

A Little Book of R for Bioinformatics 2.0

Avirl Coghlan, with contributions by Nathan L. Brouwer

2021-07-14

Contents

Preface to version 2.0	7
1 How to review this book	9
1.1 File format	9
1.2 What can you do?	9
1.3 Adding comments	9
2 Test shiny app	11
3 Test webpage insertion	13
4 Datacamp test	15
5 Downloading R	17
5.1 Preface	17
5.2 Introduction to R	17
5.3 Installing R	17
5.4 Starting <i>R</i>	18
6 Installing the RStudio IDE	19
6.1 Getting to know RStudio	19
6.2 RStudio versus RStudio Cloud	19
7 Installing <i>R</i> packages	21
7.1 Downloading packages with the RStudio IDE	21
7.2 Downloading packages with the function <code>install.packages()</code>	21
7.3 Using packages after they are downloaded	22
8 Installing Bioconductor	23
8.1 Bioconductor	23
8.2 Installing BiocManager	23
8.3 The ins and outs of package installation	24
8.4 Actually loading a package	25
9 A Brief introduction to R	27
9.1 Vocabulary	27
9.2 R functions	27
9.3 Interacting with R	28
9.4 Variables in R	28
9.5 Arguments	30
9.6 Help files with <code>help()</code> and <code>?</code>	31
9.7 Searching for functions with <code>help.search()</code> and <code>RSiteSearch()</code>	31
9.8 More on functions	31

9.9	Quitting R	32
9.10	Links and Further Reading	32
10	DNA descriptive statics - Part 1	33
10.1	Preface	33
10.2	Writing TODO:	33
10.3	Introduction	33
10.4	Vocabulary	33
10.5	Functions	33
10.6	Preliminaries	33
10.7	Converting DNA from FASTA format	34
10.8	Length of a DNA sequence	34
10.9	Acknowledgements	36
11	Programming in R: for loops	39
11.1	Preface	39
11.2	Vocab	39
11.3	Functions	39
11.4	Basic for loops in in R	39
11.5	Challenge: complicated vectors of values	40
12	Mini tutorial: Vectors in R	41
12.1	Preface	41
12.2	Vocab	41
13	Functions	43
13.1	Vectors in R	43
13.2	Math on vectors	43
13.3	Functions on vectors	44
13.4	Operations with two vectors	44
13.5	Subsetting vectors	45
13.6	Sequences of numbers	45
13.7	Vectors can hold numeric or character data	46
13.8	Regular expressions can modify character data	46
14	Plotting vectors in base R	47
14.1	Preface	47
14.2	Plotting numeric data	47
14.3	Other plotting packages	48
15	Programming in R: functions	49
15.1	Preface	49
15.2	Vocab	49
15.3	Functions	49
15.4	Functions in R	49
15.5	Comments in R	50
16	FASTA Files	51
16.1	Example FASTA file	52
16.2	Multiple sequences in a single FASTA file	52
16.3	Multiple sequence alignments can be stored in FASTA format	52
16.4	FASTQ Format	52
17	Downloading DNA sequences as FASTA files in R	55
17.1	DNA Sequence Statistics: Part 1	56

17.2	OPTIONAL: Saving FASTA files	59
17.3	Next steps	59
18	Downloading DNA sequences as FASTA files in R	61
18.1	Preliminaries	61
18.2	Convert FASTA sequence to an R variable	61
19	Downloading protein sequences in R	65
19.1	Preliminaries	65
19.2	Retrieving a UniProt protein sequence using rentrez	65
20	Sequence dotplots in R	67
20.1	Preliminaries	67
20.2	Visualizing two identical sequences	67
20.3	Visualizing repeats	67
20.4	Inversions	68
20.5	Translocations	68
20.6	Random sequence	68
20.7	Comparing two real sequences using a dotplot	68
21	Global proteins alignments in R	71
21.1	Preliminaries	71
21.2	Pairwise global alignment of DNA sequences using the Needleman-Wunsch algorithm	71
21.3	Pairwise global alignment of protein sequences using the Needleman-Wunsch algorithm	73
21.4	Aligning UniProt sequences	74
21.5	Viewing a long pairwise alignment	74
22	Local protein alignments in R	77
22.1	Preliminaries	77
22.2	Pairwise local alignment of protein sequences using the Smith-Waterman algorithm	77
23	Retrieving multiple sequences in R	79
23.1	Preliminaries	79
23.2	Retrieving a set of sequences from UniProt	79
23.3	Downloading sequences in bulk	80
24	Multiple sequence alignment in R	83
24.1	Preliminaries	83
24.2	Multiple sequence alignment (MSA)	83
24.3	Make MSA with msa()	83
24.4	Viewing your MSA	85
24.5	Discarding very poorly conserved regions from an alignment	86
25	Calculating genetic distances between sequences	87
25.1	Preliminaries	87
25.2	Introduction	87
25.3	Calculating genetic distances between DNA/mRNA sequences	88
25.4	Calculating genetic distance	89
26	Unrooted neighbor-joining phylogenetic trees	91
26.1	Preliminaries	91
26.2	Building an unrooted phylogenetic tree for protein sequences	91
26.3	Bootstrap values indicate support for clades	92
26.4	Branch lengths indicate divergence between sequences	93
26.5	Unrooted trees lack an outgroup	93

27 A complete bioinformatics workflow in R	95
27.1 Software Preliminaires	95
27.2 Downloading macro-molecular sequences	96
27.3 Prepping macromolecular sequences	97
27.4 Aligning sequences	98
27.5 The shroom family of genes	98
27.6 Downloading multiple sequences	99
27.7 Multiple sequence alignment	99
27.8 Genetic distance.	100
27.9 Phylogenetic trees (finally!)	101
 I Appendices	 103
Appendix 01: Getting access to R	107
27.10Getting Started With R and RStudio	107
 Getting started with R itself (or not)	 109
Vocabulary	109
R commands	109
27.11Help!	111
27.12Other features of RStudio	112
27.13Practice (OPTIONAL)	112

Preface to version 2.0

Welcome to *A Little Book of R for Bioinformatics 2.0!*.

This book is based on the original *A Little Book of R for Bioinformatics* by Dr. Avril Coghlan (Hereafter “ALBRB 1.0”). Dr. Coghlan’s book was one of the first and most thorough introductions to using R for bioinformatics, and was generously published under the Creative Commons 3.0 Attribution License (CC BY 3.0). In addition to describing how to do bioinformatics in R, Coghlan provided numerous functions to facilitate important tasks, practice questions, and references to further reading.

ALBRB 1.0 was extremely useful to me when I was learning bioinformatics and computational biology. In this version of the book, which I’ll refer to as ALBRB 2.0, I have adapted Dr. Coghlan’s original book to suit my own teaching needs.

