

Introduction to R

QFAB Bioinformatics





Queensland Cyber Infrastructure Foundation Ltd, ABN 13 225 133 729, Axon building, 47 – The University of Queensland - St Lucia, Qld, 4072
(QCIF incorporates QFAB Bioinformatics). 16.01.023-20170117

Contents

Abstract	5
1 Getting started	7
1.1 How to read this book	7
1.2 RStudio IDE	8
1.3 Alternatives to RStudio	9
1.4 The R workspace	9
1.5 Literate programming	9
2 Data variables (objects) in R	13
2.1 Declaring variables	13
2.2 Datatypes	14
2.3 Functions	17
2.4 Data structures	18
2.5 Finding, describing and removing objects	26
2.6 Best practise and variable naming conventions	26
3 Going further with R	27
3.1 Documentation	27
3.2 Extending R with packages	27
3.3 Vignettes	28
3.4 SessionInfo	29
3.5 Other ways of running R	30
4 Importing and exporting data from R	33
4.1 Saving session data	33
4.2 Current working directory	34
4.3 Reading and writing tabular data	35
4.4 Scan	39
4.5 Reading from a web connection	39
5 Transforming and summarising data in R	41
5.1 User defined functions	41
5.2 For loops, While loops and Vectorisation	43
5.3 Apply, Lapply	44
5.4 Basic data analysis methods	45
6 Data presentation and reporting	47
6.1 Boxplots	47
6.2 Histograms	49
6.3 Bar charts	51
6.4 Scatterplots	54
Appendix 1	57
Summary of Object Types	57
More information on R and Bioconductor for Genomics	57

Appendix 2 - Solutions to exercises	59
2.4.1 Using vectors	59
2.4.3 Using lists	59
2.4.4 Adding factor values	59
2.4.6 Creating a data frame	59
4.2 Reading in the GeneAtlas data	60
5.1 Creating and using functions	60
5.2 Loops and vectorisation	60
6.5 Venn diagram	61
Appendix 3 - The SessionInfo for this version of	63

Abstract

R is a open-sourced software programming language and software environment for statistical computing and graphics. The R language is widely used among statisticians and data analysts for developing statistical workflows for data analysis. R is an implementation of the S programming language combined with lexical scoping semantics inspired by Scheme. R was created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand, and is currently developed by the R Development Core Team.

Bioconductor is a project to develop innovative software tools for use in computational biology. It is based on the R language. Bioconductor packages provide flexible interactive tools for carrying out a number of different computational tasks. Literate programming and analytical pipeline crafting.

R and Bioconductor may not be as fast as dedicated bioinformatics software for computationally intensive processes such as mapping short sequence reads onto the genome, but the flexibility of having raw access to all of the data, methods and structure is empowering.

CONTENTS

Chapter 1

Getting started

R is a statistical environment and programming language for data analysis and graphical display. The software is open-source, freely available and has been compiled ready for use on Windows, Mac and Linux computers. The Bioconductor framework has a considerable number of packages that have been implemented for the analysis and exploration of metabolomics, proteomics, DNA microarray and Next Generation DNA sequence data.

This workshop is intended as an introduction to R that should familiarise you with the concepts that underpin the crafting of workflows in R and some of the key techniques that will aid in the crafting of reusable workflows.

Some of the topics we will cover over the course of the day include:

- Using the RStudio IDE for running a data analysis in R
- The basics of literate programming for studies using R
- Loading, creating, modifying and saving basic -omics data objects according to key formats
- Create graphs and plots of data
- Understanding how to read and understand someone else's R-code

1.1 How to read this book

In the shaded dialog boxes with the thick left border are commands that should be typed into your R *Console* window. This corresponds to **input**, the **output** will sometimes follow below with lines preceded with two hashmarks (##).

```
| print("HELLO WORLD")
## [1] "HELLO WORLD"
```

Tips, Suggestions and Traps

Comments and salient advice are provided in dialogs such as this. We'll use this format to warn you of traps that you could slip into and to provide some hints.

Warnings and Traps

Warnings are provided in dialogs such as this. We'll use this format to warn you of traps that you could slip into and to provide some hints.



Time for you to make R work

Exercises are the best way to learn. A document can provide an insight to the process but hands-on interaction with R is the only way to learn. In the exercise boxes are some suggested exercises that will apply the knowledge that is being shared.

1.2 RStudio IDE

RStudio is a free and open source integrated development environment (IDE) for R. An IDE is most typically used in software development where a user is presented with a unified environment where code can be edited, documentation can be read and debuggers can be applied to understand what is happening in the code and why it may not be working as expected.

In bioinformatics (or computational biology) an IDE provides us with an integrated approach to writing R scripts (or markdown documents), interacting with our data as we implement our script and provides us with an overview of the objects that we have created and their content. RStudio also retains a history of the commands that we have typed, a collection of the figures that we have prepared and provides access to method documentation and package vignettes. We will consider all of these aspects during this workshop. The RStudio IDE is great for the preparation of scripts and reports since the code syntax-highlighting helps you discover typos and errors in your code and even spelling mistakes.

RStudio comes in two main flavours: a *Desktop* version that runs straight off the computer that we are sitting at and a *Server* version that can be hosted on a slightly more capable computer in a data centre. The server version is sometimes preferable since it is accessed via a web-browser and can be configured with massive memory and disk allocations that might not be practical for a normal laptop. The server version is platform-agnostic (from the user perspective that is; the server itself needs to run Linux) and the Desktop version can be run on Windows, Mac and Linux computers. An example RStudio session running on the Desktop version is shown in Figure 1.1.

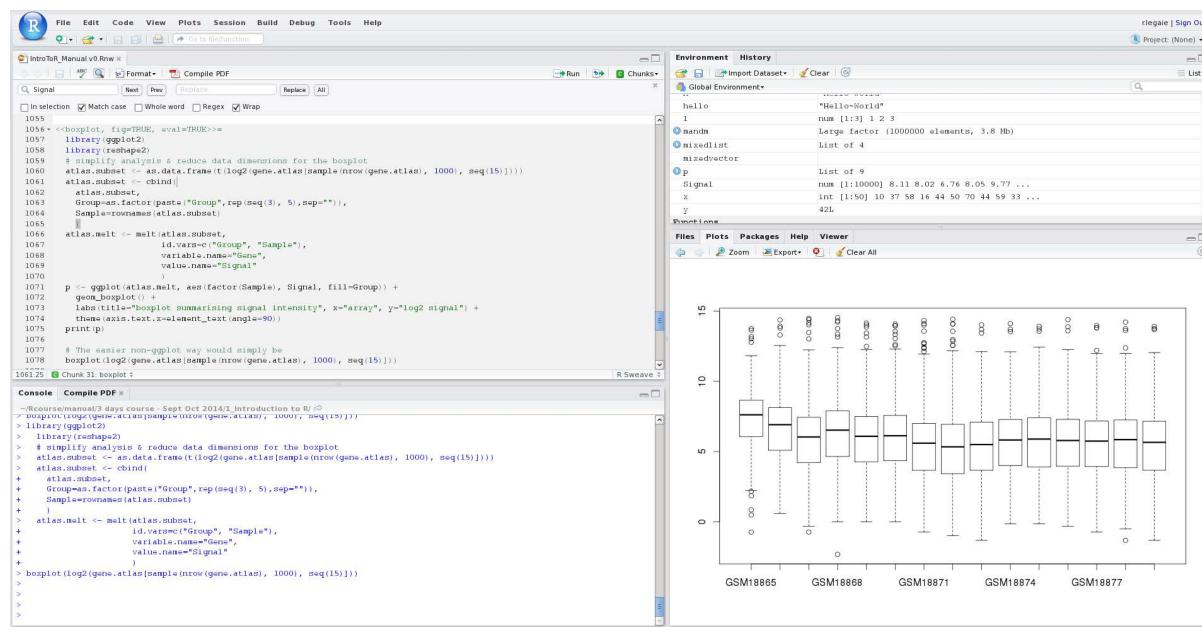


Figure 1.1: A screenshot of the RStudio session used in the preparation of this course material. The document being prepared is displayed in the top left pane. A figure is shown in the bottom right pane and the loaded environment variables are presented in the top right. Having access to all of these information in a single application is simpler than managing an active R console, a text editor and an image viewer.

RStudio integrates very cleanly with a number of best scientific working practices (data sharing, reproducible research and research documentation). RStudio allows for experiments to be performed in

projects - this allows you to create a separate workspace for each study that you are performing. Code can be shared between projects easily and Version Control (Subversion and Github) facilitate the sharing of workflows and methods with other researchers.

For this course we recommend that you use our server version that has been preconfigured for the software, packages and data that will be used.



Open up an RStudio server client and create a project

1. Open a web browser (not Internet Explorer) and connect to the server given to you by the trainer. Use the username and password that we have provided to connect and explore the options available in the top bar.
2. Explore the interface
 - a) Where is the **Files** tab? What does it show? Can you see the **data** folder?
 - b) Click the **data** folder to look inside, keep going until you get to some files
3. Click on the **Packages** tab and look at the list of packages that are installed in R. If you are ahead, scroll down and click on the *methods* name. This should bring you to the **Help** tab, see what additional functions you can do using this package.
 - We will talk about what packages are later in the workshop.
 - a) In the top-right corner of the **Help** tab is a search text field. Type in “*print*” and hit enter (if the autofinish starts to pop up click on print). This will show you the **print()** function and what parameters it accepts.

1.3 Alternatives to RStudio

RStudio is not the only way to access R. The standard R installation from CRAN (Comprehensive R Archive Network) provides a basic R graphical user interface called the **R console**. This provides access to the key R functionality and provides interfaces for the preparation of R scripts. The simplest way to access R is perhaps through the command line. With a typical R installation simply typing "R" at the console should load an R session which you can directly interact with, as seen in Figure 1.2.

1.4 The R workspace

The **workspace** refers to a R working environment and the collection of data objects that have been created by the user. A **session** refers to an instance of a workspace. In a session you populate your workspace with a collection of objects and functions. A session can be saved, this means that you can save the workspace content so that the next time that you use R you can have all of the information loaded and already available.

Session and workspace concepts are best understood within an environment such as RStudio.

1.5 Literate programming

There are two approaches to use R to perform data analysis:

- The first approach is done interactively via the **console** window. You can tell you are in this window by the > prompt and the blinking cursor. When you are in this window, all commands are executed *immediately* as soon as you hit the **[Enter]** key.
- The second approach is by writing scripts (or text file) and then running the script from top to bottom. This is done via the **editor** window (not visible by default, only when you open a script file.) The commands entered in this screen are *not* executed immediately. A script consists of a

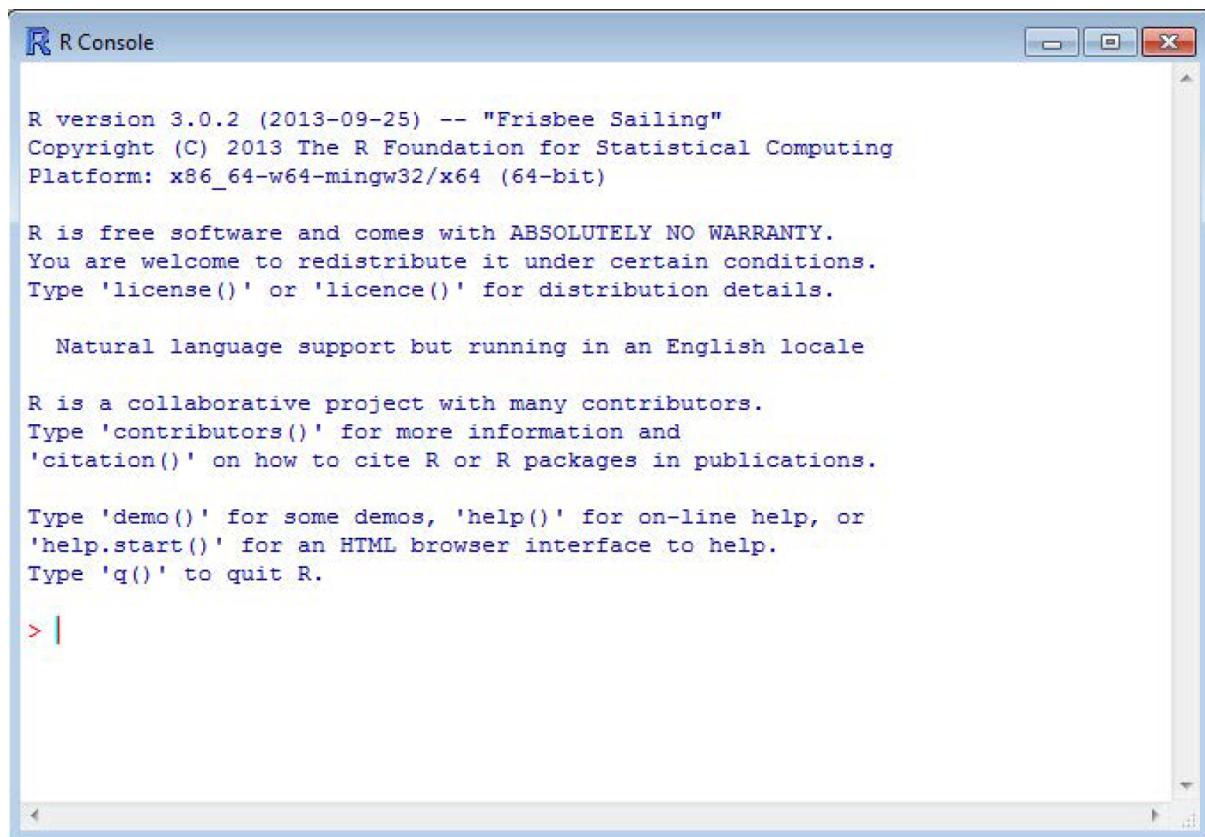


Figure 1.2: R installations typically include a basic R editor and console that provide a basic framework for interacting with the software, preparing R scripts and for the installation and management of packages. This figure shows the R console as installed by default on a Mac computer.

collection of commands that perform a larger task on the data and is generally executed from top to bottom. A quick example will be shown in Section 3.5.1.

The preferred mode of operation is up to the user. There are pros and cons using either method. The interactive mode is quick to start and you get immediate feedback about your data. However, it does not support repeated analysis very well as you would constantly have to retype the same lines of code again and again. The scripting method is very useful for repeated analysis but can be slow to set up as mistakes in the command can become hard to debug for large scripts. With practice you will become more familiar (and faster) and will develop your own style of working.

