

Introduction to R: practical

Lewis Spurgin

May 2017

Introduction

The aim of this practical is to get you comfortable using R and R studio. Some tips:

- Write all your commands in a SCRIPT and run them from there.
- Annotate your script well (using #). No amount of annotation is too much if it helps you understand what you're doing.
- Use R's help (e.g. `?mean`) if you are not sure what a function does
- Google is your friend

The data

Instructions

Our aim is to calculate the proportion of heterozygotes in each of the populations in our dataset, and to plot these proportions using a barplot. *If you think you can do this without guidance, please go ahead and do so.* Otherwise, step-bystep instructions are provided below.

Step 1. Set the working directory, read in and explore the data.

Set your working directory to the “Lectures/Intro_to_R” folder

```
setwd("path/PopGenBerlin/Lectures/Intro_to_R")
```

Now you can read in your data, and assign it to an object name. We will use the name `dd`.

```
dd <- read.csv("Data/Microsat_data.csv")
```

In R studio, click on the “environment” tab in the top right, if it is not already open. You will notice that `dd` is there, along with the description. This, as you may have guessed, tells us that our data frame has 250 rows and 4 columns. Click on the table icon to the right of this pane. You should now be able to see your data. You will see that the four variables are called `Ind_ID`, `Pop`, `Allele_1`, `Allele_2`.

We can also explore our data by typing and running some functions into our script. We will use the functions `head()` and `str()`.

```
head(dd)
```

```
##   Ind_ID      Pop Allele_1 Allele_2
## 1 Ind_1 Mainland    102     102
## 2 Ind_2 Mainland    100     100
## 3 Ind_3 Mainland    102     102
## 4 Ind_4 Mainland    100     100
## 5 Ind_5 Mainland    100     102
## 6 Ind_6 Mainland    100     102
```

```
str(dd)
```

```
## 'data.frame': 250 obs. of 4 variables:
## $ Ind_ID : Factor w/ 250 levels "Ind_1","Ind_10",...: 1 112 174 185 196 207 218 229 240 2 ...
## $ Pop : Factor w/ 2 levels "Island","Mainland": 2 2 2 2 2 2 2 2 2 2 ...
## $ Allele_1: int 102 100 102 100 100 100 100 102 102 100 ...
## $ Allele_2: int 102 100 102 100 102 102 102 102 102 100 ...
```

The function `head()` can be very useful when you want a quick look data frames with many rows. But `str()` gives us all the same information and more. From the output above, you should be able to answer the following:

- How many individuals are there in your dataset?
- How many populations are there?
- How has R stored the microsatellite allele data?

Step two - define our genotypes

In our dataset we have separate information for each of the two alleles found in each individual. For example, we can see from when we ran the `head()` function that the first individual in our dataset has an allele of length 102 for its first allele, and an allele of length 102 for its second allele. This individual is therefore a homozygote at this microsatellite locus. We can tell R to give use this information to tell us about heterozygosity. But first, we need to get an understanding of logical operators in R. Try typing the following commands directly into your terminal (I know we said not to do this, but it's ok here!):

```
4 < 3
4 > 3

4 = 3
4 == 3 #gives us an error. Why?

x = 4
x == 4

"cat" == "dog"

x <- "cat"
x == "cat"
```

Hopefully you have figured out that a single equals sign and double equals sign mean very different things to R. A single equals sign generally works in the same way as the `<-` operator (with a few subtle differences), whereas a double equals sign asks a logical question, and returns a value of `TRUE` or `FALSE`.

So coming back to our microsatellite data, we can use the `==` operator to ask whether an individual is a homozygote, by asking whether the variables `Allele_1` and `Allele_2` are the same.

```
# Get the data for the first individual and have a look at it
ind1 <- dd[1,]
ind1
```

```
## Ind_ID Pop Allele_1 Allele_2
## 1 Ind_1 Mainland 102 102
```

```
# Now ask whether Allele 1 and Allele 2 are identical
ind1$Allele_1 == ind1$Allele_2
```

```
## [1] TRUE
```

Because R works so nicely with vectors, we can do this for every individual in our dataset very simply, as follows:

```
dd$Allele_1 == dd$Allele_2
```

```
## [1] TRUE TRUE TRUE TRUE FALSE FALSE FALSE TRUE TRUE TRUE FALSE
## [12] FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
## [23] TRUE TRUE TRUE FALSE TRUE TRUE FALSE FALSE TRUE FALSE TRUE TRUE
## [34] TRUE FALSE FALSE TRUE TRUE TRUE FALSE TRUE FALSE TRUE TRUE
## [45] FALSE FALSE FALSE FALSE FALSE TRUE TRUE TRUE FALSE FALSE TRUE
## [56] TRUE TRUE FALSE FALSE FALSE TRUE FALSE TRUE TRUE TRUE FALSE
## [67] TRUE FALSE TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE TRUE TRUE
## [78] TRUE TRUE FALSE TRUE TRUE FALSE FALSE FALSE FALSE FALSE TRUE TRUE
## [89] TRUE FALSE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE FALSE
## [100] FALSE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE FALSE
## [111] TRUE FALSE TRUE FALSE TRUE TRUE TRUE TRUE TRUE FALSE FALSE TRUE
## [122] TRUE FALSE TRUE TRUE TRUE TRUE FALSE FALSE TRUE FALSE TRUE
## [133] TRUE TRUE TRUE TRUE TRUE FALSE FALSE TRUE TRUE TRUE FALSE
## [144] TRUE FALSE TRUE TRUE FALSE TRUE FALSE TRUE TRUE TRUE FALSE
## [155] FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE FALSE TRUE TRUE
## [166] TRUE FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE
## [177] TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE TRUE TRUE
## [188] FALSE TRUE TRUE TRUE FALSE TRUE TRUE FALSE TRUE TRUE TRUE
## [199] TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [210] TRUE TRUE TRUE TRUE TRUE FALSE TRUE TRUE TRUE TRUE TRUE
## [221] TRUE TRUE TRUE TRUE FALSE TRUE FALSE TRUE TRUE TRUE TRUE
## [232] TRUE TRUE TRUE TRUE TRUE FALSE FALSE TRUE TRUE FALSE FALSE
## [243] FALSE TRUE TRUE TRUE FALSE TRUE TRUE TRUE
```