Below I’ve outlined the general types of changes I’ve made to the original book. I have tried to link back to the original content that these updates are derived from and note how changes were made. Any errors or inconsistencies should be ascribed to me, not Dr. Coghlan. If you have any feedback, please email me at brouwer@gmail.com

Nathan Brouwer, June 2021

Changes implemented in ALBRB 2.0 by Nathan Brouwer

1. Converted the entire book to RMarkdown and published it via bookdown.
2. Added instructions for using RStudio and RStudio Cloud.
3. Updated instructions to reflect any changes in software, including changes to how the bioinformatics repository Bioconductor now works.
4. Split up chapters into smaller units.
5. Reorganized the order of some material.
6. Added links to the book I am developing, *Computational Biology for All*.
7. Moved most functions and datasets to my teaching package `compbio4all`.
8. Changed some plotting to `ggplot2` or `ggpubr`.
9. Added additional subheadings
10. Added vocab and function lists to the beginning of many chapters
11. At times replaced non-biological examples with biological ones.
12. Change from British to American English (Sorry! Couldn’t help myself.)
13. Provided additional links to external resources.
14. Added use of `rentrez` for querying NCBI databases

Chapter 1

How to review this book

1.1 File format

These lessons are written in RStudio using RMarkdown. Each .Rmd file is a mix of text, written in plain format, and **code chunks**, which look like this

Code chunks start with three apostrophes and `{r}`, like this: `“‘{r}`. They end with three apostrophes`”`. They will appear gray when opened up in RStudio but be white in the normal R code editor or other text editor.

1.2 What can you do?

- Read and fix typos :)
- Add comments within the file
- Email me general comments about the files (structure, topics, confusing parts etc)

1.3 Adding comments

You're welcome to add comments anywhere to the files, do the exercises and type up the key, propose your own exercises, etc.

1.3.1 HTML-tagged comments

The easiest way is to type a comment into the normal text part of the .Rmd file and then surround it with an html comment tag. A comment saying "A comment" will therefore look like this:

In RStudio you can type up a comment, highlight it then hit Shift+Control+C on a PC or Shift+Command+C on a mac.

1.3.2 Code chunk comments

Another way to add a comment is to make a RMarkdown code chunk then type up your comments in it by commenting out each line, like the one below. Key to doing this is to put `eval = F`, `echo = F` in the braces after the `r`. `eval = F`, `echo = F` tells RStudio to leave that alone when the .Rmd file gets rendered into a web page or PDF.

In RStudio you can add code chunks with a shortcut key. On a Mac the shortcut is `OPTION + COMMAND + I`.

1.3.3 Keys to exercise

Some files will have exercises at the end. I don't always include the key, and you're welcome to try the exercise and type up a key (or fix any errors in mine). As for code chunk comments include `eval = F`, `echo = F` in the braces so the key won't appear when rendered. So a problem and its key would look something like this

1.3.4 Problem - fix this code

```
print(correct answer)
```

Chapter 2

Test shiny app

Here is an example of a Shiny app embedded in a bookdown book.

```
knitr::include_app(url = "https://wrightaprilm.shinyapps.io/treesiftr_app/")
```



Please Wait






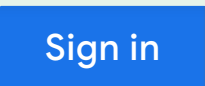
Chapter 3

Test webpage insertion

```
knitr::include_url(url ="https://docs.google.com/spreadsheets/d/18dt7cYHAaszV5Iq8Y-7Y0EFnv6hIXbD39Vn75-
```

 **test**  Saved to Drive



This version of Safari is no longer supported. Please upgrade to a [supported browser](#). [Dismiss](#)

  100%  View only 

A1

fx

	A	B	C	D
1				
2				
3		x	2	
4		y	4	
5				

 Sheet1 

Chapter 4

Datacamp test

By default, `tutorial` will convert all R chunks.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJhIDwtIDJcbmIgPC0gM1xuXG5hICsgYiJ9

Chapter 5

Downloading R

By: Avril Coghlan

Adapted, edited and expanded: Nathan Brouwer (brouwern@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

5.1 Preface

The following introduction to *R* is based on the first part of “How to install *R* and a Brief Introduction to *R*” by Avril Coghlan, which was released under the Creative Commons 3.0 Attribution License (CC BY 3.0). For additional information see the Appendices and “Getting *R* onto your computer”.

5.2 Introduction to R

R (www.r-project.org) is a commonly used free statistics software. *R* allows you to carry out statistical analyses in an interactive mode, as well as allowing programming.

5.3 Installing R

To use *R*, you first need to install the *R* program on your computer.

5.3.1 Installing *R* on a Windows PC

These instructions will focus on installing *R* on a Windows PC. However, I will also briefly mention how to install *R* on a Macintosh or Linux computer (see below).

These steps have not been checked as of 8/13/2019 so there may be small variations in what the prompts are. Installing R, however, is basically that same as any other program. Clicking “Yes” etc on everything should work.

PROTIP: Even if you have used *R* before its good to regularly update it to avoid conflicts with recently produced software.

Minor updates of *R* are made very regularly (approximately every 6 months), as *R* is actively being improved all the time. It is worthwhile installing new versions of *R* a couple times a year, to make sure that you have a recent version of *R* (to ensure compatibility with all the latest versions of the *R* packages that you have downloaded).

To install *R* on your **Windows** computer, follow these steps:

1. Go to <https://cran.r-project.org/>
2. Under “Download and Install R”, click on the “Windows” link.
3. Under “Subdirectories”, click on the “**base**” link.
4. On the next page, you should see a link saying something like “Download *R* 4.1.0 for Windows” (or *R* X.X.X, where X.X.X gives the version of the program). Click on this link.
5. You may be asked if you want to save or run a file “R-x.x.x-win32.exe”. Choose “Save” and save the file. Then double-click on the icon for the file to run it.
6. You will be asked what language to install it in.
7. The *R* Setup Wizard will appear in a window. Click “Next” at the bottom of the *R* Setup wizard window.
8. The next page says “Information” at the top. Click “Next” again.
9. The next page says “Select Destination Location” at the top. By default, it will suggest to install *R* on the C drive in the “Program Files” directory on your computer.
10. Click “Next” at the bottom of the *R* Setup wizard window.
11. The next page says “Select components” at the top. Click “Next” again.
12. The next page says “Startup options” at the top. Click “Next” again.
13. The next page says “Select start menu folder” at the top. Click “Next” again.
14. The next page says “Select additional tasks” at the top. Click “Next” again.
15. *R* should now be installing. This will take about a minute. When *R* has finished, you will see “Completing the *R* for Windows Setup Wizard” appear. Click “Finish”.
16. To start *R*, you can do one of the following steps:
17. Check if there is an “R” icon on the desktop of the computer that you are using. If so, double-click on the “R” icon to start *R*. If you cannot find an “R” icon, try the next step instead.
18. Click on the “Start” button at the bottom left of your computer screen, and then choose “All programs”, and start *R* by selecting “R” (or *R* X.X.X, where X.X.X gives the version of *R*) from the menu of programs.
19. The *R* console (a rectangle) should pop up:

5.3.2 How to install *R* on non-Windows computers (eg. Macintosh or Linux computers)

These steps have not been checked as of 8/13/2019 so there may be small variations in what the prompts are. Installing *R*, however, is basically that same as any other program. Clicking “Yes” etc on everything should work.