When using RStudio, it allows you a third mode of operation which is a hybrid of the two approaches mentioned above. That is, you write your commands in the editor window and run sections of the code (or **chunks**) as required. So the analysis is done semi-interactively.



For this workshop, we recommend that you use the script method to work through the exercises. That is, create a new text file for each section and execute the lines one at a time so you can see what happens with each command. This will allow you to come back to the exercises in your own time and review the content.



Literate programming

1. Try some basic functions in the **Console** window (the large panel to the left)
 - a) Type `1 + 1` and hit [ENTER], what is the output and where is it shown?
 - b) Type `print("HELLO WORLD")` command in the console and hit [ENTER].
 - c) Try some other functions, e.g.
 - `3^2`, which is three squared
 - `11%%2`, which is the modulo function (the remainder left after the first number is divided by the second).
2. Go to **File** and select **New File > R Script** and type the same commands in the **editor** window.
 - a) Highlight *only* line 1 and click on the **Run** button at the top right of the editor window (or hit **[Ctrl]+[Enter]** on the keyboard)
 - b) Where is the output shown?
 - c) Highlight all three lines and click on the **Run** or **[Ctrl]+[Enter]**
3. Create a New Project that we will use for this course.
 - a) Under the **File** menu, select **New Project**
 - b) Click on **New Directory**, then **Empty Project**
 - c) Select a name for your project. Generally this should reflect the analysis that you will be performing. Remember, you may want to come back to your projects months or even years in the future.
 - d) You should now see an `.Rproj` file in your **Files** tab, with the name of the project you have just chosen. This file will contain your R session.

Below are some helpful tips when using RStudio to help you speed up your analysis.



The up arrow key

When typing in the **console** window, you can recall previous commands by hitting the **[up arrow]** key. This will cycle through the history of your previous commands starting with the most recent. Use the **[up arrow]** and **[down arrow]** keys to move through the commands. Hit the **[esc]** key to escape and return to the `>` prompt.



The history window

The history window on the top right hand corner (next to Environment) holds a history of all the commands you have executed in the past. This is a faster way of navigating your previous commands to rerun them. But this window is only available in these types of development environment software like RStudio. You even:

- * search your history using the search box
- * highlight the desired lines and click on either the **To Console** or **To Source** button

Chapter 2

Data variables (objects) in R

As with all other computer languages information needs to be stored in a way that can be recalled, displayed and analysed. A **variable name** is used to point to a data object in memory (or on disk).

The conventions for a variable name:

- should not start with a number,
- should not contain spaces,
- should not be the same as function names as this can be ambiguous,
- should not contain special characters such as #, %, &. These characters require special handling and often lead to errors that require debugging.

2.1 Declaring variables

Declaring a variable (or object) is simple in R, we just need to give it a name and the content for that variable. This can be done either directly in the console or in the editor window and saved as a script.

Remember that code entered into the console window will be executed immediately when you press the [Enter] key. But do not expect output for every line you enter, many R commands do not print anything to the screen. Code in the editor document will only be executed when you run the lines of commands or run the entire script.



Declare an variable

1. Enter the example code below into your RStudio interface. You can type it directly into the console and pressing the **[Enter]** key after each line or in the editor window.
 - a) Which method are you using? console or editor?
 - b) Did you get any output after executing line [1]?
2. If you entered above code in the console, click on the **Environment** tab in the top-right corner window. Do you see the value x listed? If you do not then try again (or select the `x <- 42` line in your editor window and click on **Run**).
 - a) Try other commands, e.g. `y <- 555`, do you see a new variable y listed in the *Environment* tab?

```
x <- 42  
# this has set x to the value of 42  
x  
## [1] 42
```

```
| print(x)
## [1] 42
```

Setting variable values and viewing variable contents

Command	Description
<-	This marker is used to set a variable's value. The more standard = operator will also work but they have different levels of precedence. As best practice it is better to use <- in R.
#	The hashmark at the start of a line is a way to comment your code. This directs R to ignore the rest of the current line. Commenting your good is very good practice and is also used lots to remove commands that you may use in debugging a more complex workflow.
print()	Is the print function which writes out the value of the variable. Simply typing the name of the variable will also write out the value contained within, as seen in the examples above: x and print(x)

2.2 Datatypes

In common with many programming languages, every variable in R is assigned a **data type**. The primary data types in R are **Numeric**, **Integer**, **Character**, **Logical**, and **Complex** (which is used for specialist mathematical applications). The type of a variable determines what data can be stored within it and often, what functions can be applied to it.

2.2.1 Numeric

Decimal values are called **numerics** in R. In the previous section we declared a variable called x, to which we assigned the value of 42. R recognised this as a number and automatically set the type of x to be numeric. We can test whether a variable is a numeric using the **is.numeric()** function. Or, we can ask directly what type it is using **class()**.

```
| is.numeric(x)
## [1] TRUE
```

```
| class(x)
## [1] "numeric"
```

2.2.2 Integer

Integers are a separate type to numerics and contain *whole number* values only. Just as we can test whether a variable is a numeric, we can also test whether one is an integer:

```
| is.integer(x)
## [1] FALSE
```

Strangely, although 42 is indeed an integer, this reports FALSE. This is because x has been **cast** as a numeric and a variable can only be of a single type. Indeed even **is.integer(42)** returns FALSE because R by default, assigns any numerical value as type numeric.

2.2.2.1 Typecasting

If we want to perform integer functions on `x`, we can force it from one type to another, in this case using the `as.integer()` function. This is known as **typecasting**. The `as.integer()` function also forces non-whole numbers into integers by rounding them down to the nearest whole number.

```
y <- as.integer(x)
is.integer(y)
## [1] TRUE

class(y)
## [1] "integer"

as.integer(3.1415)
## [1] 3
```



Writing to variables or the screen

- In the example above, you will notice that the line `y <- as.integer(x)` generates no output, while `as.integer(3.1415)` prints a result to screen. This is because in the first case, the output of `as.integer()` is passed into the variable `y`, while in the second case, no destination is given so by default it prints to screen.
- This also means that the output from the first command can be accessed and used again later by calling `y` but the output from the second is lost and cannot be used by future R commands.

2.2.3 Logical

Logical datatypes can have one of only two values: `TRUE` or `FALSE`. **Note**, these are case-specific, `True` or `true` are not valid logical values. However, the shortcut: `T` and '`F`' will work.

```
TRUE
## [1] TRUE

TRUE == TRUE
## [1] TRUE

TRUE == FALSE
## [1] FALSE

TRUE != FALSE
## [1] TRUE

class(FALSE)
## [1] "logical"

# Logicals can also be typecast
as.integer(TRUE)
## [1] 1
```



** Simple datatypes **

1. Now try typecasting FALSE to an integer.
2. Next try typecasting in the reverse, that is, converting the digits (1, 0) into TRUE or FALSE. *Hint: as.logical()*
3. [Optional] If you are ahead, have a look at typecasting other values (e.g. >1 and <0) into logical types. Can you see what is happening?

2.2.4 Character

A **character** variable is used to represent strings or text values in R. We convert variables into character values with the `as.character()` function:

```
| h <- "Hello world"
| class(h)
|
## [1] "character"
|
| as.character(x)
|
## [1] "42"
|
| as.character(FALSE)
|
## [1] "FALSE"
```



Setting, retrieving and changing data types

1. Enter the code blocks from the sections above (from Numerics through to Character) in RStudio (as before, either into the console or the editor window), and run them to generate the output.
 - Remember, you can use the *Environments* tab to keep track of the variables being created.
 - **Note**, below all the outputs are shown altogether for readability. If you are running them in console, you will see the output immediately as you hit [Enter].
2. Experiment with different values so you understand better how the types and typecasting behave.
3. Characters are much more interesting than just storing a piece of text. Try the commands below, which explore a number of fundamental character manipulating functions. The functions are explained in the table below.

```

letters <- c("A", "B", "C")
paste(letters, collapse=" ")
hello <- paste("Hello", "World", sep="~")
hello
nchar(hello)
substr(hello, 3, 7)
strsplit(hello, '~')

gsub("World", "R", hello)
grep("Hello", hello)
grep("Mouse", hello)
grepl("Stephen", hello)

## [1] "A B C"
## [1] "Hello~World"
## [1] 11
## [1] "llo~W"
## [[1]]
## [1] "Hello" "World"
##
## [1] "Hello~R"
## [1] 1
## integer(0)
## [1] FALSE

```

Characters, strings and regular expressions

Function	Description
c()	means <i>combined</i> , is used to join a number of single values separated by a comma together into a data structure. We will touch on this in the next section (2.4)
paste(variable1, variable2, ..., sep)	concatenate two or more character variables into a single variable. The concatenation can be performed on individual specified character variables or complex variables (see next section). The <code>sep</code> argument, short for separate, is used to specify the delimiter that will be used to link the variables together.
nchar(variable)	returns an integer value that corresponds to the number of characters contained within the variable.
substr(variable, start, stop)	to prepare a substring from a character variable using the start and stop coordinates of the substring to retain.
strsplit(variable, delimiter)	character variable into a list of character variables based on a delimiter. If an empty delimiter argument is provided the individual characters will be returned.
gsub(pattern, replacement, variable)	is used to substitute a substring in the given character variable with a different (or empty) substring.
grep(pattern, variable)	searches for the pattern in the character variable and returns 1 if found or <code>integer(0)</code> if not found; while <code>grepl()</code> provides a logical response as to whether the pattern exists.

2.3 Functions

Before we move onto data structures in the next section, we will take an aside to look at **functions**. Functions are a group of commands that perform a specific action. By itself it is not a complete executable program but form part of a larger workflow. Like variables that hold data, functions are also given a name so that they can be *reused*. They generally follow the `input -> process -> output` model, that is, they take input, do something with the input and produce output. The output can either be printed

to the screen or be *returned* and assigned to a variable that holds the new data for further reuse.

So far we have already started using some of the in-built functions in R like `print()`, `str()`, `as.integer()` etc. They take the form of `output_variable <- function_name(input_variable)` where the `output_variable` is optional. If not `output_variable` is provided, by default the output from the function will be printed to the screen.

2.3.1 Nesting functions

Other than assigning the output of a function to a variable, functions can also be **nested** within other functions. For example, in the previous exercise, there was the line:

```
| letters <- c("A", "B", "C")
| paste(letters, collapse=" ")
```

Since the variable `letters` is not going to be used anywhere else in our program, we can simply nest the two lines together such that:

```
| paste(c("A", "B", "C"))
```

Just like in mathematics, the **order of operations** (or *operator precedence*) is the same in R. Operations are evaluated (or executed) from the innermost brackets first. So in the command above, the inner `c("A", "B", "C")` is executed first and the output (which is held temporarily in memory) is then passed to the outer function `paste()`.

With this shortcut, you can reduce the resources that a program takes (not to mention the lines of code) and create very complex nested function calls. Keep an eye out for nested functions throughout the workshop.



Break it up first, then reassemble

When new to R and coming across a complex nested function call, the tip is to **break it up** into sections. Break the functions into smaller pieces *starting from the innermost function* and assign them to variables first. Execute each line, one at a time and examine the output. Then when you are comfortable with what each line is doing, put the pieces together again. This will help you understand what each function is doing first before trying to decipher the entire line in one go. Once you become familiar with more functions, reading nested functions will become a breeze!

2.4 Data structures

We do not typically expect to work with atomic variables during an *-omics* analysis. To make the best use of R we want to use hundreds, thousands and millions of data points. We arrange these in one of a number of different data structures: vectors, lists, matrices and data frames.

2.4.1 Vectors

A **vector** is a group of components of the *same* type. Elements of a vector can be accessed by their numeric position (referred to as **indexing**) starting at position 1. They can also be named so that we can access the variable component using the namespace. The `c()` function (for *combine*) is used to join a number of single values together into a vector.

```
| l <- c(1,2,3)
| is.vector(l)
## [1] TRUE
```

Access the second element of the vector:

```
| l[2]
## [1] 2
```

What happens if we try to mix data types in a vector? `str()` is a function that describes the structure of an variable.

```
| mixedVector <- c(animal="sheep", meaningOfLife=as.integer(42), pi=3.1415)
| str(mixedVector)
## # Named chr [1:3] "sheep" "42" "3.1415"
## - attr(*, "names")= chr [1:3] "animal" "meaningOfLife" "pi"
```



Vectors and Collections of data

- Because elements of a vector must be the same type, when we investigate the content of `mixedVector` we can observe that the string, integer and numeric have all been cast to character.
- `mixedVector` is a named vector, so we can access elements by name as well as position.

`c()` can be used to combine vectors, not just individual values:

```
| c(l, mixedVector)
##
##           "1"          "2"          "3"      animal meaningOfLife
##           pi                  "sheep"        "42"
```



Using vectors

1. Enter the chunks above into your RStudio interface and review the output.
2. After executing the last code block, check on the values that are in variables `l` and `mixedVector` again.
 - a) How many elements does each variable hold? _____
 - b) Do any of them hold 6 elements? _____
 - c) What do you need to do to retain the combined vector with 6 elements?
3. Access the `mixedVector` element by name using `mixedVector['animal']`. What is the output? _____

Challenge, continue if you are ahead

4. Now enter the next chunk of code below and review the output.
 - a) Can you determine the difference between `mixedVector[2]` and `mixedVector[[2]]`?
5. What would the equivalent command to `mixedVector[[2]]` be, if you were accessing the vector by name?

```
#Access some elements of mixedVector by position or name
mixedVector[2]
mixedVector[[2]]
mixedVector['animal']
names(mixedVector)
```

2.4.2 Lists

Lists are similar to vectors in that they are one-dimensional structures for storing data and like vectors can be accessed by position or namespace. They differ from vectors in that they can contain multiple *different* data types. This makes them more flexible for storing data but more limited in the analyses that can be performed on them. Lists are created using the `list()` command, which has a similar format to `c()`.

```
mixedList <- list(animal="sheep", meaningOfLife=as.integer(42), pi=3.1415)
str(mixedList)
## List of 3
## $ animal      : chr "sheep"
## $ meaningOfLife: int 42
## $ pi           : num 3.14
```



Checking for membership

`%in%` is a way to test the membership of a single element in a list or vector of items. For example, try typing in the command: `"sheep" %in% mixedList`

2.4.3 Changing Vectors and Lists

While a list or vector is useful, their utility is much greater when we can add or remove elements from them, or change the value of existing elements.

```
# Add an item named 'pens' with a numeric value of 3 to the list
mixedList <- append(mixedList, c(pens=3))

# Now add another item named 'papers' with a chr value of 'bioinformatics'
mixedList <- append(mixedList, c(papers="bioinformatics"))

# Display just the second element of the list
mixedList[2]

## $meaningOfLife
## [1] 42

# Display all but the second element of the list
mixedList[-2]

## $animal
## [1] "sheep"
##
## $pi
## [1] 3.1415
##
## $pens
## [1] 3
```

```

##  

## $papers  

## [1] "bioinformatics"  
  

# Change the value of the pens item  

mixedList$pens <- 4  
  

# Remove the papers item  

mixedList$papers <- NULL

```



Using Lists

1. Enter the two `mixedList` code chunks into RStudio and review the output.
2. Try and retrieve different elements of the list
 - a) Can you get the second **and** fourth elements? *Hint:* use the combine `c()` function when indexing
 - b) What about everything **BUT** the second and fourth?
3. **[Optional]** If you are ahead, now try to getting a range of values, that is get the second to the forth elements. *Hint:* try using the colon (`:`) symbol to specify a range.

