using the function `mlg.id`. This function will return a list where each element represents a unique MLG. You can use this data to find out which individuals correspond to specific MLGs. Each element in the list is named with the MLG, but the index does not necessarily match up, so it is important to convert your query MLGs to strings:

```
virus.id <- mlg.id(virus)
virus.id[as.character(UGNN)]

## $`315`
## [1] "AB243868" "DQ883618" "DQ883628" "EU501642" "EU501643" "EU501735" "EU501742"
##
## $`317`
## [1] "DQ883619" "DQ883620" "EU501513" "EU501609" "EU502513"
##
## $`620`
## [1] "CY026119" "CY026120" "CY026121" "CY026122" "CY026131" "CY026132" "CY026135"
## [8] "EU502462" "EU502463" "EU502464"
```

We can also use the vector of MLGs to subset `mlg.table` with the `mlgsub` flag.

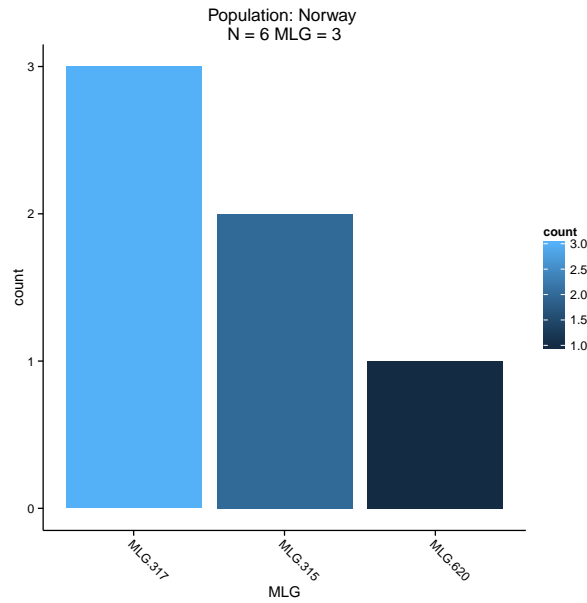
```
mlg.table(virus, mlgsub = UGNN, bar = TRUE)
```

```
##      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

That showed us exactly which populations these three MLGs came from in our data set.

Figure 6: An example of the same bar-chart as *Figure 1*, but focusing on three MLGs.



4 Appendix

4.1 General hierarchy method use

To reiterate, there are currently 5 methods that manipulate population hierarchies in genclone objects:
NOTE: Refer to [ACCESSING HIERARCHIES](#) for more details on how to access hierarchies.

Method	Function	Input	Result
split	splithierarchy	formula	defined hierarchical levels
set	sethierarchy	data frame	new hierarchy
get	gethierarchy	formula	data frame
name	namehierarchy	formula	new hierarchy names
add	addhierarchy	vector or data frame	new hierarchical level

These functions all have a syntax that looks like this:

```
newobject <- FUNCTION(object, input)
```

Let's take a data set of *Phytophthora infestans* collected from North America and South America and use that as an example. It has two population hierarchies defined, Continent and Country:

```
data(Pinf)
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
##    72 multilocus genotypes
##    86 tetraploid individuals
##    11 codominant loci
##
## Population information:
##
##    2 hierarchical levels - Continent Country
##    2 populations defined - South America North America
```

Let's say I wanted to change the hierarchy names to Spanish. I can do that using `namehierarchy`.

```
elPinf <- namehierarchy(Pinf, ~continente/pais) # Don't forget the formula syntax!
elPinf

##
## This is a genclone object
## -----
## Genotype information:
##
##    72 multilocus genotypes
##    86 tetraploid individuals
##    11 codominant loci
##
## Population information:
##
##    2 hierarchical levels - continente pais
##    2 populations defined - South America North America

Pinf
```

```
##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
##
##      2 hierarchical levels - Continent Country
##      2 populations defined - South America North America
```

The original data set has stayed the same and we now have a new data set with the names we want.

Of course, it would be silly to create a new data set every time we wanted to do something like change names. This is why all of the above functions (with the exception of `gethierarchy`) all have the replacement syntax of:

```
FUNCTION(object) <- input
```

NOTE: This is the preferred syntax.

This allows the object to be edited *in place* and makes things generally easier. Let's revisit our previous example:

```
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
##
##      2 hierarchical levels - Continent Country
##      2 populations defined - South America North America

namehierarchy(Pinf) <- ~continente/pais
Pinf

##
## This is a genclone object
## -----
## Genotype information:
##
##      72 multilocus genotypes
##      86 tetraploid individuals
##      11 codominant loci
##
## Population information:
##
##      2 hierarchical levels - continente pais
##      2 populations defined - South America North America
```

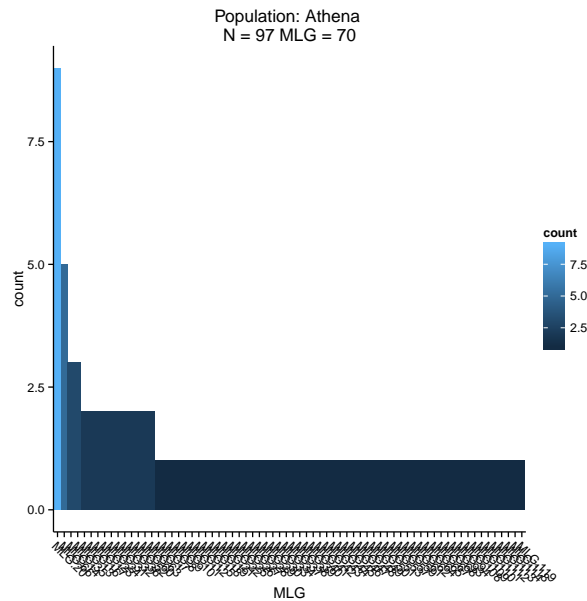
While we will mainly be using the replacement syntax in this vignette, the advantage to having both systems is that with the `function(object, input)` syntax, you can test manipulation without affecting your object.

4.2 Manipulating Graphics

Poppr utilizes *ggplot2* to produce many of its graphs. One advantage it gives the user is the ability to manipulate these graphs. With base R graphs, the only manipulation that can be performed is by adding elements to the graph. It is a static image. The *ggplot* graphs are actually represented as objects in your R environment. We can use the function `last_plot()` from *ggplot2* to be able to grab the plot that was plotted last in our window. Let's illustrate this using a MLG bar graph from the Athena population of the Aeut data set.

```
library("poppr")
library("ggplot2")
data(Aeut)
Athena.tab <- mlg.table(Aeut, sublist = "Athena")
```

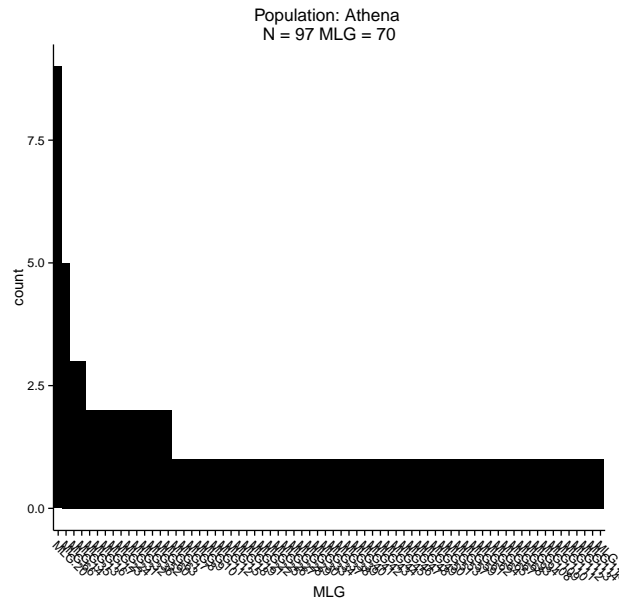
```
## | Athena
```



```
p <- last_plot()
```

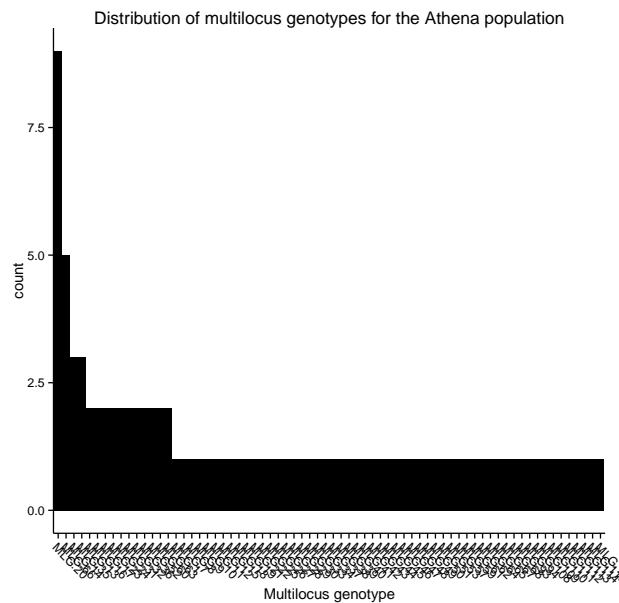
We've captured our plot using `last_plot()` and now we can manipulate it. Let's say we didn't like the scale going from blue to black and we wanted make it monotone without a color guide. (NOTE: the `()` around the call allows us to show the result immediately.)