The instructions above are for installing *R* on a Windows PC. If you want to install *R* on a computer that has a non-Windows operating system (for example, a Macintosh or computer running Linux, you should download the appropriate *R* installer for that operating system at <https://cran.r-project.org/> and follow the *R* installation instructions for the appropriate operating system at https://cran.r-project.org/doc/FAQ/R-FAQ.html#How-can-R-be-installed_003f.

5.4 Starting *R*

To start *R*, Check if there is an *R* icon on the desktop of the computer that you are using. If so, double-click on the *R* icon to start *R*. If you cannot find an *R* icon, try the next step instead.

You can also start *R* from the Start menu in Windows. Click on the “Start” button at the bottom left of your computer screen, and then choose “All programs”, and start *R* by selecting “R” (or *R* X.X.X, where X.X.X gives the version of *R*, e.g.. *R* 2.10.0) from the menu of programs.

Say “Hi” to *R* and take a quick look at how it looks. Now say “Goodbye”, because we will never actually do any work in this version of *R*; instead, we’ll use the **RStudio IDE (integrated development environment)**.

Chapter 6

Installing the RStudio IDE

By: Nathan Brouwer

The name “R” refers both to the programming language and the program that runs that language. When you download *itR** there is also a basic **GUI** (graphical user interface) that you can access via the *R* icon.

Other GUIs are available, and the most popular currently is **RStudio**. RStudio a for-profit company that is a main driver of development of R. Much of what they produce has free basic versions or is entirely free. They produce software (RStudio), cloud-based applications (**RStudio Cloud**), and web server infrastructure for business applications of R.

A brief overview of installing RStudio can be found here “Getting RStudio on to your computer”

6.1 Getting to know RStudio

For a brief overview of RStudio see “Getting started with RStudio”

A good overview of what the different parts of RStudio can be seen in the image in this tweet: <https://twitter.com/RLadiesNCL/status/1138812826917724160?s=20>

6.2 RStudio versus RStudio Cloud

RStudio and RStudio cloud work almost identically, so anything you read about RStudio will apply to RStudio Cloud. RStudio is easy to download and use, but RStudio Cloud eliminates even the minor hiccups that occur. Free accounts with RStudio Cloud allow up to 15 hours per month, which is enough for you to get a taste for using R.

Chapter 7

Installing *R* packages

By: Avril Coghlan.

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

R is a programming language, and **packages** (aka **libraries**) are bundles of software built using *R*. Most sessions using *R* involve using additional *R* packages. This is especially true for bioinformatics and computational biology.

NOTE: If you are working in an RStudio Cloud environment organized by someone else (e.g. a course instructor), they likely are taking care of many of the package management issues. The following information is still useful to be familiar with.

7.1 Downloading packages with the RStudio IDE

There is a point-and-click interface for installing *R* packages in RStudio. There is a brief introduction to downloading packages on this site: <http://web.cs.ucla.edu/~gulzar/rstudio/>

I've summarized it here:

1. “Click on the”Packages” tab in the bottom-right section and then click on “Install”. The following dialog box will appear.
2. In the “Install Packages” dialog, write the package name you want to install under the Packages field and then click install. This will install the package you searched for or give you a list of matching package based on your package text.

7.2 Downloading packages with the function `install.packages()`

The easiest way to install a package if you know its name is to use the *R* function `install.packages()`. Note that it might be better to call this “download.packages” since after you install it, you also have to load it!

Frequently I will include `install.packages(...)` at the beginning of a chapter the first time we use a package to make sure the package is downloaded. Note, however, that if you already have downloaded the package, running `install.packages(...)` will download a new copy. While packages do get updated from time to time, but its best to re-run `install.packages(...)` only occassionaly.

We'll download a package used for plotting called `ggplot2`, which stands for “Grammar of Graphics.” `ggplot2` was developed by Dr. Hadley Wickham, who is now the Chief Scientists for RStudio.

To download `ggplot2`, run the following command:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpbmN0YWxsLnBhY2thZ2VzKFwiZ2dwbG90MlwiKSAjIG5vdGUgdGhlIFwiIF
```

Often when you download a package you'll see a fair bit of angry-looking red text, and sometime other things will pop up. Usually there's nothing of interest here, but sometimes you need to read things carefully over it for hints about why something didn't work.

7.3 Using packages after they are downloaded

To actually make the functions in package accessible you need to use the `library()` command. Note that this is *not* in quotes.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5J5KGdncGxvdDIpICMgdm90ZTogTk8gXCIgXCIifQ==
```

Chapter 8

Installing Bioconductor

By: Avril Coghlan.

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0), including details on install Bioconductor and common prompts and error messages that appear during installation.

8.1 Bioconductor

R **packages** (aka “libraries”) can live in many places. Most are accessed via **CRAN**, the **Comprehensive R Archive Network**. The bioinformatics and computational biology community also has its own package hosting system called Bioconductor. *R* has played an important part in the development and application of bioinformatics techniques in the 21th century. Bioconductor 1.0 was released in 2002 with 15 packages. As of winter 2021, there are almost 2000 packages in the current release!

NOTE: If you are working in an RStudio Cloud environment organized by someone else (eg a course instructor), they likely are taking care of most of package management issues, including setting up Bioconductor. The following information is still useful to be familiar with.

To interface with Bioconductor you need the BiocManager package. The Bioconductor people have put BiocManager on CRAN to allow you to set up interactions with Bioconductor. See the BiocManager documentation for more information (<https://cran.r-project.org/web/packages/BiocManager/vignettes/BiocManager.html>).

Note that if you have an old version of R you will need to update it to interact with Bioconductor.