Now let us imagine that QFAB is actually performing a study on the course participants. To perform subsequent statistical testing we will need to know a little something about each of our consenting subjects and we are going to build a data structure to store this information.



3. Create three vectors that contain, respectively, the names, sex and age of the research cohort on this course. Let's have the vectors of type:
 - character for `name`
 - logical for `sex` (e.g. have female be TRUE) and
 - integer for `age`
 - a) When you have these vectors created, place them all in a single list variable called `subject`.

2.4.4 Factors

Factors are vector objects that contain grouping (classification) information of its components. Functionally they are similar to vectors in that they are a collection of objects of a *single* type. Factors are best applied when there are a limited number of different values. They are often used when categorical values are used in modelling or presenting data, e.g. phenotypic or study class data, such as diabetic, pre-diabetic, healthy.

For the example that we explore in the following exercise, we are going to consider a massive bucket of M&Ms. Assuming that the different colours are randomly mixed with an equal probability we would like to explore a million different sweets.



Using Factors

Use the code below to create a factorised `m and m` vector containing 1 million M&Ms of various colours.

Note, the output is again grouped together at the end for readability. The table below explains the new functions and what they mean.

```
colours <- c("red", "yellow", "green", "blue", "orange")
mandm <- sample(colours, 1000000, replace=TRUE)
object.size(mandm)
length(mandm)
mandm <- as.factor(mandm)
str(mandm)
object.size(mandm)
table(mandm)
head(mandm == "blue", 10)
which(mandm == "white")
length(which(mandm == "blue"))

## 8000280 bytes
## [1] 1000000
##  Factor w/ 5 levels "blue","green",...: 1 1 3 4 5 3 3 2 5 1 ...
## 4000688 bytes
## mandm
##   blue  green orange    red yellow
## 199868 199927 200481 200461 199263
## [1] TRUE TRUE FALSE FALSE FALSE FALSE FALSE FALSE TRUE
## integer(0)
## [1] 199868
```

What is happening here?

Function	Description
sample()	is a really useful method that allows you to sample elements from your data collection. The replace variable defined whether a value that is sampled should be replaced following its selection.
object.size()	describes the amount of memory that an object is using - this can be useful to identify really humungous objects that can be cleaned from your workspace.
length()	describes the number of elements present in the vector.
table()	prepares a tabular summary of values and the number of times that they occur - this can be used to summarise data effectively.
head()	shows only the first n elements of the vector, here we specified show the top 10 elements.
which()	returns a vector of the positions that means the condition, in this example which of the 1 million mandm's are 'white'? Since there are none, this returns an empty vector, represented by integer(0).
length(which())	nests the two functions together. It performs the inner which() function first, then passes the result to the length() function. As the names of the function suggest, this is finding out how many elements in mandm are equal to 'blue'?

The real challenge with factors is not creating them or using them but adding novel content to them, at least in terms of new categories. If we find an M&M that is not one of our five original colours, we can't just change the colour of that entry in our vector, we need to specifically add that colour to our factor list first. As a result, factors are suited more for *immutable* or static content. Other data structures provide simpler mechanisms for data manipulation.



Adding factor values

1. Follow the code chunk below to try to change the colour of the first M&M to the value “white”. The comments explain what each line is doing. You may find it best entering these commands directly in the console rather than the editor window, so you can see what is happening at each stage.
2. Add another colour option (e.g. brown) and change one or more of the vector elements to that new colour (e.g. elements from position 100 to 200)
3. **Challenge** add yet another colour (“purple”) and this time *randomly* change 20 elements to the new colour purple.

```
# First, just try to change the first element to white
mandm[1] <- "white"

## Warning in `<-factor`(`*tmp*`, 1, value = "white"): invalid factor level,
## NA generated
```

```
# Check the factor values, have we successfully added white?
levels(mandm)

# Now explicitly add white as a factor value
levels(mandm) <- c(levels(mandm), "white")

# Try again to change the colour of the first M&M
mandm[1] <- "white"

# And check the factor values again
levels(mandm)
table(mandm)
```

2.4.5 Matrices

A **matrix** is a collection of data elements arranged in a two-dimensional rectangular layout, effectively a table. Similar to the vector a matrix can contain only a *single* type of data.

It is often easiest to create a matrix from a vector of data. The next exercise will show you how to do this, follow the code below to generate a variety of matrices. R cannot guess the dimension of the data so you should specify either the number of rows (**nrow**) or columns (**ncol**) that the final data should have.



Creating a matrix

Note, this time the output follows *immediately* after each command, so that you can compare with what you get.

1. Enter the code below to generate a variety of matrices. *Hint:* Use the **Help** tab to search for **matrix** to find out more about the function and what the **nrow** and **byrow** options do.
2. Rather than printing the output directly to screen, store the matrix into a variable named **testMatrix**.

```
matrix(c(1,2,3,4,5,6,7,8,9,10,11,12))
##      [,1]
## [1,]    1
```

```
## [2,]    2
## [3,]    3
## [4,]    4
## [5,]    5
## [6,]    6
## [7,]    7
## [8,]    8
## [9,]    9
## [10,]   10
## [11,]   11
## [12,]   12
```

```
| matrix(c(1,2,3,4,5,6,7,8,9,10,11,12), nrow=4)
```

```
##      [,1] [,2] [,3]
## [1,]    1    5    9
## [2,]    2    6   10
## [3,]    3    7   11
## [4,]    4    8   12
```

```
| matrix(c(1,2,3,4,5,6,7,8,9,10,11,12), nrow=4, byrow=TRUE)
```

```
##      [,1] [,2] [,3]
## [1,]    1    2    3
## [2,]    4    5    6
## [3,]    7    8    9
## [4,]   10   11   12
```

```
| letter.mat <- matrix(c("A","B","C","D","E","F","G","H","I","J","K","L"),nrow=4, byrow=TRUE)
```



- Instead of typing out the letters above, look at the built in variable `LETTERS` in R (all caps). You can simply type this in the console and see what is returned. Use `class(LETTERS)` to find out the datatype.
- Also try the following command: `LETTERS[1:8]`, what do you get? _____
- You can also specify a list of sequence numbers in a similar manner, try executing the command `1:10`. What did you get?
 - * Assign this to a variable named `numRange`.

2.4.5.1 Indexing a matrix

To access the elements of a matrix you use the following methods:

```
letter.mat[2,] # returns the 2nd row
letter.mat[,3] # returns the 3rd column
letter.mat[2,3] # returns the element at row=2 and column=3
letter.mat[,-2] # returns all but the 2nd column
letter.mat[3:4,] # returns rows 3 to 4
```

2.4.6 Data frame

The **data frame** is to a matrix what a list is to a vector. Like a matrix it is a two dimensional table for storing data and is made up of equal length rows and columns, but like a list, the columns can contain data of *different* data types. Although within a column, the data must all be of the *same* type.



Follow the code chunk below to start exploring data frames. The comments explain what is happening in each stage. The first step tries creating a data frame using the same format as creating a matrix. This doesn't work though, see if you can work out what is actually happening here.

```
# Try creating a data frame using the same terminology as a matrix
data.frame(c(LETTERS[1:12]), nrow=4, byrow=TRUE)

# Now build a data frame from a temporary matrix
as.data.frame(matrix(c(LETTERS[1:12]), nrow=4, byrow=TRUE))
exampleDf <- as.data.frame(matrix(c(LETTERS[1:12]), nrow=4, byrow=TRUE))

# It's often useful to name the columns and rows of a data frame (or matrix)
colnames(exampleDf) <- c("x", "y", "z")
rownames(exampleDf)

# What is the size and structure of our data frame
dim(exampleDf) # Dimension of the data frame
ncol(exampleDf) # Number of columns
nrow(exampleDf) # Number of rows
str(exampleDf) # Structure information

# Retrieve a 'data slice' from the data frame
# First, get the contents of the first column
exampleDf[,1]

# Then the contents of the first row
exampleDf[1,]

# Finally, add a fourth column of a different type to the data frame
exampleDf[,4] <- c(1,2,3,4)
exampleDf
class(exampleDf[,3])
class(exampleDf[,4])
```



Creating a data frame

Previously, we created a list variable called `subject` that described the subjects within this group. A `data.frame` would be more convenient to use.

1. Convert the earlier list of vectors into a data frame.
2. If you gave names to the elements of your list earlier, the columns of your data frame should have inherited these names. Check that this is the case, and try naming the columns if not.



Accessing data.frame columns by name

An easier way to access data.frame columns by name is to use the \$ annotation. In the example code chunk about, you can access 'x' by using `exampleDf$x`.

2.5 Finding, describing and removing objects

In the previous data types section we have created a couple of lists, integers, numerics and a data.frame. This is the beginning of a data analysis workflow. While you are typing at your keyboard it is easy to become distracted and to forget the name of the variable that you created, or to create many temporary variables that you no longer need.



Tidying up objects

1. Use the `ls()` function to list all the variables in your workspace.
2. Use `rm()` to get rid of a variable or variables that you no longer want.
 - To remove multiple objects, just enter the names of all those objects separated by commas
3. Delete everything from your workspace using `rm(list=ls())`
4. Review your history (all the commands you have entered into the console in this session) using the RStudio History tab (top right window).

```
ls()  
rm(h)  
rm(x, y)  
rm(list=ls())  
history()
```

2.6 Best practise and variable naming conventions

Some quick best practice tips:

- Best practice suggests that you should not assign meaningless names to variables; give a variable a succinct but meaningful name that will be of value to understand which information you have saved.
 - `df` may be an obvious name for a data.frame, but `phenotypeDf` or `phenotype.df` may be more memorable later on in your analysis.
 - you can use **camelCase** (or camelNotation), underscore (`_`) or fullstop (`.`) for variable names in R.
- You can create thousands of variables in your R workspace but are they all necessary? It is worthwhile performing housekeeping on your data collection if you are creating many variables.
- If you have many variables in your R workspace then you may not be using the most ideal method to store your data. Consider data structures like `data.frame` and `lists` to better structure your data.
- RStudio allows you to remove all objects from your workspace by using the **Clear** button in the **Environment** tab as an alternative to the command: `rm(list=ls())`. (You can follow this with the `gc()` function for *garbage collection* to further ensure that memory has been freed up.)

Chapter 3

Going further with R

This section is a quick glance at the other features in R that will help you in your future data analysis. We will quickly look at:

- **Documentation** with R
- **Extending R with packages** - what are packages and how to install and load packages in R
- **Vignettes** that come with R packages including citing packages
- **SessionInfo** information to reproduce your data analysis, and
- **Other ways to run R**

3.1 Documentation

As you may now appreciate, there is a massive amount of information tied to R. There are functions already implemented for most of the typical data transformations that you may wish to perform. One of hardest challenges with R is finding the method that does what you need. It's out there!

Function	Description
? will give you information about how to use the <code>t.test</code> function in R.	is used to find out information about a specific
<code>apropos()</code>	will return a vector of all the objects or function

containing the specified search string. E.g. `apropos("test")` will find around 50 different statistical test functions, while `apropos("mixed")` will find `mixedList` and `mixedVector` variables that we created earlier (or at least it would if we hadn't just deleted all our objects!) | `help.search()` | searches the help documentation for entries matching the search string. Items will match if the search term is included in the function description or keywords, not just the name as is the case for `apropos`. |

```
?t.test  
apropos("test")  
help.search("topic")
```

3.2 Extending R with packages

One of the reason that R has such a following in the statistical and biological fields is the range of available packages that can be used to extend the utility of the software. **Packages** are collections of functions, data, and compiled code written in R and other languages that are encapsulated and distributed in the package format. Packages are generally domain or analysis specific, thus are only installed on a as needed basis. The directory where packages are stored is called the library.

R comes with a standard set of packages. Others are available for download and installation from repositories that include CRAN, Bioconductor and R-forge. To use a package, you must first install it and then load it into your session using either the `require` or `library` functions.

Because of the configuration of the training server we are using today, you can only install packages to *your own personal library* and not a system wide installed. This means that the packages you install will not be available for use by another participant in this class.

The code below outlines the process involved:

```
# Look to see which packages are already installed
library()

# To install one that is not there, use install.packages
# N.B. You will get a message "Would you like to use a personal library instead?"
install.packages("tidyR")

# Once it's installed, we need to load it
require("tidyR")
# or
library("tidyR")

# library() shows the installed packages. search() shows which of those are loaded.
search()

# If you no longer need a package, then you may want to unload it with detach()
# You need to specify "package:" before the name, as per the listing from search()
detach("package:tidyR")
```



Using packages

1. Use `library()` and `search()` to see what packages are available on the training server and loaded into your workspace.
2. Load the `limma` library using either `require()` or `library()`.
3. Use `search()` again to check it has loaded correctly.

3.3 Vignettes

Just now, we explored some of the help functions built into R. As well as this builtin help, most R packages also come with their own, often extensive, documentation. This is normally in the form of a vignette, a document that provides a task-oriented description of package functionality. Vignettes contain executable examples and are intended to be used interactively. You can also download the vignette for a package that isn't installed on your system by visiting the CRAN, Bioconductor project webpages, or via an internet search.