This gives us 250 TRUE or FALSE values, and tells us whether each individual in our dataset is a homozygote or not. Now all we need to do is turn this into a column of our data frame, indicating whether each individual is a homozygote or a heterozygote. We can do this using the `ifelse()` function. Check out what this function does using `?ifelse`

```
#Add a Heterozygosity column
dd$Het <- ifelse(dd$Allele_1 == dd$Allele_2,"Homozygote","Heterozygote")

#Check that it's worked
head(dd)
```

```
## Ind_ID Pop Allele_1 Allele_2 Het
## 1 Ind_1 Mainland 102 102 Homozygote
## 2 Ind_2 Mainland 100 100 Homozygote
## 3 Ind_3 Mainland 102 102 Homozygote
## 4 Ind_4 Mainland 100 100 Homozygote
## 5 Ind_5 Mainland 100 102 Heterozygote
## 6 Ind_6 Mainland 100 102 Heterozygote
```

Great - we can see that the homozygotes and heterzygotes are being correctly identified. Now we can do some calculations and make a plot.

Step three - calculate and plot genotype frequencies

Now we want to calculate the proportion of heterozygotes in each of our two populations. To do this, we need to know i) the number of heterozygotes in each population and ii) the total number of individuals in each population. We can get both of these pieces of information using the function `table()`.

```
table(dd$Het) #Number of each genotype
```

```
##  
## Heterozygote   Homozygote  
##           85           165
```

```
table(dd$Pop) #Number of individuals in each population
```

```
##  
##   Island Mainland  
##     93      157
```

```
table(dd$Het,dd$Pop) #Both combined
```

```
##  
##           Island Mainland  
## Heterozygote     19      66  
## Homozygote      74      91
```

We can assign the output of the `table()` function to an object. Then using this we can extract subsets of data and perform calculations using them. Using this approach, we can divide the number of heterozygotes by the sample size to calculate observed heterozygosity.

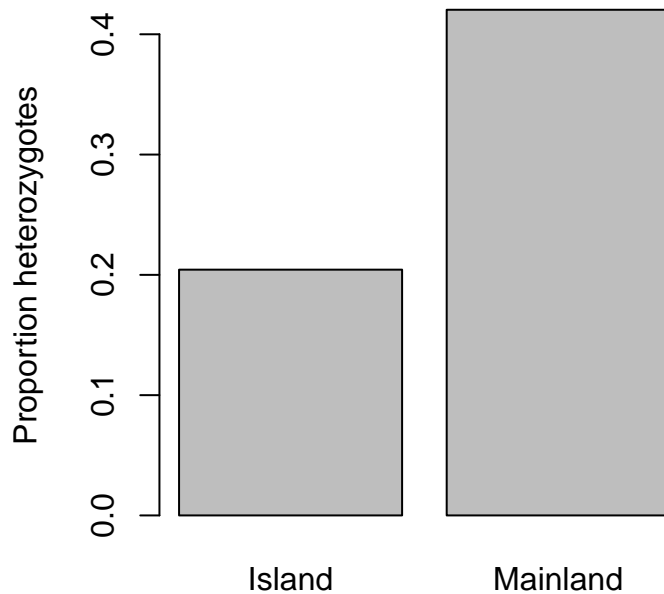
```
#Get the number of heterozygotes and the sample size for each population  
het_counts <- table(dd$Het,dd$Pop)["Heterozygote",] #pulls out the heterozygote row  
samplesizes <- table(dd$Pop) #Calculate sample sizes
```

```
#Now divide the two together  
het_freqs <- het_counts/samplesizes  
het_freqs
```

```
##  
##   Island Mainland  
## 0.2043011 0.4203822
```

Now we can plot the frequencies. We will get into plotting in much more detail later in the course. For now, we will use the function `barplot()`.

```
barplot(het_freqs,ylab = "Proportion heterozygotes")
```



Finally, we can test whether the proportion of heterozygotes differs between mainland and island population using a Chi-squared test (`chisq.test`).

```
chisq.test(dd$Het, dd$Pop)
```

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: dd$Het and dd$Pop
## X-squared = 11.208, df = 1, p-value = 0.0008143
```

Some challenges, if you have finished

- Calculate and plot the frequency of the two alleles in each population.
- Calculate expected genotype frequencies for each population, according to Hardy-Weinberg. Compare your observed and expected frequencies using chi-squared tests.
- I generated the microsatellite dataset for this practical in R. However, this dataset is very poorly annotated. Go through the script “Microsat_data.R”, figure out what is being done with each line of code, and annotate the script accordingly.