8.2 Installing BiocManager

BiocManager can be installed using the `install.packages()` packages command.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpbmN0YWxsLnBhY2thZ2VzKFwiQmlvY01hbmFnZXJcIikgYBSZW1lbWJlcil
```

Once downloaded, BioManager needs to be explicitly loaded into your active R session using `library()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KEJpb2NNYW5hZ2VzKSAjIG5vIHF1b3Rlc3sgYWdhaW4sIGlnb
```

Individual Bioconductor packages can then be downloaded using the `install()` command. An essential packages is `Biostrings`. To do this ,

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJCZW9jTWFuYWdlcj06aW5zdGFsbChcIkJpb3N0cmduZ3NcIikifQ==
```

8.3 The ins and outs of package installation

IMPORANT Bioconductor has many **dependencies** - other packages which it relies on. When you install Bioconductor packages you may need to update these packages. If something seems to not be working during this process, restart R and begin the Bioconductor installation process until things seem to work.

Below I discuss the series of prompts I had to deal with while re-installing Biostrings while editing this chapter.

8.3.1 Updating other packages when downloading a package

When I re-installed Biostrings while writing this I was given a HUGE blob of red text that contained this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiInZ2V0T3B0aW9uKFWicmVwb3NcIiknIHJlcGxhY2VzIEJpb2NvbmlR1Y3RvcilE
```

Hidden at the bottom was a prompt: “Update all/some/none? [a/s/n]:”

It's a little vague, but what it wants me to do is type in a, s or n and press enter to tell it what to do. I almost always chose “a”, though this may take a while to update everything.

8.3.2 Packages “from source”

You are likely to get lots of random-looking feedback from R when doing Bioconductor-related installations. Look carefully for any prompts as the very last line. While updating Biostrings I was told: “*There are binary versions available but the source versions are later:*” and given a table of packages. I was then asked “*Do you want to install from sources the packages which need compilation? (Yes/no/cancel)*”

I almost always chose “no”.

8.3.3 More on angry red text

After the prompt about packages from source, R proceeded to download a lot of updates to packages, which took a few minutes. Lots of red text scrolled by, but this is normal.


```

Console Terminal x R Markdown x Jobs x
~/google_backup_sync_nlb24/lbrb/ ↗
trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/tidyselect_1.
Content type 'application/x-gzip' length 198800 bytes (194 KB)
=====
downloaded 194 KB

trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/tidytree_0.3.
Content type 'application/x-gzip' length 220977 bytes (215 KB)
=====
downloaded 215 KB

trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/tinytex_0.32.
Content type 'application/x-gzip' length 121221 bytes (118 KB)
=====
downloaded 118 KB

trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/tufte_0.10.tg
Content type 'application/x-gzip' length 269000 bytes (262 KB)
=====
downloaded 262 KB

trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/UniprotR_2.0.
Content type 'application/x-gzip' length 236776 bytes (231 KB)
=====
downloaded 231 KB

trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.0/units_0.7-2.t
Content type 'application/x-gzip' length 1260870 bytes (1.2 MB)
=====

```

8.4 Actually loading a package

Again, to actually load the `Biostrings` package into your active R sessions requires the `library()` command:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KEJpb3N0cmIuZ3MpIn0=
```

As you might expect, there's more red text scrolling up my screen!

```

~/google_backup_sync_nlb24/lbrb/ ↗
The following objects are masked from 'package:base':

  anyDuplicated, append, as.data.frame, basename, cbind, colnames, dirname,
  do.call, duplicated, eval, evalq, Filter, Find, get, grep, grepl, intersect,
  is.unsorted, lapply, Map, mapply, match, mget, order, paste, pmax, pmax.int,
  pmin, pmin.int, Position, rank, rbind, Reduce, rownames, sapply, setdiff,
  table, tapply, union, unique, unsplit, which, which.max, which.min

Loading required package: S4Vectors
Loading required package: stats4

Attaching package: 'S4Vectors'

The following object is masked from 'package:base':

  expand.grid

Loading required package: IRanges
Loading required package: XVector

Attaching package: 'Biostrings'

The following object is masked from 'package:base':

  strsplit

> |

```

I can tell that is actually worked because at the end of all the red stuff is the R prompt of “>” and my cursor.



Chapter 9

A Brief introduction to R

By: Avril Coghlan.

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

This chapter provides a brief introduction to R. At the end of are links to additional resources for getting started with R.

9.1 Vocabulary

- scalar
- vector
- list
- class
- numeric
- character
- assignment
- elements of an object
- indices
- attributes of an object
- argument of a function

9.2 R functions

- `<-`
- `[]`
- `$`
- `table()`
- `function`
- `c()`
- `log10()`
- `help()`, `?`
- `help.search()`
- `RSiteSearch()`
- `mean()`
- `return()`
- `q()`

9.3 Interacting with R

You will type *R* commands into the RStudio **console** in order to carry out analyses in *R*. In the RStudio console you will see the R prompt starting with the symbol “>”. “>” will always be there at the beginning of each new command - don’t try to delete it! Moreover, you never need to type it.



We type the **commands** needed for a particular task after this prompt. The command is carried out by *R* after you hit the Return key.

Once you have started *R*, you can start typing commands into the RStudio console, and the results will be calculated immediately, for example:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIyKjMifQ==
```

Note that prior to the output of “6” it shows “[1]”.

Now subtraction:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIxMC0zIn0=
```

Again, prior to the output of “7” it shows “[1]”.

R can act like a basic calculator that you type commands in to. You can also use it like a more advanced scientific calculator and create **variables** that store information. All variables created by *R* are called **objects**. In *R*, we assign values to variables using an arrow-looking function <- the **assignment operator**. For example, we can **assign** the value 2*3 to the variable x using the command:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ4IDwtIDIqMyJ9
```

To view the contents of any *R* object, just type its name, press enter, and the contents of that *R* object will be displayed:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ4In0=
```

9.4 Variables in R

There are several different types of objects in *R* with fancy math names, including **scalars**, **vectors**, **matrices** (singular: **matrix**), arrays, dataframes, tables, **and** lists. The scalar** variable x above is one example of an *R* object. While a scalar variable such as x has just one element, a **vector** consists of several elements. The elements in a vector are all of the same **type** (e.g.. numbers or alphabetic characters), while **lists** may include elements such as characters as well as numeric quantities. Vectors and dataframes are the most common variables you’ll use. You’ll also encounter matrices often, and lists are ubiquitous in *R* but beginning users often don’t encounter them because they remain behind the scenes.

9.4.1 Vectors

To create a vector, we can use the `c()` (combine) function. For example, to create a vector called **myvector** that has elements with values 8, 6, 9, 10, and 5, we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcjA4LSBjKDgsIDYsIDksIDUwLCA1KSAjIG5vdGU6IGNvbW1hcy99
```

To see the contents of the variable **myvector**, we can just type its name and press enter:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcjE9
```

9.4.2 Vector indexing

The `[1]` is the **index** of the first **element** in the vector. We can **extract** any element of the vector by typing the vector name with the index of that element given in **square brackets** `[...]`.

For example, to get the value of the 4th element in the vector `myvector`, we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteXZlY3Rvcls0XSJ9
```