Using vignettes (optional)

1. The `vignette()` command provides access to vignettes for installed packages.
2. Run `vignette()` at the console to list all the vignettes available.
 - Notice that a package can, but does not always, have a vignette with the package name
- for example the `annotate` package has an `annotate` vignette, but `BiocParallel` only

has the `IntroductionToBiocParallel` vignette

3. Type `vignette("biomaRt")` to open the vignette for the biomaRt package, which links to the Biomart gene annotation database.
4. Some vignettes have non-unique names - for example, several packages have an intro vignette.
 - a) Try `vignette("intro")`
 - b) To get the limma package vignette, use the command `vignette("intro", package = "limma")`

3.3.1 Citations

The R core development team and the very active community of package authors have invested a lot of time and effort in creating R as it is today. Please give credit where credit is due and cite R and R packages when you use them for data analysis.

```
citation("VennDiagram")
citation("limma")
```

3.4 SessionInfo

Congratulations - you have now conquered the basics of R and you are ready to start breaking things! One of the most useful aspects of R is the rich ecosystem of supplementary packages that can aid, facilitate and empower your research. The Bioconductor framework contains hundreds of packages of methods of relevance to biologists working with biological data (population genetics, SNPs, NGS reads, microarrays, mass spectrometry etc). As we discovered in the Vignette section, there are well documented tutorials and guides that will work us through the application of fabulous methods and workflows. You will come across some unexpected 'error' messages.

R packages are written by people like the QFAB bioinformaticians. Well intentioned, but busy. Code is crafted lovingly and tested in the scenarios that are implicit in our daily development. You are likely to be running R on a different computer with combinations of installed packages that are subtly different to those that we wrote the software with. If you identify a "bug" in that the package gives a wrong result, fails with a hairy error message or misbehaves it is worthwhile to let the developer know that you are having a problem. Reproducibility of the error is critical - we would like to know which version of R is being used and all of the packages that are loaded in memory so that the developer can agree that yes, there is a problem and can help in the resolution.

It is trivial to report the data on the R environment that you are working in - the `SessionInfo` function reports installed packages and their versions.

```
require(limma)
sessionInfo()

## R version 3.3.0 (2016-05-03)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 14393)
##
## locale:
## [1] LC_COLLATE=English_Australia.1252 LC_CTYPE=English_Australia.1252
## [3] LC_MONETARY=English_Australia.1252 LC_NUMERIC=C
## [5] LC_TIME=English_Australia.1252
##
## attached base packages:
```

```
## [1] stats      graphics   grDevices utils      datasets   methods    base
##
## other attached packages:
## [1] tidyverse_0.4.1      BiocInstaller_1.22.3 bookdown_0.1.4
## [4] RevoUtilsMath_3.2.5  reshape2_1.4.1       ggplot2_2.1.0
## [7] limma_3.28.21
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.5          formatR_1.4        plyr_1.8.3
## [4] tools_3.3.0          digest_0.6.9       evaluate_0.9
## [7] tibble_1.1           gtable_0.2.0      shiny_0.13.2
## [10] DBI_0.4-1           rstudioapi_0.5    yaml_2.1.13
## [13] parallel_3.3.0      stringr_1.0.0     dplyr_0.4.3
## [16] knitr_1.13          grid_3.3.0       R6_2.1.2
## [19] rmarkdown_1.0.9005   magrittr_1.5     scales_0.4.0
## [22] codetools_0.2-14    htmltools_0.3.5  assertthat_0.1
## [25] mime_0.4            xtable_1.8-2      colorspace_1.2-6
## [28] httpuv_1.3.3       labeling_0.3     stringi_1.1.1
## [31] miniUI_0.1.1        munsell_0.4.3
```

3.5 Other ways of running R

Today we are using RStudio, but R can be run in a number of different ways; for example, you could embed some R functionality in a computer program written in Python or Java. As discussed in sections 1.2 and 1.3, both R and RStudio provide an interactive environment where you are prompted to enter a single command at the time. In this training course we are providing you with snippets of information that could be used within a data analysis workflow. These are intended to be run interactively.

In a bioinformatics laboratory the scientists who use R craft a mixture of packages and workflow scripts into analytical pipelines. R may be even be called by other software environments. Bioinformaticians will typically use the R console for refining, tuning and modifying existing scripts that they have written. The workflow is run through master scripts. In this section we will look at a simple R script that will create a table of information and write it out to file. We will be having a more complete look at the process of writing data to file in a later section (Section 4.3).

The advantage of running R scripts as a batch analysis is that larger and more complex analyses (such as the mapping of short DNA sequence reads) can be run overnight and each of the commands will run successively as prior commands complete.

3.5.1 Preparing an R script

An R script is a container for multiple R commands. It is intended to be run largely hands-off.

The RStudio software provides us with a very convenient way to create an R script - in the file dialog there is the option to create file and an R Script is a primary file type. R scripts typically have the extension .R. The file is a plain text file and needs to know the packages that should be loaded, the objects that should be set and the working environment where we should be working.

3.5.2 Running an R script

The following command is how you run an R script from a terminal window. This is outside of RStudio.
R CMD BATCH [options] my_script.R [outfile]
R.exe" CMD BATCH
--vanilla --slave "c:\my projects\my_script.R"

We will demonstrate the next exercise to show you the method to run a script file using RStudio. You can repeat this exercise in your own time.



Running a script in RStudio

- Create a new script file by going to **File > New File > R Script**.
- Enter some commands in the editor window, for example the code chunk below.
- Save the script with a file name e.g. “area-rectangle.R”
- In RStudio, there are 3 ways you can run this script from top to bottom:
 1. In the **console** window, type in the command: `source("area-rectangle.R")`. What is the output?
 2. In the **editor** window, top right corner click on the **arrow** next to the **Source** button, select **Source**. What is the output?
 3. Repeat step 2 but this time select **Source with Echo**. What is the difference between step 2 and 3?

```
print(date())
print("This is an R script!")
length <- 10.25
width <- 4.35
area.of.rectangle <- length * width
print(paste('The area of a rectangle with length="',
            length,
            "cm by width=",
            width,
            "cm is",
            area.of.rectangle,"square cm."))
```



Use [Tab] to autocomplete

You can use the tab key to autocomplete a variable name, a function name or a filename (if the file is in the working directory). This will save you a lot of time and prevent typing errors when you are analysing your data.

Chapter 4

Importing and exporting data from R

It is likely that you will not always be able to complete your R analysis in a single sitting. Rather than having to recreate all your variables each time, there are various ways that you can save data for reuse. This section covers how to save your session data, importing and exporting data.

4.1 Saving session data

If you started a new Project for your work, at the end of the day you can select **File > Close Project**. This will save all of your variables, your history, your open files and a range of other information. When you then re-open R, you can double click on the project file (which will have the suffix **.Rproj**) to continue where you left off.

As well as saving your entire project, you can save just specified variables, using the `save()` command, for later reloading using `load()`.



Load and save variables

Follow the code chunk below to:

- create a new variable
- save that variable to a file
- delete that variable from the active workspace, and
- load it again from the file.

```
# Create a new variable
hw <- "HelloWorld"

# Save that variable to a file called "myfile.RData"
save(hw,file="myfile.RData")

# Delete hw from the workspace
remove(hw)

# And confirm that it's gone
print(hw)

# Reload the data from the file
load("myfile.RData")
print(hw)

# You can save multiple variables
testVector <- c(1,2,3)
save(hw, testVector, file = "myfile.RData")
remove(hw, testVector)
load("myfile.RData")
print(hw)
print(testVector)
```

4.2 Current working directory

An aside before we move onto importing data. When working with R, you need to be aware of your **current working directory** (CWD). As the name suggests, this is the directory (folder) in which you are currently working in. Any files you wish to import into or export from R will be with respect to this directory, if no paths are specified.



The *Files* tab in the bottom-right corner is *not* always set to your CWD. Think of this as a Windows explorer (or Finder in Mac) built-in with Rstudio. This tab allows you to navigate your file system and check on files. This means that during your session, you may navigate away from your CWD.

To check your CWD use the command:

```
getwd()  
## [1] "G:/QFAB/16_Training/16_01_Workshops/16_01_023 Introduction to  
R/2017017_UQ_StLucia/Manual"
```

and to change your CWD to another location use the `setwd(file-path)` command.

4.2.1 File paths

File paths come in two forms:

* **relative filepath** is relative to the current working directory. For example, using
a Windows system, if my currently working directory is `C:\Users\john.smith\Documents`
and I save a script called "area-rectangle.R" without specifying a specific location,
then R will save the file in the location
`C:\Users\john.smith\Documents\area-rectangle.R`.

* **absolute filepath** is the fullpath of the location starting with the drive letter in Windows (e.g. C:\ or D:\); or in Linux and Mac starting with the `\\` symbol, meaning the *root* directory (the top).



R uses the forward slash (/) to specify filepaths by default. If you use the Windows format with the backwards slash (\), you will get an error. You can use two backward slashes (\\) to escape the error but it's faster and easier to read to use the / slash.

4.3 Reading and writing tabular data

Data is most commonly shared in tabular formats (Excel presents data in a tabular format). These are most-often text files that contain columns of data separated with comma, tab or other field-delimiters. In R there are some simple methods for quickly reading in massive amounts of data.

For the next few steps in this course we are going to explore a public dataset of human tissue gene expression data that was originally produced by Novartis. The dataset is known as the Novartis GeneAtlas data (now called BioGPS) and can be found in the GEO database under accession GSE1133 (<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE1133>). We will look at some Affymetrix gene expression data from the study. These data have been prepared from microarray hybridisations using the Affymetrix Human Genome U133A Array. The data is in the file data/Intro_to_R/GSE1133-GPL96_series_matrix.txt.

```
# read data from the training server
data.path <- "data/Intro_to_R"

# Of course, you will need a different path if the data is on your local P.C.
# e.g. data.path <- "C:/Intro_to_R/Data"
# Create a variable pointing to the file location.

gene.atlas.file <- file.path(data.path, "GSE1133-GPL96_series_matrix.txt")

# Then look at the start and finish of that file
head(readLines(gene.atlas.file))
tail(readLines(gene.atlas.file))
```

What is happening here?

Function	Description
<code>readLines()</code>	reads each line from a target file into a vector. See also <code>read.table</code> , <code>read.csv</code> functions.
<code>file.path()</code>	<i>DOES NOT</i> read in the contents of a file. It just point to the location of a files on your system. You need to use <code>readLines</code> , <code>read.table</code> or similar to access the contents of that file.

Using the `head` and `tail` commands helps us understand the nature of the file. It seems that the ! prefix indicates metadata rather than data.



Beware of printing large variables with RStudio Server, you can freeze your session. Judiciously use methods to limit the output:

- read into a variable first

- limit the output
- use `head()` / `tail()` functions

We have started to explore the contents of the expression data file using the code above, but at the moment have not imported it into a variable in R. Before we get to that stage, we will first load the annotation data associated with this expression dataset.



Reading in tabular data

1. If you have not already done so, use the previous code chunk to set up your `data.path` and take an initial look at the expression data in the variable `gene.atlas.file`.
2. Follow the annotated code below to use `read.table()` to import the annotation data, or metadata, for this experiment

```
# As before, define the location of the file
annotation.file <- file.path(data.path, "GSK_RNA.sdrf")

# Look at the first six lines of the data
# sep = "\t" means that the data is tab-separated
head(read.table(annotation.file, sep="\t"))

# The last command included the header as the first line of data.
# Specify (header = TRUE) that the first line contains header information
head(read.table(annotation.file, sep="\t", header=TRUE))

# Use read.table to import the data into an variable
expressionMetadata <- read.table(annotation.file, sep="\t", header=TRUE)

# Find out some information about the structure and contents of that new variable
tail(expressionMetadata)
dim(expressionMetadata)
colnames(expressionMetadata)
str(expressionMetadata$Material.Type)

# By default, read.table imports text fields as factors. That's probably
# not what we want, so disable it with stringsAsFactors=FALSE
expressionMetadata <- read.table(annotation.file,
                                    sep="\t", header=TRUE,
                                    stringsAsFactors=FALSE)
str(expressionMetadata$Material.Type)
```

Note that the lines starting with the default comment symbol (#) have been omitted from the import.



Read the expression data into a variable called `gene.atlas`

In the last couple of sections we have used `readLines()` to look at the content of a file and we have used `read.table()` to read some experimental metadata into a data frame. Your challenge is now to import the “GSE1133-GPL96_series_matrix.txt” data into a data frame variable called `gene.atlas`.

Hints: * Use the `read.table` documentation to discover the nuances of the function. * This import is not as simple as the annotation file. You will need to make some inferences about the structure of the file (*Hint:* use `head()` and `tail()` and look at some of the options for