```
(pb <- p + scale_fill_continuous(low = "black", high = "black", guide = "none"))
```



We could also change the title

```
(pbt <- pb + ggtitle("Distribution of multilocus genotypes for the Athena population")) + xlab("Multilocus genotype")
```



This allows you to produce publication quality graphs directly in R. Please see Hadley Wickham's *ggplot2* package for more details [19]. Note that if you don't like using *ggplot2*, you can access the data in the *ggplot2* object and plot the data yourself:

```
head(p$data)
```

```
##      MLG count
## 15 MLG.20     9
## 57 MLG.66     5
##  9 MLG.14     3
```

```
## 30 MLG.35      3
## 8  MLG.13      2
## 11 MLG.16      2
```

4.3 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/>.

4.3.1 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.

4.3.2 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>

4.3.3 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.

4.3.4 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$other$x$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.

4.4 Table of Functions

Below is a table of functions found in *poppr*. These functions are linked within the document. If a function name is blue, simply click on it to go to its definition and description.

Table 10: Functions available in *poppr*

Function	Description
IMPORT/EXPORT	
<code>getfile</code>	Provides a quick GUI to grab files for import
<code>read.genalex</code>	Read <i>GenAlEx</i> formatted csv files to a genind object
<code>genind2genalex</code>	Converts genind objects to <i>GenAlEx</i> formatted csv files
<code>as.genclone</code>	Converts genind objects to genclone objects
MANIPULATION	
<code>setpop</code>	Set the population using defined hierarchies
<code>splithierarchy</code>	Split a concatenated hierarchy imported as a population
<code>sethierarchy</code>	Define a population hierarchy of a genclone object
<code>gethierarchy</code>	Extract the hierarchy data frame
<code>addhierarchy</code>	Add a vector or data frame to an existing hierarchy
<code>namehierarchy</code>	Rename a population hierarchy
<code>missingno</code>	Handles missing data
<code>clonecorrect</code>	Clone censors at a specified population hierarchy
<code>informloci</code>	Detects and removes phylogenetically uninformative loci
<code>popsub</code>	Subsets genind objects by population
<code>shufflepop</code>	Shuffles genotypes at each locus using four different shuffling algorithms
<code>splitcombine*</code>	Manipulates population hierarchy *Deprecated
<code>recode.polyploids</code>	recode polyploid data sets with missing alleles imported as “0”
DISTANCES	
<code>bruvo.dist</code>	Bruvo’s distance
<code>diss.dist</code>	Absolute genetic distance (see <code>provesti.dist</code>)
<code>nei.dist</code>	Nei’s 1978 genetic distance
<code>rogers.dist</code>	Rogers’ euclidean distance
<code>reynolds.dist</code>	Reynolds’ coancestry distance
<code>edwards.dist</code>	Edwards’ angular distance
<code>provesti.dist</code>	Provesti’s absolute genetic distance
BOOTSTRAPPING	
<code>aboot</code>	Creates a bootstrapped dendrogram for any distance measure
<code>bruvo.boot</code>	Produces dendrograms with bootstrap support based on Bruvo’s distance
ANALYSIS	
<code>poppr.amova</code>	Analysis of Molecular Variance (as implemented in ade4)
<code>ia</code>	Calculates the index of association
<code>mlg</code>	Calculates the number of multilocus genotypes
<code>mlg.crosspop</code>	Finds all multilocus genotypes that cross populations
<code>mlg.table</code>	Returns a table of populations by multilocus genotypes
<code>mlg.vector</code>	Returns a vector of a numeric multilocus genotype assignment for each individual
<code>mlg.id</code>	Identifies individuals associated with each MLG
<code>poppr</code>	Returns a diversity table by population
<code>poppr.all</code>	Returns a diversity table by population for all compatible files specified
<code>private.alleles</code>	Tabulates the occurrences of alleles that only occur in one population
<code>locus.table</code>	Creates a table of summary statistics per locus

Continued on next page...

Table 10 – continued from previous page

Function	Description
VISUALIZATION	
<code>plot_poppr_msn</code>	Plots minimum spanning networks produced in poppr with scale bar and legend
<code>greycurve</code>	Helper to determine the appropriate parameters for adjusting the grey level for msn functions
<code>bruvo.msn</code>	Produces minimum spanning networks based off Bruvo’s distance colored by population
<code>poppr.msn</code>	Produces a minimum spanning network for any pairwise distance matrix related to the data
<code>info_table</code>	Creates a heatmap representing missing data or observed ploidy
<code>genotype_curve</code>	Creates a series of boxplots demonstrating how many loci are needed to represent the diversity of your data.

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] Erica M. Goss, Javier F. Tabima, David E. L. Cooke, Silvia Restrepo, William E. Fry, Gregory A. Forbes, Valerie J. Fieland, Martha Cardenas, and Niklaus J. Grünwald. The irish potato famine pathogen *phytophthora infestans* originated in central mexico rather than the andes. *Proceedings of the National Academy of Sciences*, 111(24):8791–8796, 2014.
- [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] 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.
- [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.
- [19] Hadley Wickham. *ggplot2: elegant graphics for data analysis*. Springer New York, 2009.