9.4.3 Character vectors

Vectors can contain letters, such as those designating nucleic acids

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteS5zZXEgPC0gYyhcIkFclxcIIReIixcIkNcIixcIkdcIikifQ==
```

They can also contain multi-letter **strings**:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteS5vbGlnb3MgPC0gYyhcIkFUQ0dDXClXCJUVFRDR0NcIixcIkNDQ0dDR0
```

9.4.4 Lists

NOTE: *below is a discussion of lists in R. This is excellent information, but not necessary if this is your very very first time using R.*

In contrast to a vector, a **list** can contain elements of different types, for example, both numbers and letters. A list can even include other variables such as a vector. The `list()` function is used to create a list. For example, we could create a list `mylist` by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3QgPC0gbGldChuYW1lPVwiQ2hhcmxlcYBEYXJ3aW5cIiwgXG4gI
```

We can then print out the contents of the list `mylist` by typing its name:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3QifQ==
```

The **elements** in a list are numbered, and can be referred to using **indices**. We can extract an element of a list by typing the list name with the index of the element given in double **square brackets** (in contrast to a vector, where we only use single square brackets).

We can extract the second element from `mylist` by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3RbWzJdXSAgIyBub3RlIHRoZSBkb3VibGUgc3F1YXJlIGJyYWNrZ
```

As a baby step towards our next task, we can wrap index values as in the `c()` command like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3RbW2MoMildXSAgIyBub3RlIHRoZSBkb3VibGUgc3F1YXJlIGJyY
```

The number 2 and `c(2)` mean the same thing.

Now, we can extract the second AND third elements from `mylist`. First, we put the indices 2 and 3 into a vector `c(2,3)`, then wrap that vector in double square brackets: `[c(2,3)]`. All together it looks like this.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3RbYygyLDMpXSAjIG5vdGUgdGhlIGRvdWJsZSBicmFja2V0cyJ9
```

Elements of lists may also be named, resulting in a **named lists**. The elements may then be referred to by giving the list name, followed by “\$”, followed by the element name. For example, `mylist$name` is the same as `mylist[[1]]` and `mylist$wife` is the same as `mylist[[2]]`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteWxpc3Qkd2lmZSJ9
```

We can find out the names of the named elements in a list by using the `attributes()` function, for example:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JteXZlY3Rvcl9kaHRyaWJ1dGVzKG15bGldCkifQ==
```

When you use the `attributes()` function to find the named elements of a list variable, the named elements are always listed under a heading “\$names”. Therefore, we see that the named elements of the list variable

`mylist` are called “name” and “wife”, and we can retrieve their values by typing `mylist$name` and `mylist$wife`, respectively.

9.4.5 Tables

Another type of object that you will encounter in R is a **table**. The `table()` function allows you to total up or tabulate the number of times a value occurs within a vector. Tables are typically used on vectors containing **character data**, such as letters, words, or names, but can work on numeric data data.

9.4.5.1 Tables - The basics

If we made a vector variable “nucleotides” containing the of a DNA molecule, we can use the `table()` function to produce a **table variable** that contains the number of bases with each possible nucleotides:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJiYXNlcyA8LSBjKFwiQVwiLCBcIlRcIiwgXCJBXCIsIFwiQVwiLCBcIlRcIiwg
```

Now make the table

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJ0YWJsZShiYXNlcykifQ==
```

We can store the table variable produced by the function `table()`, and call the stored table “bases.table”, by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJiYXNlcy50YWJsZSA8LSB0YWJsZShiYXNlcykifQ==
```

Tables also work on vectors containing numbers. First, a vector of numbers.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJudW1lcmljLnZlY3RlciA8LSBjKDEsMSwxLDEsMyw0LDQsNCw0KSJ9
```

Second, a table, showing how many times each number occurs.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJ0YWJsZShudW1lcmljLnZlY3RlcikifQ==
```

9.4.5.2 Tables - further details

To access elements in a table variable, you need to use double square brackets, just like accessing elements in a list. For example, to access the fourth element in the table `bases.table` (the number of Ts in the sequence), we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJiYXNlcy50YWJsZVtbNF1dICAjIGRvdWJsZSBicmFja2V0cyEifQ==
```

Alternatively, you can use the name of the fourth element in the table (“John”) to find the value of that table element:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJiYXNlcy50YWJsZVtbXCJUXCJdXSJ9
```

9.5 Arguments

Functions in R usually require **arguments**, which are input variables (i.e.. objects) that are **passed** to them, which they then carry out some operation on. For example, the `log10()` function is passed a number, and it then calculates the log to the base 10 of that number:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJsb2cxMCGxMDApIn0=
```

There’s a more generic function, `log()`, where we pass it not only a number to take the log of, but also the specific **base** of the logarithm. To take the log base 10 with the `log()` function we do this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJsb2coMTAwLCBiYXNlID0gMTApIn0=
```

We can also take logs with other bases, such as 2:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUOiJsb2coMTAwLCBiYXNlID0gMikifQ==
```

9.6 Help files with `help()` and ?

In *R*, you can get help about a particular function by using the `help()` function. For example, if you want help about the `log10()` function, you can type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJoZWxwKFwibG9nMTBcLikifQ==
```

When you use the `help()` function, a box or web pag will show up in one of the panes of RStudio with information about the function that you asked for help with. You can also use the `?` next to the function like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiI/bG9nMTAifQ==
```

Help files are a mixed bag in *R*, and it can take some getting used to them. An excellent overview of this is Kieran Healy’s “How to read an *R* help page.”

9.7 Searching for functions with `help.search()` and `RSiteSearch()`

If you are not sure of the name of a function, but think you know part of its name, you can search for the function name using the `help.search()` and `RSiteSearch()` functions. The `help.search()` function searches to see if you already have a function installed (from one of the *R* packages that you have installed) that may be related to some topic you’re interested in. `RSiteSearch()` searches *all* *R* functions (including those in packages that you haven’t yet installed) for functions related to the topic you are interested in.

For example, if you want to know if there is a function to calculate the standard deviation (SD) of a set of numbers, you can search for the names of all installed functions containing the word “deviation” in their description by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJoZWxwLnNlYXJjaChcImRldmlhdGlvblwiKSJ9
```

Among the functions that were found, is the function `sd()` in the `stats` package (an *R* package that comes with the base *R* installation), which is used for calculating the standard deviation.