`read.table()` in the Help.) * Make sure that the proper column names are included in the data frame.

```
##      ID_REF GSM18865 GSM18866 GSM18867 GSM18868 GSM18869 GSM18870 GSM18871
## 1 1007_s_at    2082.0   1577.2   231.8   328.4   83.6   73.0   53.2
## 2 1053_at     598.4    256.2    57.0   102.1   39.1   43.6   73.7
## 3 117_at      194.2   135.9   512.3   713.4   108.6   75.7  165.7
## 4 121_at     1910.5   1516.2   529.8   990.2   383.4   396.4  227.7
## 5 1255_g_at    138.4    53.8    30.7    75.7   16.9   22.8    4.1
## 6 1294_at     263.4    287.2   347.8   345.6   300.8   377.7  234.2
##   GSM18872 GSM18873 GSM18874 GSM18875 GSM18876 GSM18877 GSM18878 GSM18879
## 1    70.0    92.2    75.0    78.8   60.1   117.7   149.4  105.8
## 2    32.9    37.6    38.4    43.0   82.6   30.2     9.0   29.6
## 3   152.6    71.1   105.3   140.2   82.4   50.0   11.1   38.2
## 4   251.6   260.0   260.1   287.3   270.6   269.8   395.2  340.2
## 5     4.9    23.3    14.7    27.0     4.2    27.8     9.0   17.6
## 6   256.3   245.1   171.7   369.3   287.0   382.3   424.4  499.9
##   GSM18880 GSM18881 GSM18882 GSM18883 GSM18884 GSM18885 GSM18886 GSM18887
## 1    97.6   109.6   113.9   42.3   83.9   76.4   93.1  247.8
## 2    15.4    38.0    57.5   76.5   71.5   55.1   87.5  133.1
## 3    19.4    38.8     8.5   34.4    9.6   10.6   32.6   83.5
## 4   342.9   215.3   159.7   257.2   334.5   227.4   216.5  627.6
## 5     7.3     6.9    9.3     6.4     6.4     4.2    19.1   39.0
## 6   440.8   263.9   159.2   171.4   104.7   258.8   221.6  255.6
##   GSM18888 GSM18889 GSM18890 GSM18891 GSM18892 GSM18893 GSM18894 GSM18895
## 1   277.9   405.0   351.5   296.1   463.6   225.6   252.9  210.2
## 2   140.9   196.8   214.7   63.9   50.6   58.6   40.1   67.1
## 3   85.3    67.4   72.9   139.3   79.7   60.7   66.5   64.3
## 4   806.5   245.0   312.7  1038.8   951.5   632.7   620.8  797.9
## 5    60.5     6.8   19.3   53.0   40.2   34.0   35.7   41.5
## 6   287.4   188.6   183.3   201.4   193.9   135.0   103.7  106.3
##   GSM18896 GSM18897 GSM18898 GSM18899 GSM18900 GSM18901 GSM18902 GSM18903
## 1    98.8   302.0   358.2  1008.2   806.8   499.6   554.9  387.5
## 2   106.3   177.6   120.5   72.0   93.0   89.1  111.1   58.7
## 3    79.0    84.6   93.3   145.3   163.9   217.2   115.6  159.4
## 4   601.5   703.2   749.1   991.6  1042.0  1100.6  1022.6  1163.6
## 5    44.8    60.1   45.4   54.6   59.3   47.2   67.6   40.7
## 6   103.5   151.9   189.0   396.3   430.4   397.6   302.8  421.0
##   GSM18904 GSM18905 GSM18906 GSM18907 GSM18908 GSM18909 GSM18910 GSM18911
## 1   404.8  1195.7   713.4   91.3   30.9   518.4   598.4  1442.0
## 2    65.6   108.8   53.3   141.7   90.7   119.3   51.9   23.5
## 3   235.2   403.9   212.1   99.2   85.2   278.6   148.8  223.4
## 4  1314.8  2879.3  1847.8   722.2   541.6  1924.4  2362.2  2004.0
## 5    58.1   145.5   18.7   45.4     1.7   86.4   95.5   410.7
## 6   408.0   463.3   358.8   135.6   118.9   451.9   407.7  475.6
##   GSM18912 GSM18913 GSM18914 GSM18915 GSM18916 GSM18917 GSM18918 GSM18919
## 1  1400.8  2076.0  2644.9  1978.3  2182.1  2784.9  2738.3  2855.6
## 2    26.5    32.5   94.0   108.0   51.5   33.0   96.8   61.3
## 3    71.2   351.3   427.5   306.2   310.6   213.6   365.9  205.8
## 4  2451.0  2449.2  2026.1  1901.9  1950.4  3687.7  3131.7  2105.6
## 5   450.0   104.6   360.4   109.2   140.8   110.3   196.5  242.4
## 6   388.5   873.8   973.5   406.1   422.3   624.2   518.2  527.4
##   GSM18920 GSM18921 GSM18922 GSM18923 GSM18924 GSM18925 GSM18926 GSM18927
## 1  3351.4  2390.5  2837.6  3499.5  3183.4  2725.0  2825.7  2865.1
## 2    78.4    69.8   93.0   50.8   72.4   87.7   27.3   103.6
## 3  120.3   188.8   183.4   458.4   396.5  179.7  300.7  222.0
## 4 2231.3  1538.6  1975.7  2533.2  2679.1  2100.0  2223.5  2146.3
## 5   319.3    87.9   192.3   230.1   244.3   261.7   349.6  226.5
```

4.3. READING AND WRITING TABULAR DATA

```

## 6   406.6   270.5   361.1   581.4   637.2   572.9   543.6   296.6
## GSM18928 GSM18929 GSM18930 GSM18931 GSM18932 GSM18933 GSM18934 GSM18935
## 1   2901.3  1857.2  3347.2  2625.7  2532.0  3173.0  4019.9  2662.7
## 2    81.6    16.3    83.4    93.8    18.6    14.6    112.1    69.1
## 3   222.5   195.3   234.1   342.1   169.6   130.3   255.9   394.9
## 4   2457.2  1585.6  1709.5  1855.0  1618.3  1183.5  1854.4  2054.8
## 5   321.2   130.0   161.2   173.2   158.8   192.3   195.5   138.9
## 6   424.6   307.9   445.3   551.3   526.5   285.1   443.9   508.1
## GSM18936 GSM18937 GSM18938 GSM18939 GSM18940 GSM18941 GSM18942 GSM18943
## 1   2727.0  2119.1  1769.4  2550.2  2158.4  2694.2  2497.1  4398.5
## 2    25.9    136.2   18.9    17.2    153.0    29.0    190.7   108.8
## 3   276.6   411.8   73.2    317.9   367.9   239.3   409.2   341.1
## 4   1289.8  1412.5  2304.5  2885.6  2003.1  2591.7  3007.3  1614.5
## 5   102.3   321.6   296.3   332.8   216.3   374.5   308.8   104.7
## 6   479.7   449.6   416.0   882.8   407.5   957.9   650.1   350.0
## GSM18944 GSM18945 GSM18946 GSM18947 GSM18948 GSM18949 GSM18950 GSM18951
## 1   4808.6  2141.3  2010.9  1494.8  1319.3  1238.5  1807.9  1240.9
## 2    30.8    18.8    97.4    9.7    30.1    140.3    86.9   112.5
## 3   330.1   188.8   205.0   229.1   48.3    203.2   167.3   90.8
## 4   1882.9  1484.9  1442.7  2284.5  2051.4  1455.4  1347.1  2623.5
## 5   160.7   141.4   147.7   130.0   102.4   57.1    83.0   201.2
## 6   396.0   310.6   267.7   629.5   457.9   442.8   414.7   377.8
## GSM18952 GSM18953 GSM18954 GSM18955 GSM18956 GSM18957 GSM18958 GSM18959
## 1   925.5   1084.8  691.6   1982.6  1786.9  2476.4  1822.6  717.9
## 2   127.0   106.3   24.5    13.1    46.3    20.2    62.0   52.8
## 3   132.8   301.7   286.0   238.1   160.3   181.1   145.4   108.9
## 4   1778.8  3045.8  1471.3  5618.5  6562.1  1715.6  1250.7  1376.9
## 5   98.4    134.9   114.0   116.0   132.5   142.5   75.8    103.9
## 6   560.3   665.4   653.0   870.8   597.4   326.6   549.3   244.5
## GSM18960 GSM18961 GSM18962 GSM18963 GSM18964 GSM18965 GSM18966 GSM18967
## 1   817.6   2152.8  2450.7  1632.0  1734.4  874.3   822.7  2202.5
## 2    15.5    85.9    88.8    48.6    53.8    64.0    46.0   79.2
## 3   114.8   70.6    154.3   214.0   93.3    222.7   261.8   165.1
## 4   1298.4  7545.0  6722.3  8545.9  9865.6  1147.4  935.1   1734.1
## 5   38.3    83.9    76.2    90.8    46.8    26.6    35.3   102.8
## 6   557.0   410.7   269.9   391.6   437.6   276.1   276.3   235.2
## GSM18968 GSM18969 GSM18970 GSM18971 GSM18972 GSM18973 GSM18974 GSM18975
## 1   2383.5  469.5   889.9   417.3   374.9   984.3   1183.0  476.8
## 2    13.1    73.7   161.5   60.8    41.6    48.1    27.1   12.7
## 3   268.3   105.2   202.6   64.0    45.0    21.5    47.1   170.6
## 4   1415.8  924.1   2068.1  787.0   647.0   491.2   645.3   1480.7
## 5   86.9    67.5   137.7   49.6    40.5    27.0    39.8   85.9
## 6   383.6   212.9   472.2   255.6   232.3   119.6   84.4   331.5
## GSM18976 GSM18977 GSM18978 GSM18979 GSM18980 GSM18981 GSM18982 GSM18983
## 1   408.9   1083.1  903.0   940.6   2153.8  1071.2  1029.1  1249.0
## 2    60.3    18.1    46.6    38.8    29.7    95.4    67.1   54.3
## 3   140.5   137.1   189.1   91.9    104.7   322.5   326.8   317.5
## 4   1224.4  1324.6  1678.1  1345.9  1196.6  2295.8  2019.3  2608.5
## 5   15.7    4.6     34.7    77.6    19.8    560.1   584.6   430.7
## 6   426.7   218.7   345.6   210.9   243.3   330.8   389.2   617.4
## GSM18984 GSM18985 GSM18986 GSM18987 GSM18988 GSM18989 GSM18990 GSM18991
## 1   1433.6  997.3   656.2   902.5   809.9   1302.8  1062.4  1061.3
## 2   159.8   120.4   151.9   75.7    24.9    24.2    101.6   6.8
## 3   233.9   262.5   238.8   278.5   465.6   404.1   309.1   78.7
## 4   2923.3  2187.5  1549.6  2372.4  1900.4  2860.4  2122.6  990.1
## 5   378.0   584.2   120.1   537.2   596.3   602.4   299.2   80.9
## 6   692.6   514.2   309.2   541.9   302.0   432.7   440.5   446.5
## GSM18992 GSM18993 GSM18994 GSM18995 GSM18996 GSM18997 GSM18998 GSM18999

```

```

## 1   760.0  2448.7  1484.5   878.1   813.7   353.1   439.1  1161.8
## 2    93.8     36.0    13.1    72.6    33.7    37.2     8.6   139.2
## 3   122.6   177.8   186.9   165.8   286.5   158.7   87.2   334.4
## 4  1055.1  2063.2  1373.0  1912.1  1939.8  1344.2  1126.5  2506.8
## 5   132.9    65.3   128.7   145.3   280.8   64.6   117.2   192.1
## 6   263.4   373.2   480.1   680.1   535.7   206.7   197.5   496.9
##   GSM19000 GSM19001 GSM19002 GSM19003 GSM19004 GSM19005 GSM19006 GSM19007
## 1   761.0  1233.3  1468.4  1201.7  1655.4  1481.1  2163.5  1310.4
## 2    35.2     51.2    34.2    53.1   193.7    46.7   137.5    37.6
## 3    84.9   248.2   123.6   491.9   317.2   915.7   788.7   902.9
## 4  2278.1  2459.3  2234.4  3817.7  3342.6  4585.7  5395.7  3728.7
## 5   401.6   180.7   188.4   384.2   541.8   640.4   882.2   441.7
## 6   334.8   376.1   542.4   984.5  1168.3   862.2  1015.1   901.8
##   GSM19008 GSM19009 GSM19010 GSM19011 GSM19012 GSM19013 GSM19014 GSM19015
## 1   1205.3  1564.5  1461.3  2925.4  3369.9  1657.9  1322.6   732.2
## 2    37.9    83.3   126.7   163.9    50.3    52.4   143.8    77.6
## 3   838.6   470.8   456.1  1297.6   708.1   616.0   658.0   107.7
## 4  4631.8  3678.2  3594.2  8048.9  7283.1  4009.4  3086.5  1574.5
## 5   426.9   226.2   471.3  1114.2  1115.9   343.7   285.7   159.4
## 6   1124.5   649.7  1710.2  2061.3   1637.6   901.2  1008.4   330.3
##   GSM19016 GSM19017 GSM19018 GSM19019 GSM19020 GSM19021 GSM19022
## 1   841.7  1557.3  1506.3  1791.0  1697.2   923.0  1078.4
## 2    42.6    95.2    16.9    17.9    82.9    23.2    13.6
## 3   438.1   228.1   201.1   383.3   327.9   309.4   420.6
## 4  1720.2  2405.3  1707.4  1695.3  1317.5  1430.1  1734.7
## 5   266.0   105.8   190.9    40.2   154.2    61.9   111.3
## 6   487.9   459.3   487.8   535.3   458.6   363.0   307.2

```

4.4 Scan

`scan` is a primitive method that can be used to import data from a variety of sources. In most instances it would be preferable to read data from file using `read.table` or `readLines`. `scan` is however magnificent for reading in data from the standard input or from other software streams. As a quick demonstration we can read in data as you type it:

```
| scan(what=numeric())
```

4.5 Reading from a web connection

Most of the time we will be reading information into R that is stored on our local filesystem, but R can also import data directly from the web. The `read.table()` and `readLines()` functions are happy to read in from a web socket. This is great but requires that the data be present on a web-page to download.

The `RCurl` package provides a much more flexible approach to accessing data that is on the web and is worth reviewing if you wish to scrape a web-accessible database in a more automated fashion.

Of greatest interest however is the ability to download pre-structured biological data from the web. This can be managed using packages such as `bioMart` and `GEOquery`.



As our server is behind a firewall, the following code chunk will not work as it needs to access data from the Internet. In such cases, you might need to set up your http/ftp proxy using the `Sys.setenv` function.

You can try these commands again on your local RStudio.

```
web.data <- readLines("http://data.princeton.edu/wws509/datasets/effort.dat")
class(web.data)
head(web.data)

# ho-hum - we can read in the data but it is neither up-to-date or biorelevant...
# we can pull in the GeneAtlas data straight from GEO
library(GEOquery)
geo.gene.atlas <- getGEO(myGSE, destdir=".")  
  
# or load the GEO data previously downloaded in text format
gene.atlas.gpl96 <- getGEO(filename=file.path(data.path,
                                              "GSE1133-GPL96_series_matrix.txt"))
class(gene.atlas.gpl96)
str(gene.atlas.gpl96)
colnames(pData(gene.atlas.gpl96))

# w00t - we have an ExpressionSet
head(exprs(gene.atlas.gpl96))
gene.atlas <- exprs(gene.atlas.gpl96)
```

What is happening here?

- This code is a little scarier than some of the code that we have run thus far. `getGEO()` imports data from the Gene Expression Omnibus (GEO) database. Since there are three different platforms used within the study we need to select for our Affymetrix platform of interest - we select the appropriate element from the list.
- `pData` is a method that selects the phenotype data from the header of the GEO file (this is contained within the ! flagged fields from our earlier data import). This phenotypic data is rather crucial to the analysis.
- The data has been imported as an `ExpressionSet` rather than as e.g. a `data.frame`. The `ExpressionSet` is a crucial data-type that is used in most gene expression data workflows. The `ExpressionSet` contains the phenotypic data, the expression signal intensities and other data that could be of value to the data analysis. The expression data can be extracted using the `exprs()` function.

Chapter 5

Transforming and summarising data in R

So far we have looked at creating data, importing data and we have applied a couple of functions so that we can view the contents of the data. Now is where things can start getting more interesting. Once you have loaded or created your data, you want to start processing and making sense of your data.

In this section we cover how to write your own **user-defined functions** (Section 5.1), which as the name suggests are functions you write yourself instead of using the built-in functions. Sections 5.2 and 5.3 will cover methods for looping through your data and we finish this chapter off with **basic data analysis methods** in Section 5.4.

5.1 User defined functions

While there are hundreds of packages that can do everything that you need in R, it is often quicker to write a function to do something specific than it is to find someone else's function. For purposes of scientific reproducibility, for automation and for clean code, wrapping chunks of code into functions is a must!