Now, instead of searching just the packages we’ve have on our computer let’s search all *R* packages on CRAN. Let’s look for things related to DNA. Note that `RSiteSearch()` doesn’t provide output within RStudio, but rather opens up your web browser for you to display the results.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJSU2l0ZVNlYXJjaChcIkROQVwiKSJ9
```

The results of the `RSiteSearch()` function will be hits to descriptions of *R* functions, as well as to *R* mailing list discussions of those functions.

9.8 More on functions

We can perform computations with *R* using objects such as scalars and vectors. For example, to calculate the average of the values in the vector `myvector` (i.e.. the average of 8, 6, 9, 10 and 5), we can use the `mean()` function:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKG15dmVjdG9yKSAjIG5vdGU6IG5vIFwiIFwiIn0=
```

We have been using built-in *R* functions such as `mean()`, `length()`, `print()`, `plot()`, etc.

9.8.1 Writing your own functions

NOTE: *Writing your own functions is an advanced skills. New users can skip this section.

We can also create our own functions in *R* to do calculations that you want to carry out very often on different input data sets. For example, we can create a function to calculate the value of 20 plus square of some input number:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteWZ1bmN0aW9uIDwtIGZ1bmN0aW9uKHgpIHsgcmV0dXJuKDIwICsgKHg
```

This function will calculate the square of a number (x), and then add 20 to that value. The `return()` statement returns the calculated value. Once you have typed in this function, the function is then available for use. For example, we can use the function for different input numbers (e.g., 10, 25):

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteWZ1bmN0aW9uKDEwKSJ9
```

9.9 Quitting R

To quit R either close the program, or type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJxKCkifQ==
```

9.10 Links and Further Reading

Some links are included here for further reading.

For a more in-depth introduction to R, a good online tutorial is available on the “Kickstarting R” website, cran.r-project.org/doc/contrib/Lemon-kickstart.

There is another nice (slightly more in-depth) tutorial to R available on the “Introduction to R” website, cran.r-project.org/doc/manuals/R-intro.html.

Chapter 3 of Danielle Navarro’s book is an excellent intro to the basics of R.

Chapter 10

DNA descriptive statics - Part 1

By: Avril Coghlan

Adapted, edited and expanded: Nathan Brouwer (brouwern@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

10.1 Preface

This is a modification of “DNA Sequence Statistics (1)” from Avril Coghlan’s *A little book of R for bioinformatics..* The text and code were originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

10.2 Writing TODO:

- Add biology introduction
- Work on flow
- organize intial sections (intro, vocab, preliminaries)

10.3 Introduction

10.4 Vocabulary

- GC content
- DNA words
- scatterplots, histograms, piecharts, and boxplots

10.5 Functions

- `seqinr::GC()`
- `seqinr::count()`

10.6 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKVxubGlicmFyeShzZXFpbnIpIn0=

10.7 Converting DNA from FASTA format

In a previous exercise we downloaded and examined DNA sequence in the FASTA format. The sequence we worked with is also stored as a data file within the `compbio4a11` package and can be brought into memory using the `data()` command.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX2Zhc3RhXCIPIn0=
```

We can look at this data object with the `str()` command

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX2Zhc3RhKSJ9
```

This isn't in a format we can work with directly so we'll use the function `fasta_cleaner()` to set it up.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

Now check it out.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

What we have here is each base of the sequence in a separate slot of our vector.

The first four bases are "AGTT"

We can see the first one like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

The second one like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

The first and second like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

and all four like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

10.8 Length of a DNA sequence

Once you have retrieved a DNA sequence, we can obtain some simple statistics to describe that sequence, such as the sequence's total length in nucleotides. In the above example, we retrieved the DEN-1 Dengue virus genome sequence, and stored it in the vector variable `dengueseq_vector`. To obtain the length of the genome sequence, we would use the `length()` function, typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

The `length()` function will give you back the length of the sequence stored in variable `dengueseq_vector`, in nucleotides. The `length()` function actually gives the number of **elements** (slots) in the input vector that you passed to it, which in this case is the number of elements in the vector `dengueseq_vector`. Since each element of the vector `dengueseq_vector` contains one nucleotide of the DEN-1 Dengue virus sequence, the result for the DEN-1 Dengue virus genome tells us the length of its genome sequence (ie. 10735 nucleotides long).

10.8.1 Base composition of a DNA sequence

An obvious first analysis of any DNA sequence is to count the number of occurrences of the four different nucleotides ("A", "C", "G", and "T") in the sequence. This can be done using the `table()` function. For example, to find the number of As, Cs, Gs, and Ts in the DEN-1 Dengue virus sequence (which you have put into vector variable `dengueseq_vector`, using the commands above), you would type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JkYXRhKFwiZGVuZ3Vlc2VxX3ZlY3RvcikifQ==
```

This means that the DEN-1 Dengue virus genome sequence has 3426 As occurring throughout the genome, 2240 Cs, and so forth.

10.8.2 GC Content of DNA

One of the most fundamental properties of a genome sequence is its **GC content**, the fraction of the sequence that consists of Gs and Cs, ie. the $\%(G+C)$.

The GC content can be calculated as the percentage of the bases in the genome that are Gs or Cs. That is, GC content = (number of Gs + number of Cs)100/(genome length). For example, if the genome is 100 bp, and 20 bases are Gs and 21 bases are Cs, then the GC content is $(20 + 21)100/100 = 41\%$.

You can easily calculate the GC content based on the number of As, Gs, Cs, and Ts in the genome sequence. For example, for the DEN-1 Dengue virus genome sequence, we know from using the `table()` function above that the genome contains 3426 As, 2240 Cs, 2770 Gs and 2299 Ts. Therefore, we can calculate the GC content using the command:

eyJsYW5ndWFnZSI6InRlcjZlYW1wbGUiOiIoMjI0MCsyNzcxKSoxMDAvKDM0MjI0MCsyNzcxKzIyOTkpb30=

Alternatively, if you are feeling lazy, you can use the `GC()` function in the `SeqinR` package, which gives the fraction of bases in the sequence that are Gs or Cs.

eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzZXFpbmI6OkdDKGRlbmd1ZXNlcV92ZWNO0b3IpIn0=

The result above means that the fraction of bases in the DEN-1 Dengue virus genome that are Gs or Cs is 0.4666977. To convert the fraction to a percentage, we have to multiply by 100, so the GC content as a percentage is 46.66977%.

10.8.3 DNA words

As well as the frequency of each of the individual nucleotides (“A”, “G”, “T”, “C”) in a DNA sequence, it is also interesting to know the frequency of longer **DNA words**, also referred to as **genomic words**. The individual nucleotides are DNA words that are 1 nucleotide long, but we may also want to find out the frequency of DNA words that are 2 nucleotides long (ie. “AA”, “AG”, “AC”, “AT”, “CA”, “CG”, “CC”, “CT”, “GA”, “GG”, “GC”, “GT”, “TA”, “TG”, “TC”, and “TT”), 3 nucleotides long (eg. “AAA”, “AAT”, “ACG”, etc.), 4 nucleotides long, etc.

To find the number of occurrences of DNA words of a particular length, we can use the `count()` function from the R `SeqinR` package.