Crafting a function is pretty simple. We use the `function()` command, and provide it with a name for our new function and some logic for that function to perform.

```
# Create a function called addOrminus
addOrminus <- function(a, b=2, add=TRUE) {
  if (add) { # If the logical add is TRUE, then carry out this section
    a + b    # Add a and b and return the result
  } else {   # or if add is FALSE, then do this instead
    a - b    # Subtract b from a and return that result
  }
}

# Test what we get with some different input options
addOrminus(1)
addOrminus(2, 4)

# Assign the answer from the function to a new variable.
answer <- addOrminus(2, add=FALSE, b=4)
```

What is happening here?

In the example above, we create a function called `addOrminus`. This function takes the following three arguments:

- (i) a *required* numerical value named **a**,
- (ii) an *optional* numerical value named **b** and
- (iii) an *optional* logical value named **add**.

Arguments **b** and **add** are given **default values** of 2 and TRUE respectively by using the equal (=) sign in the function declaration. These defaults values are used in the logic of the function if no other values are provided when the function is called.

You call the function using its name, `addOrminus(1)` and you can assign the answer to a new variable.

```
# Create a second function to report absolute magnitude,
# calculated as the square root of the square of the input
absMagnitude <- function(x) {
  return(sqrt(x^2))
}
absMagnitude(-123.2)
```

What is happening here?

The second example named `absMagnitude`, takes a single required numerical value only and **returns** the absolute magnitude of that value.

Functions can contain any data-type as input, they can be simple data types (Section 2.2) or even the more complex data structures (Section 2.4). You can have a mixture of required and optional values. By convention, required arguments are listed first, then followed by any optional arguments.

Creating and using functions

1. Use the code chunks above to create and experiment with the two example functions.
2. Create a new function based on the `addOrMinus` code that can be used to **multiply** or **divide** two numbers.
3. Using our character manipulation commands from Section 2.2.4, create a function that accepts two inputs: (i) a character string and (ii) an integer. This function will return the letter at that number position in the string in UPPER CASE.
4. How about a function that takes a row count and a column count and creates a matrix of those dimensions filled with either numbers or letters (as defined by a third input parameter).

Functions and Packages

- Like variables, give functions a meaningful and self-explanatory name.
- Keep things in functions simple and structured.
- Use some form of indentation to keep your code readable so that you can follow what is happening. Have a look at the R coding conventions and links contained within at http://journal.rproject.org/archive/2012-2/RJournal_2012-2_Baaaath.pdf.
- An explicit `return()` is not required in a function. By default a function will return the last executed command. However, it is worthwhile to include a return statement if you have a complicated function that has multiple endpoints and logical decisions. This also makes debugging easier.
 - * You can also create and return complex data structures e.g. `return(list(...))`
- If you think that writing a function is great, the next step will be to write whole packages that automate, streamline and focus your research. This is not difficult but unfortunately does not fit in the scope of this workshop.

5.2 For loops, While loops and Vectorisation

for loops are controversial in R. They are slow and inefficient and due to the way that R has been crafted (the concept of vectorisation) there are typically faster and more efficient ways to process data. This is an introductory course and we don't have time to go into the depth required for the various `apply` and `dplyr` methods that will guarantee amazing performance. The for loop is suitable for getting proof-of-concept workflows running!

To create a loop, we need to set a looping condition. There are two types in R:

1. a `for()` loop will run on each element of the input variable, and
2. a `while()` loop repeats until an exit condition is reached. **Note**, you must ensure that you update the exit condition so it eventually becomes true, otherwise you will get an **infinite looping** error where your code *never* exists until it uses up all available resources (usually the memory).

We also need to wrap that loop in braces (`{` and `}`) to define the start and end of the code block to be iterated.

Examples of a for loop: * `seq(10)` - creates a list of integers from starting from 5 to 35 in intervals of 5.

```
# a simple for loop along a sequence of numbers generated by seq()
mysequence <- seq(5,35,by=5)
for (i in mysequence) {
  print(i)
}

## [1] 5
## [1] 10
## [1] 15
## [1] 20
## [1] 25
## [1] 30
## [1] 35
```

```
# You can also loop through a vector, or even a matrix
for (colour in c("red", "green", "blue")) {
  print(colour)

}

## [1] "red"
## [1] "green"
## [1] "blue"
```

Example of a while loop:

```
# A while loop is similar in structure
# Start with a value of i = 1
i <- 1
# Loop while i is less than or equal to 100
while (i <= 5) {
  print(i);
  # Remember to increment i each loop, otherwise it will loop forever
  i<-i+1;
}

## [1] 1
## [1] 2
## [1] 3
## [1] 4
## [1] 5
```

Advantages of vectorisation

Why don't we use loops often in R? Because **vectorisation** is a fundamental approach in R. This means that a function can be applied to every element of a vector or matrix without having to explicitly loop through each element of the whole variable.

```
# vectorisation means that for loops are not quite as essential in R
# Generate the square root for every number from 1 to 100
sqrt(seq(100))

# Create a temporary vector, and multiply every element in it by pi
c(1,2,3,4,5,6,7,8,9,8,7,6,5,4,3,2,1) * pi

# Perform a log transformation on every element of our expression
# dataset and print out the first few lines
head(log10(gene.atlas))
```



Loops and Vectorisation

1. Experiment with the loops and vectorisation examples above.
2. Try creating a loop that counts in 2 down from 100 to 0.
3. Create a vector containing the numbers from 1 to 20. Use the concept of vectorisation to create a second vector containing the squares of those numbers.
4. **Challenge:** multiply *every third* element in the new vector from step 3 by 20. (*Hint:* try to combine your answer from step 2.)

Performing a vectorised function on a vector or matrix does not change the contents of that variable. To perform an in-place calculation, you need to explicitly overwrite the variable with the output of the function, e.g. `gene.atlas <- log10(gene.atlas)`.



5. Examine the values in `gene.atlas` after you have performed the `log10()` function as in the last line of the code chunk above.
 - a) Now try explicitly overwriting the `gene.atlas` object with the `log10()` values and examine the object again.
 - b) **Challenge:** if you are ahead, can you get the original `gene.atlas` values back? (*Hint:* $\log_b(x) = y$ and $b^y = x$)



Going Loopy

- Try to avoid `for` loops in production code, use them for testing ideas and exploring data.
- R objects are often vectors and can be analysed in bulk using normal mathematical and statistical transformations. This applies for the vectorised data in `data.frames` and `matrices`.
- `factors` cannot be used quite as easily in such analyses, so beware.
- `seq()` generates a regular sequence of numbers and is useful in testing.

5.3 Apply, Lapply

In the previous section we considered the simplicity with which we could perform simple transformations on vectorised data. In these examples a discrete transformation was applied to each value.

If we consider our chunk of gene expression data in the `gene.atlas` variable, we can imagine a number of simple transformations that we might wish to perform such as log2 transformations. There are more cases when we might wish to perform an analysis by row or column within the data. While this could be managed using a for loop to iterate over the data there are a number of simpler ways to access the data. One way to do this is with `apply()`.

`apply()` takes at least three arguments:

- the matrix of input data,
- the MARGIN, whether to perform the calculation by row (1) or column (2) and
- the function to apply to the margin

This function can be one that is predefined in R (such as `mean` or `sd` in the exercise below) or one you have created yourself as per the previous function section. In the latter case remember that the function should expect a vector as input and return a single value (e.g. it receives a vector of numbers and returns their mean).

The example below finds the standard deviation (`sd`) and the `mean` of each of the first few rows of the `gene.atlas` variable using `apply`.

```
| apply(head(gene.atlas), MARGIN=1, sd)
##   1007_s_at    1053_at     117_at    121_at  1255_g_at    1294_at
## 1026.28428   63.23077 193.86279 1642.96339 193.79413 286.80079

| apply(gene.atlas[,1:6], MARGIN=2, mean)
## GSM18865 GSM18866 GSM18867 GSM18868 GSM18869 GSM18870
## 387.8136 329.7084 285.1957 317.8131 266.4186 285.6001
```

	apply
	<ul style="list-style-type: none"> – <code>apply</code> is a crucial method for exploring data within the rows and columns of a <code>data frame</code> or <code>matrix</code>. – <code>lapply</code> is similar to <code>apply</code> but it returns results as a list rather than a matrix. – There is also the <code>dplyr</code> package, which is even more powerful but is beyond the scope of this workshop. See https://cran.r-project.org/web/packages/dplyr/dplyr.pdf

5.4 Basic data analysis methods

We have already been introduced to the `mean()` and `sd()` functions and as you might expect R provides many other numerical summary functions. These include:

Function	Description
<code>sum()</code>	total sum of a numeric vector or matrix
<code>mean()</code>	mean value of a numeric vector or matrix
<code>sd()</code>	standard deviation, if your object has missing values (NA) then see the <code>na.rm</code> option
<code>min()/max()</code>	minimum/maximum value in your object, also see the <code>na.rm</code> option
<code>range()</code>	range (min and max) for a numeric vector or matrix
<code>sqrt()</code>	is the square root of a numeric object. Can be applied to a vector or matrix, in which case it will return the results in vector or matrix form
<code>sample()</code>	select a specified number of items at random from a vector. Can be used without replacement, where an item can only be selected once (for example, to select five test subjects out of a pool of twenty) or with replacement, where items can be selected multiple times (as we did for the M&Ms exercise)

Function	Description
<code>duplicated()</code>	provides a vector of logicals indicating which elements of a vector have already been seen in that vector. In other words, the first time a value is seen it will return FALSE, and then if it occurs again it will return TRUE)
<code>unique()</code>	provides a non-redundant list of all values in a vector; any duplicated values will be output only the first time they are seen
<code>table()</code>	provides a frequency table of the counts of each combination of the

present values. Is generally used with `factor` data types or factor columns of a data frame. |



Data summary functions

Follow the commented code below to explore some of the above functions.

```
# Take the 101st line of the gene.atlas data as our test dataset
x <- as.numeric(gene.atlas[101,])

# Find various statistical summaries of that test data. The semicolon
# separates commands on the same line.
sum(x); mean(x); sd(x); sqrt(x)
min(x); max(x); range(x)

# randomly pick 50 numbers between 1 and 100 with replacement.
x <- sample.int(100, 50, replace=TRUE)

# since replace = TRUE in the line above, some numbers may have been picked more than once.
# duplicated() will tell us which elements have been well duplicated.
duplicated(x)

# This isn't very helpful, since it just gives us a list of TRUE and FALSE
# By taking a vector 'slice', we can get the values of the duplicated items
x[duplicated(x)]

# show the unique(or rather, non-redundant) numbers
unique(x)

# Show the number of times that any value is present in the data
# In this case, we're sampling 500 integers between 1 and 10
table(sample.int(10, 500, replace=TRUE))
```



Analytical methods

For this exercise, the numeric data taken from the `gene.atlas` data frame has been cast into a vector using the `as.numeric()` function. Slicing a data frame will usually return a data frame and some methods will only function with a vector.

Chapter 6

Data presentation and reporting

In the previous sections we have looked at a number of ways to import data, access data, summarise data. We should now consider how we best present data in a graphical way. R is endowed with a number of native and add-on packages that facilitate the preparation of figures, diagrams and graphs. In this section we will explore plots prepared using mainly the builtin functions and the `ggplot2` package.

`ggplot2` is not the simplest library to use to prepare plots at the start, but it is the most configurable and powerful. They also have very comprehensive documentation and examples on their website <http://docs.ggplot2.org/current/>.

6.1 Boxplots

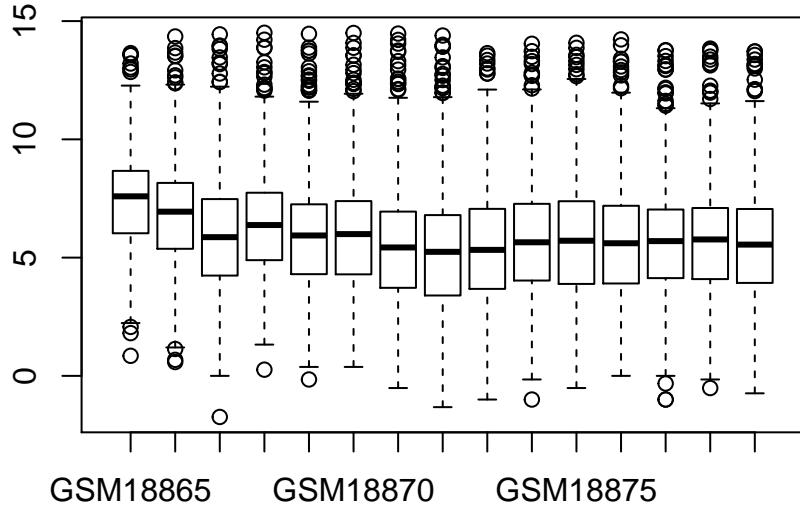
The boxplot is a widely used plot that can summarise the distribution of data within a collection. We routinely use boxplots to show a trend within data during e.g. signal intensity of DNA microarrays or read-counts for a single gene across different experimental condition.

6.1.1 Native

In this example we will prepare a boxplot for signal intensity across different arrays in our `gene.atlas` data set. The simplest way is to call the native `boxplot()` function. This produces a reasonable summary view of the data, but not one that would be suitable for publication purposes.

The command below generates a boxplot of 1000 gene expression values for the first 15 samples. Remember, `gene.atlas` contains over 22,000 expression values and 158 samples

```
| boxplot(log2(gene.atlas[sample(nrow(gene.atlas), 1000), seq(15)]))
```



There are other arguments you can use to add to the plot for example:

- `main='Expr distribution'` to add a title to the plot
 - `las=2` will rotate the x-axis labels to 90 degrees. (0=parallel, 1=horizontal, 2=perpendicular to the axis, 3=vertical)
 - `cex.axis = 0.8` to change the font size of the tick labels
 - `xlab='sample'` to label x-axes
 - `col=rainbow(15)` to add colours to the bars. `rainbow(15)` is a builtin colours palette in R, this returns 15 colours.
 - and more, use `?barplot` to use Help to find out more about the function.

6.1.2 ggplot2::geom_boxplot()

Using ggplot2 we can produce a much more attractive plot, although doing so is considerably more complicated.

```

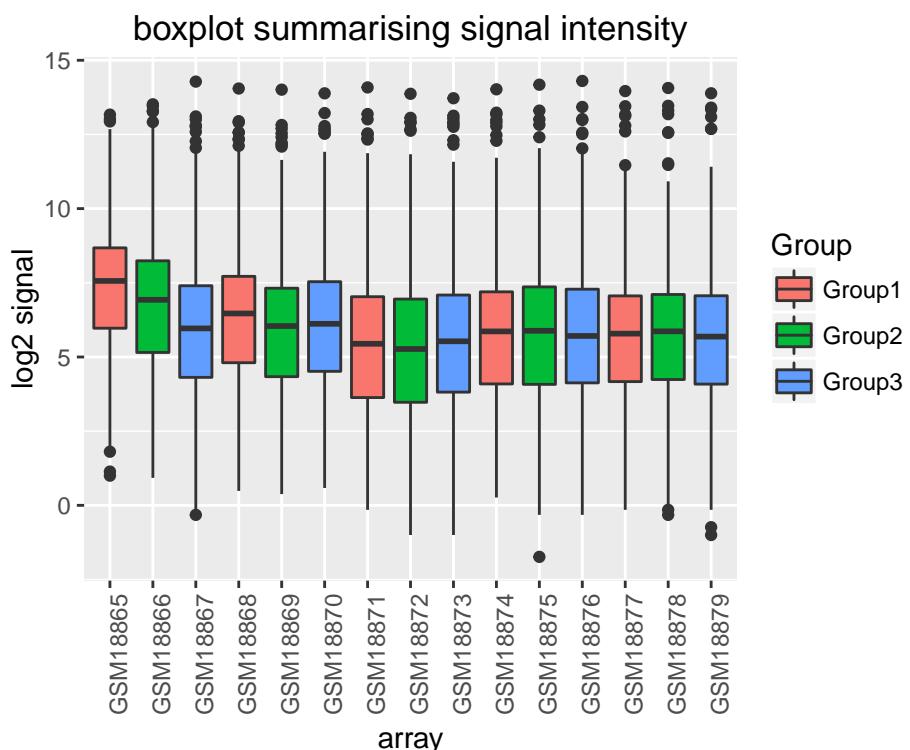
# Using ggplot, we first need to load the required libraries
library(ggplot2)
library(reshape2)

# Once again, get 1000 values for 15 samples. We need to transpose t() the
# dataset so that we have samples as rows and genes as columns.
atlas.subset <- as.data.frame(t(log2(gene.atlas[sample(nrow(gene.atlas), 1000),
                                         seq(15)])))

# We now need to transform the data from 'wide' to 'long' format
# Your instructor will explain this process
atlas.subset <- cbind(atlas.subset, Group=as.factor(paste0("Group", rep(seq(3), 5))),
                      Sample=rownames(atlas.subset))
atlas.melt <- melt(atlas.subset,
                     id.vars=c("Group", "Sample"),
                     variable.name="Gene",
                     value.name="Signal")

# Now that the data is in the correct format, we can draw the graph
# Again, the instructor will explain what is happening here
# The '+' symbols are required, and tell R that the command is not finished
p <- ggplot(atlas.melt, aes(factor(Sample), Signal, fill=Group)) +
  geom_boxplot() +
  labs(title="boxplot summarising signal intensity", x="array", y="log2 signal") +
  theme(axis.text.x=element_text(angle=90))
print(p)

```



6.2 Histograms

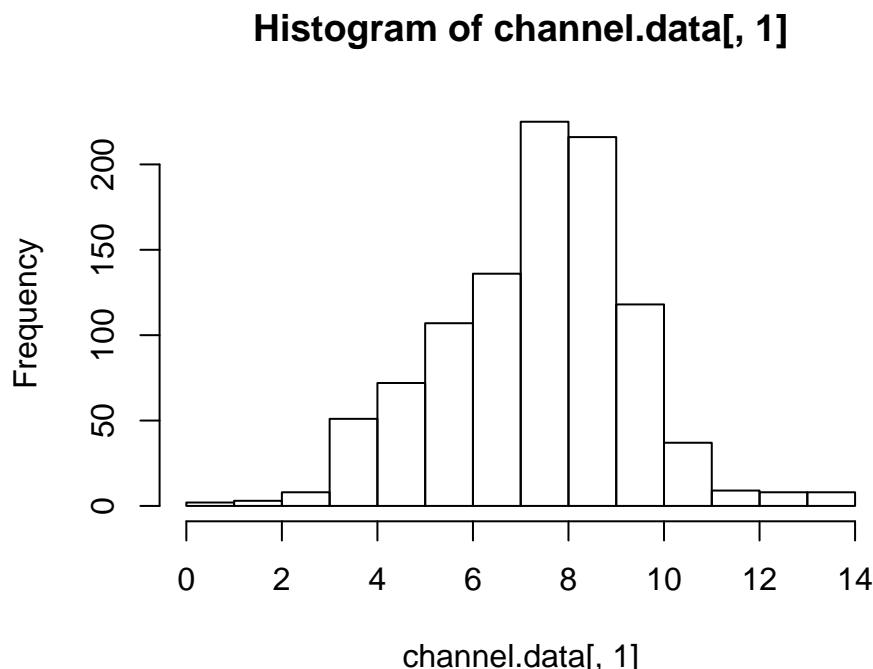
A histogram shows the distribution of data for typically a single sample.

6.2.1 Native

Let's use the same data that we used in the boxplot example but concentrating on the signal intensity for the first sample only.

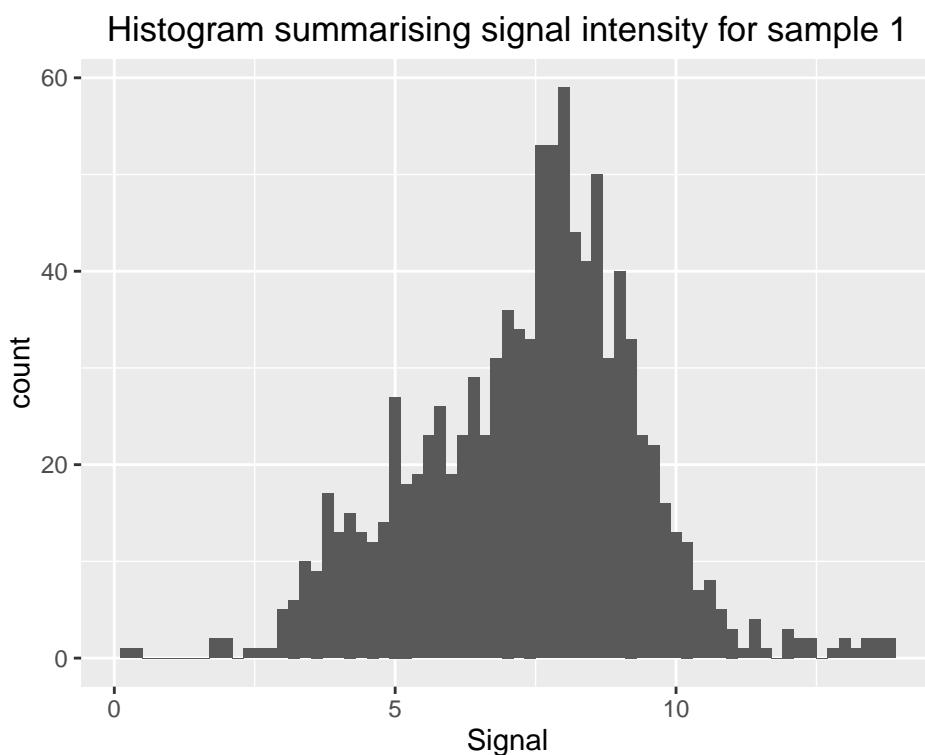
```
# The base graphics way
# Create a data frame slice containing our target expression data
channel.data <- data.frame(Signal=log2(gene.atlas[sample(nrow(gene.atlas),1000),1]))

# hist() is the built-in histogram function
# [, 1] selects the first (and in this case, only) column of data from channel.data
hist(channel.data[, 1])
```



6.2.2 ggplot2 - geom_histogram()

```
# Now with ggplot2. The data is already in long format, so no need to 'melt'
# this time. Also ggplot2 is already loaded, so no need to do that again.
p <- ggplot(channel.data, aes(x=Signal)) +
  geom_histogram(binwidth=0.2) +
  labs(title="Histogram summarising signal intensity for sample 1")
print(p)
```



 See also `plot(density())` for a more continuous view of the data distribution.

6.3 Bar charts

A barchart is superficially similar to a histogram in the bars of data are displayed. Bar charts are ideal for displaying counts associated with categorical data.

6.3.1 Native

The example chart here simply presents the number of diamonds in a database of diamonds that have been assigned to particular classes of cut. The `diamonds` data comes from one of the many datasets built into R that you can use for testing and experimentation. For a full list of these, use `library(help="datasets")`.

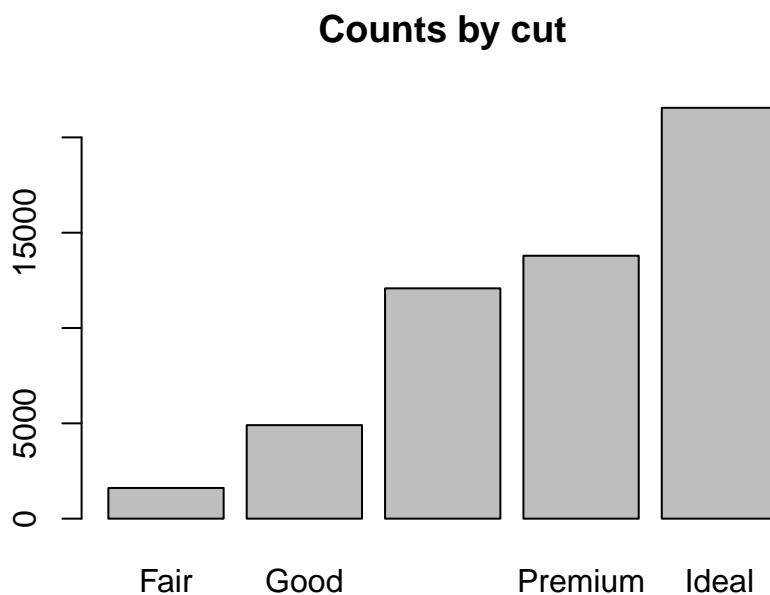
```
# Let's start by having a look at the data
head(diamonds)

## Source: local data frame [6 x 10]
##
##   carat      cut  color clarity depth table price     x     y     z
##   (dbl)    (fctr) (fctr) (fctr) (dbl)  (dbl) (int) (dbl) (dbl) (dbl)
## 1  0.23    Ideal     E    SI2  61.5    55    326  3.95  3.98  2.43
## 2  0.21  Premium     E    SI1  59.8    61    326  3.89  3.84  2.31
## 3  0.23     Good     E    VS1  56.9    65    327  4.05  4.07  2.31
## 4  0.29  Premium     I    VS2  62.4    58    334  4.20  4.23  2.63
## 5  0.31     Good     J    SI2  63.3    58    335  4.34  4.35  2.75
## 6  0.24  Very Good     J   VVS2  62.8    57    336  3.94  3.96  2.48
```

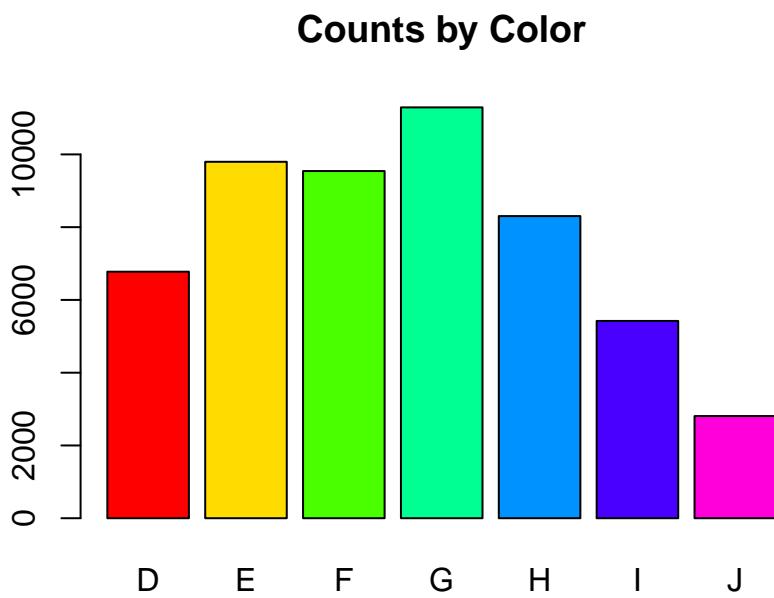
```
| str(diamonds)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 53940 obs. of 10 variables:  
## $ carat : num 0.23 0.21 0.23 0.29 0.31 0.24 0.24 0.26 0.22 0.23 ...  
## $ cut   : Ord.factor w/ 5 levels "Fair" < "Good" < ... : 5 4 2 4 2 3 3 3 1 3 ...  
## $ color  : Ord.factor w/ 7 levels "D" < "E" < "F" < "G" < ... : 2 2 2 6 7 7 6 5 2 5 ...  
## $ clarity: Ord.factor w/ 8 levels "I1" < "SI2" < "SI1" < ... : 2 3 5 4 2 6 7 3 4 5 ...  
## $ depth  : num 61.5 59.8 56.9 62.4 63.3 62.8 62.3 61.9 65.1 59.4 ...  
## $ table  : num 55 61 65 58 58 57 57 55 61 61 ...  
## $ price  : int 326 326 327 334 335 336 336 337 337 338 ...  
## $ x      : num 3.95 3.89 4.05 4.2 4.34 3.94 3.95 4.07 3.87 4 ...  
## $ y      : num 3.98 3.84 4.07 4.23 4.35 3.96 3.98 4.11 3.78 4.05 ...  
## $ z      : num 2.43 2.31 2.31 2.63 2.75 2.48 2.47 2.53 2.49 2.39 ...
```

```
| # Barplot using base version of barplot, grouping by the 'cut' category  
| barplot(table(diamonds[, "cut"]), main="Counts by cut")
```