The count() function only works with lower-case letters, so first we have to use the tolower() function to convert our upper class genome to lower case

evJsYW5ndWFnZSI6InLiLCJzYW1wbGUiOiJkZW5ndWVzZXFFdmVjdG9yIDwtdG9sb3dlcihkZW5ndWVzZXFFdmVjdG9yYk

Now we can look for words. For example, to find the number of occurrences of DNA words that are 1 nucleotide long in the sequence `dengueseq_vector`, we type:

eyJsYW5ndWFnZSI6InRlcjZlYW1wbGUuOiJzZXZpbnI6OmNvdW50KGRlbmd1ZXNlcV92ZWNo3IsIDEpIn0=

As expected, this gives us the number of occurrences of the individual nucleotides. To find the number of occurrences of DNA words that are 2 nucleotides long, we type:

eyJsYW5ndWFnZSI6InRlcjZlYW1wbGUuOiJzZXNlbnI6OmNvdW50KGRlbmd1ZXNlcV92ZWNoY3IsIDIpIn0=

Note that by default the `count()` function includes all overlapping DNA words in a sequence. Therefore, for example, the sequence “ATG” is considered to contain two words that are two nucleotides long: “AT” and “TG”;

If you type `help('count')`, you will see that the result (output) of the function `count()` is a table object. This means that you can use double square brackets to extract the values of elements from the table. For example,

to extract the value of the third element (the number of Gs in the DEN-1 Dengue virus sequence), you can type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZW5ndWV0YWJsZV8yIDwtIHNlcWlucjo6Y291bnQoZGVuZ3Vlc2VxX3ZlY
```

The command above extracts the third element of the table produced by `count(dengueseq_vector,1)`, which we have stored in the table variable `denguetable`.

Alternatively, you can find the value of the element of the table in column “g” by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZW5ndWV0YWJsZV8yW1tcImFhXCJdXSJ9
```

Once you have table you can make a basic plot

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJiYXJwbG90KGRlbmd1ZXRhYmxlXzIpIn0=
```

We can sort by the number of words using the `sort()` command

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzb3J0KGRlbmd1ZXRhYmxlXzIpIn0=
```

Let’s save over the original object

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZW5ndWV0YWJsZV8yIDwtIHNvcnQoZGVuZ3VldGFibGVfMikifQ==
```

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJiYXJwbG90KGRlbmd1ZXRhYmxlXzIpIn0=
```

R will automatically try to optimize the appearance of the labels on the graph so you may not see all of them; no worries.

R can also make pie charts. Piecharts only really work when there are a few items being plots, like the four bases.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZW5ndWV0YWJsZV8xIDwtIHNlcWlucjo6Y291bnQoZGVuZ3Vlc2VxX3ZlY
```

Make a piechart with `pie()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwaWUoZGVuZ3VldGFibGVfMSkifQ==
```

10.8.4 Summary

In this practical, have learned to use the following R functions:

`length()` for finding the length of a vector or list `table()` for printing out a table of the number of occurrences of each type of item in a vector or list. These functions belong to the standard installation of R.

You have also learnt the following R functions that belong to the `SeqinR` package:

`GC()` for calculating the GC content for a DNA sequence `count()` for calculating the number of occurrences of DNA words of a particular length in a DNA sequence

10.9 Acknowledgements

This is a modification of “DNA Sequence Statistics (1)” from Avril Coghlan’s *A little book of R for bioinformatics..* Almost all of text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

In “A little book...” Coghlan noted: “Many of the ideas for the examples and exercises for this chapter were inspired by the Matlab case studies on *Haemophilus influenzae* (www.computational-genomics.net/case_studies/haemophilus_demo.html) and Bacteriophage lambda (http://www.computational-genomics.net/case_studies/lambdaphage_demo.html) from the website that accompanies the book *Introduction to Computational Genomics: a case studies approach* by Cristianini and Hahn (Cambridge University Press; www.computational-genomics.net/book/).”

10.9.1 License

The content in this book is licensed under a Creative Commons Attribution 3.0 License.

<https://creativecommons.org/licenses/by/3.0/us/>

10.9.2 Exercises

Answer the following questions, using the R package. For each question, please record your answer, and what you typed into R to get this answer.

Model answers to the exercises are given in Answers to the exercises on DNA Sequence Statistics (1).

1. What are the last twenty nucleotides of the Dengue virus genome sequence?
2. What is the length in nucleotides of the genome sequence for the bacterium *Mycobacterium leprae* strain TN (accession NC_002677)? Note: *Mycobacterium leprae* is a bacterium that is responsible for causing leprosy, which is classified by the WHO as a neglected tropical disease. As the genome sequence is a DNA sequence, if you are retrieving its sequence via the NCBI website, you will need to look for it in the NCBI Nucleotide database.
3. How many of each of the four nucleotides A, C, T and G, and any other symbols, are there in the *Mycobacterium leprae* TN genome sequence? Note: other symbols apart from the four nucleotides A/C/T/G may appear in a sequence. They correspond to positions in the sequence that are not clearly one base or another and they are due, for example, to sequencing uncertainties. For example, the symbol 'N' means 'any base', while 'R' means 'A or G' (purine). There is a table of symbols at www.bioinformatics.org/sms/iupac.html.
4. What is the GC content of the *Mycobacterium leprae* TN genome sequence, when (i) all non-A/C/T/G nucleotides are included, (ii) non-A/C/T/G nucleotides are discarded? Hint: look at the help page for the `GC()` function to find out how it deals with non-A/C/T/G nucleotides.
5. How many of each of the four nucleotides A, C, T and G are there in the complement of the *Mycobacterium leprae* TN genome sequence? *Hint*: you will first need to search for a function to calculate the complement of a sequence. Once you have found out what function to use, remember to use the `help()` function to find out what are the arguments (inputs) and results (outputs) of that function. How does the function deal with symbols other than the four nucleotides A, C, T and G? Are the numbers of As, Cs, Ts, and Gs in the complementary sequence what you would expect?
6. How many occurrences of the DNA words CC, CG and GC occur in the *Mycobacterium leprae* TN genome sequence?
7. How many occurrences of the DNA words CC, CG and GC occur in the (i) first 1000 and (ii) last 1000 nucleotides of the *Mycobacterium leprae* TN genome sequence? 1. How can you check that the subsequence that you have looked at is 1000 nucleotides long?