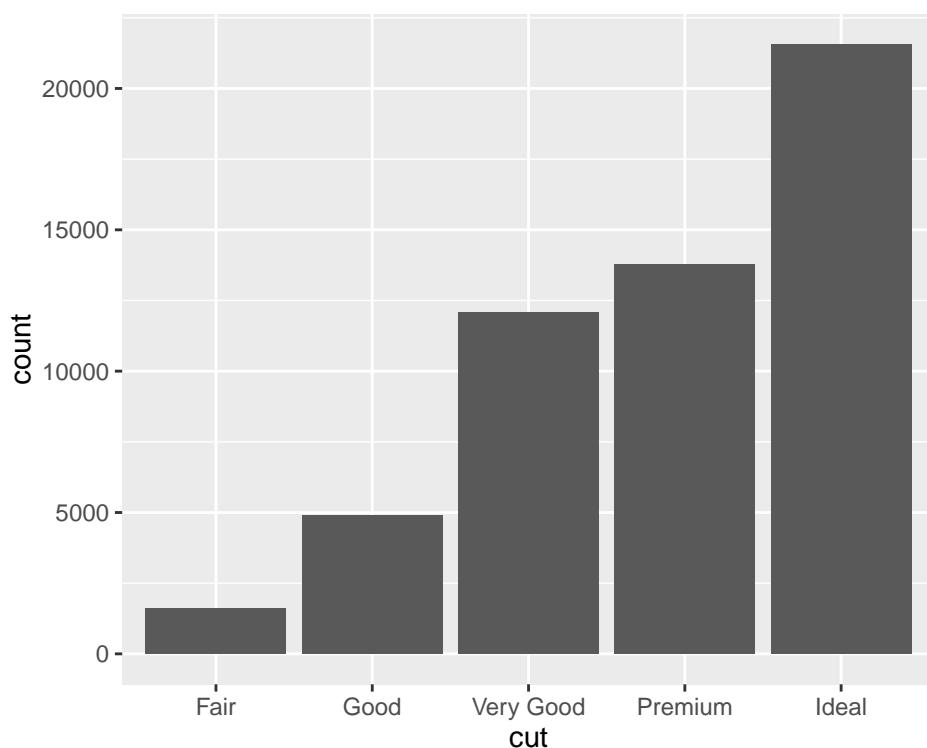


```
| barplot(table(diamonds[, "color"]), col=rainbow(7), main="Counts by Color")
```

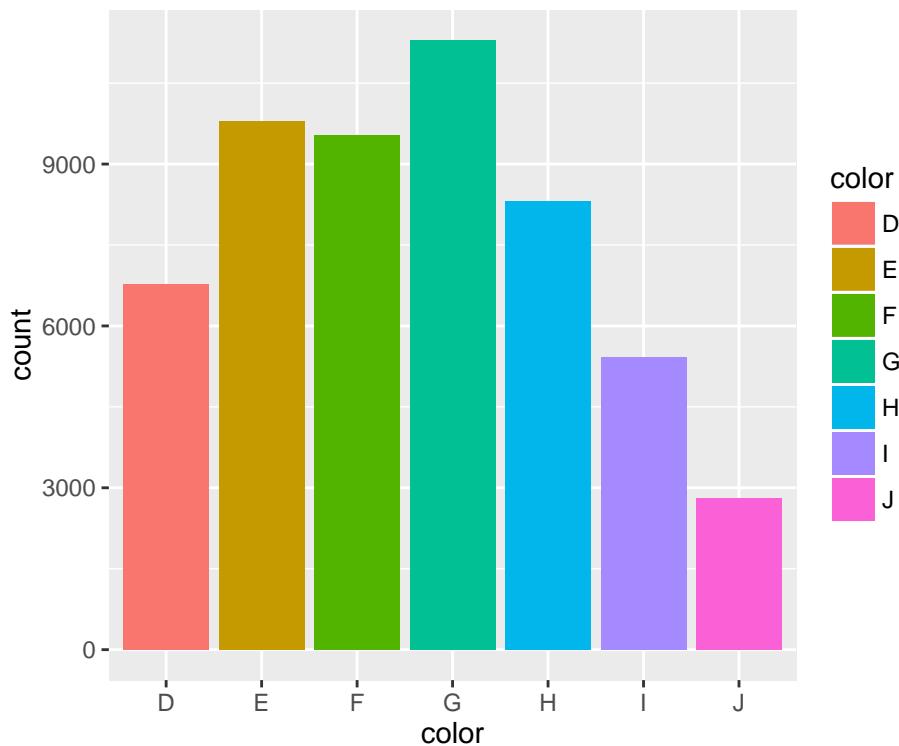


6.3.2 ggplot2::geom_bar()

```
# The ggplot version in this case can be quite simple. This is thanks to the
# data again being in long format, so not needing melting
p <- ggplot(diamonds, aes(cut)) + geom_bar()
print(p)
```



```
# Similarly, group by 'color' (and appropriately enough, let's colour it in
# this time)
p <- ggplot(diamonds, aes(color, fill = color)) + geom_bar()
print(p)
```



6.4 Scatterplots

Scatterplots are widely used for plotting data where an object has multiple variables. This is a quick way to look for any correlation in the dataset.

6.4.1 Native

The `iris` dataset, another built-in dataset in R, describes three varieties of the iris flower and a number of attributes that include the sepal and petal lengths and widths. In this case we are going to plot sepal length against width for three species of iris.

```
head(iris)

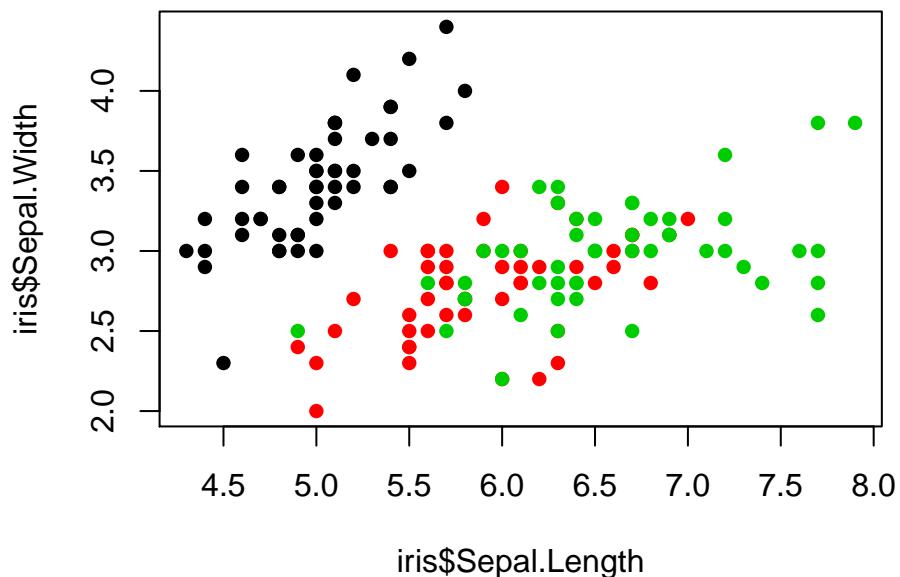
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          5.1         3.5          1.4         0.2  setosa
## 2          4.9         3.0          1.4         0.2  setosa
## 3          4.7         3.2          1.3         0.2  setosa
## 4          4.6         3.1          1.5         0.2  setosa
## 5          5.0         3.6          1.4         0.2  setosa
## 6          5.4         3.9          1.7         0.4  setosa
```

```
str(iris)

## 'data.frame': 150 obs. of 5 variables:
## $ Sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
## $ Sepal.Width : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
```

```
## $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
## $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
## $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 ...
```

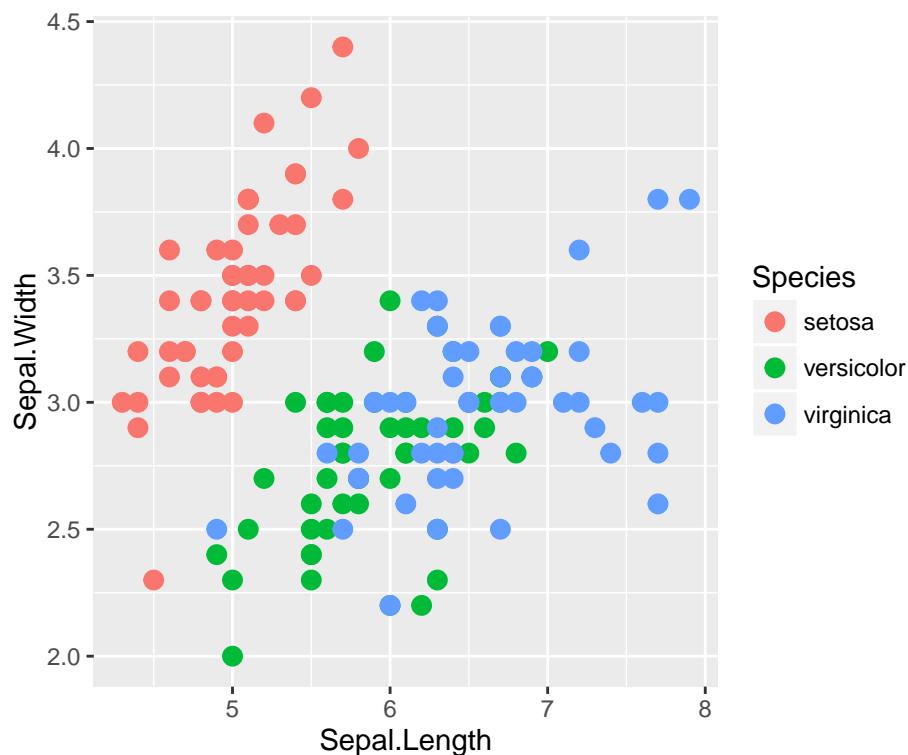
```
| plot(iris$Sepal.Length, iris$Sepal.Width, col=iris$Species, pch=16)
```



The `pch` option specifies the type of symbol that can be used in a plot. Search `?pch` which will bring up the Help documentation and the list of symbols that are available.

6.4.2 ggplot2::geom_point()

```
# With ggplot, we can specify that we want to colour the points by Species
p <- ggplot(iris, aes(Sepal.Length, Sepal.Width)) +
  geom_point(aes(color = Species), size=3)
print(p)
```



Now the variation in sepal shape between the species is clear



Plotting with R

Work your way through the previous code examples to draw boxplots, histograms, bar charts and scatter plots as described. In each case, the built-in plot function is described first, then the ggplot2 version.

Appendix 1

Summary of Object Types

Type	Description
vectors	ordered collection of numeric, character, complex and logical values.
factors	special type vectors with grouping information of its components
data	two dimensional structures with different data types
frames	
matrices	two dimensional structures with data of same type
arrays	multidimensional arrays of vectors
lists	general form of vectors with different types of elements
functions	piece of code

More information on R and Bioconductor for Genomics

- Thomas Girke's manuals and guides to R <http://manuals.bioinformatics.ucr.edu/home>
- <http://www.r-bloggers.com/>

Appendix 2 - Solutions to exercises

2.4.1 Using vectors

`mixedVector[2]` returns the name and value of the second element of the vector.

`mixedVector[[2]]` just returns the value

```
| mixedVector[['meaningOfLife']]
```

2.4.3 Using lists

```
# Retrieve second and fourth elements
mixedList[c(2,4)]

# Retrieve everything BUT 2nd and 4th elements
mixedList[c(-2,-4)]

# Create a `subject` variable holding the names, sex and age
names <- c("A", "B", "C", "D", "E", "F", "G", "H", "I", "J")
sex <- c(TRUE, TRUE, TRUE, TRUE, TRUE, FALSE, FALSE, FALSE, FALSE)
age <- c(1,2,3,4,5,6,7,8,9,10)
subjects <- list(names=names, sex=sex, age=age)
```

2.4.4 Adding factor values

```
levels(mandm) <- c(levels(mandm), "brown")

# Change vector values 100 to 200 to brown
mandm[100:200] = 'brown'
```

2.4.6 Creating a data frame

```
subject.data <- as.data.frame(subjects)
colnames(subject.data)

# If they need naming
colnames(subject.data) <- c('names', 'sex', 'age')
```

4.2 Reading in the GeneAtlas data

```
# note the the comment/metadata delimiter is the ! symbol
# the column values are separated by tabs
# the column names are included in the first row - we can import Header
head(read.table(gene.atlas.file, comment.char="!", sep="\t", header=TRUE))
gene.atlas <- read.table(gene.atlas.file, comment.char="!", sep="\t", header=TRUE)

# read the data again, this time set the row.names with the first data column
gene.atlas <- read.table(gene.atlas.file, comment.char="!", sep="\t",
                           header=TRUE, row.names=1)
```

5.1 Creating and using functions

```
# A function to multiply or divide
multiplyFunction <- function(a, b=2, multiply=TRUE) {
  if (multiply) {
    a*b
  } else {
    a/b
  }
}

# A function to return a letter at a certain position
letterFunction <- function(word, position) {
  toupper(substr(word, position, position))
}

# A function to generate a filled matrix
makeMatrix <- function (rows, columns, numbers = TRUE) {
  if (numbers) {
    matrix(seq(rows*columns), nrow = rows, ncol = columns)
  } else {
    matrix(LETTERS[1:(rows*columns)], nrow = rows, ncol = columns)
  }
}
```

5.2 Loops and vectorisation

```
# Count down from 100 in 2s
i <- 100
while (i >= 0) {
  print(i)
  i <- i-2
}

# Square of the numbers from 1 to 20
startVec <- seq(20)
squareVec <- startVec^2
```

6.5 Venn diagram

```
| str(Vstem)
| Vstem@IntersectionSets$"1111"
```


Appendix 3 - The SessionInfo for this version of

```
| sessionInfo()
## R version 3.3.0 (2016-05-03)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 14393)
##
## locale:
## [1] LC_COLLATE=English_Australia.1252  LC_CTYPE=English_Australia.1252
## [3] LC_MONETARY=English_Australia.1252 LC_NUMERIC=C
## [5] LC_TIME=English_Australia.1252
##
## attached base packages:
## [1] stats      graphics   grDevices utils      datasets   methods    base
##
## other attached packages:
## [1] tidyR_0.4.1          BiocInstaller_1.22.3 bookdown_0.1.4
## [4] RevoUtilsMath_3.2.5  reshape2_1.4.1    ggplot2_2.1.0
## [7] limma_3.28.21
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.5           formatR_1.4        plyr_1.8.3
## [4] tools_3.3.0           digest_0.6.9       evaluate_0.9
## [7] tibble_1.1             gtable_0.2.0       shiny_0.13.2
## [10] DBI_0.4-1            rstudioapi_0.5    yaml_2.1.13
## [13] parallel_3.3.0        stringr_1.0.0     dplyr_0.4.3
## [16] knitr_1.13            grid_3.3.0        R6_2.1.2
## [19] rmarkdown_1.0.9005    magrittr_1.5      scales_0.4.0
## [22] codetools_0.2-14     htmltools_0.3.5   assertthat_0.1
## [25] mime_0.4              xtable_1.8-2      colorspace_1.2-6
## [28] httpuv_1.3.3         labeling_0.3     stringi_1.1.1
## [31] miniUI_0.1.1         munsell_0.4.3
```