Chapter 11

Programming in R: for loops

By: Avril Coghlan

Adapted, edited and expanded: Nathan Brouwer (brouwer.n@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

11.1 Preface

This is a modification of “DNA Sequence Statistics (1)” from Avril Coghlan’s *A little book of R for bioinformatics*. The text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

11.2 Vocab

- for loop
- curly brackets

11.3 Functions

- for()
- print()

11.4 Basic for loops in R

In R, just as in programming languages such as **Python**, it is possible to write a **for loop** to carry out the same command several times. For example, say we have a pressing need to calculate the square the square of each number between 1 and 4. We could write for lines of code like this to do it:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIxXjJcbjJeMlxuM14yXG40XjIifQ==
```

If we know how to write a for loop, we could do the same think like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJmb3IgKGkgaW4gMT00KSB7IFxuICBwcmVudCAoaSppKSBcbiAgfSJ9
```

In the for loop above, the variable **i** is a counter or **index** for the number of cycles through the loop. In the first cycle through the loop, the value of **i** is 1, and so **i * i = 1** is printed out. In the second cycle through the loop, the value of **i** is 2, and so **i * i = 4** is printed out. In the third cycle through the loop, the value of **i** is 3, and so **i * i = 9** is printed out. The loop continues until the value of **i** is 4.

Note that the commands that are to be carried out at each cycle of the for loop must be enclosed within **curly brackets** (“{” and “}”).

You may be thinking “*ok, so it took four lines of code to do 1^2 through 4^2 each on their own, and three lines to do it with the loop; what’s the big deal?*”. What if you need to do 1 through 100 squared for some reason? Now the for loop is a lot less work.

You can also give a for loop a vector of numbers containing the values that you want the counter `i` to take in subsequent cycles. For example, you can make a vector containing the numbers 1, 2, 3, and 4, and write a for loop to print out the square of each number in vector `avector`:

```
## [1] 1
## [1] 4
## [1] 9
## [1] 16
```

The results should be the same as before.

11.5 Challenge: complicated vectors of values

Here’s a more complex example. If you don’t understand it don’t worry, it’s not something you’d probably do in practice.

Challenge: How can we use a for loop to print out the square of every second number between, say, 1 and 10? The answer is to use the `seq()` function with “`by = 2`” to tell the for loop to take every second number between 1 and 10:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJmb3IgKGkgaW4gc2VxKDEsIDFwLCBieSA9IDIpKSB7IFxuICBwcmVudCAo
```

In the first cycle of this loop, the value of `i` is 1, and so `i * i = 1` is printed out. In the second cycle through the loop, the value of `i` is 3, and so `i * i = 9` is printed out. The loop continues until the value of `i` is 9. In the fifth cycle through the loop, the value of `i` is 9, and so `i * i = 81` is printed out.

Chapter 12

Mini tutorial: Vectors in R

By: Avril Coghlan

Adapted, edited and expanded: Nathan Brouwer (brouwern@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

12.1 Preface

This is a modification of part of “DNA Sequence Statistics (2)” from Avril Coghlan’s *A little book of R for bioinformatics*.. Most of text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

12.2 Vocab

- base R
- scalar, vector, matrix
- regular expressions

Chapter 13

Functions

- `seq()`
- `is()`, `is.vector()`, `is.matrix()`
- `gsub()`

13.1 Vectors in R

Variables in R include **scalars**, **vectors**, and **lists**. **Functions** in R carry out operations on variables, for example, using the `log10()` function to calculate the log to the base 10 of a scalar variable `x`, or using the `mean()` function to calculate the average of the values in a vector variable `myvector`. For example, we can use `log10()` on a scalar object like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpIHNB3JlIHZhbHVlIG9iamVjdFxaeCA8LSAxMDBCblxuIyB0YWtlIGxv
```

Note that while mathematically `x` is a single number, or a scalar, R considers it to be a vector:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy52ZWNB3IoeCkifQ==
```

There are many “is” commands. What is returned when you run `is.matrix()` on a vector?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy5tYXRyaXgoeCkifQ==
```

Mathematically this is a bit odd, since often a vector is defined as a one-dimensional matrix, e.g., a single column or single row of a matrix. But in *R* land, a vector is a vector, and matrix is a matrix, and there are no explicit scalars.

13.2 Math on vectors

Vectors can serve as the input for mathematical operations. When this is done *R* does the mathematical operation separately on each element of the vector. This is a unique feature of *R* that can be hard to get used to even for people with previous programming experience.

Let’s make a vector of numbers:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvciA8LSBjKDMwLDE2LDMwMyw5OSwxMSwxMTEpIn0=
```

What happens when we multiply `myvector` by 10?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcioxMCJ9
```

R has taken each of the 6 values, 30 through 111, of `myvector` and multiplied each one by 10, giving us 6 results. That is, what R did was

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIzMCoxMCAgICAgIGZpcnN0IHZhbHVlIG9mIG15dmVjdG9yXG4xNioxMCAg
```

The normal order of operations rules apply to vectors as they do to operations we're more used to. So multiplying `myvector` by 10 is the same whether you put the 10 before or after `vector`. That is `myvector*10` is the same as `10*myvector`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcioxMFxuMTAqbXl2ZWNOB3IifQ==
```

What happens when you subtract 30 from `myvector`? Write the code below.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvciozMCI9
```

So, what R did was

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIzMCIzMCAgICAgIGZpcnN0IHZhbHVlIG9mIG15dmVjdG9yXG4xNi0zMCAg
```

You can also square a vector

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3Rvcj4yIn0=
```

Which is the same as

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIzMCIzMCAgICMgZmlyczQgdmFsdWUgb2YgbXl2ZWNOB3JcbjE2XjIgICAgIy
```

Also you can take the square root of a vector...

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzcXJ0KG15dmVjdG9yKSJ9
```

...and take the log of a vector...

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsb2cobXl2ZWNOB3IpIn0=
```

...and just about any other mathematical operation. Here we are working on a separate vector object; all of these rules apply to a column in a matrix or a dataframe. This attribution of R is called **vectorization**.

13.3 Functions on vectors

We can use functions on vectors. Typically these use the vectors as an input and all the numbers are processed into an output. Call the `mean()` function on the vector we made called `myvector`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKG15dmVjdG9yKSJ9
```

Note how we get a single value back - the mean of all the values in the vector. R saw that we had a vector of multiple and knew that the mean is a function that doesn't get applied to single number, but sets of numbers.

The function `sd()` calculates the standard deviation. Apply the `sd()` to `myvector`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzcHteXZlY3RvcikifQ==
```

13.4 Operations with two vectors

You can also subtract one vector from another vector. This can be a little weird when you first see it. Make another vector with the numbers 5, 10, 15, 20, 25, 30. Call this `myvector2`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcjIgPC0gYyglLCxMCwgMTUsIDIwLCAYNSwgMzApIn0=
```

Now subtract `myvector2` from `myvector`. What happens?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3Rvcj1teXZlY3RvcjIifQ==
```

13.5 Subsetting vectors

You can extract an **element** of a vector by typing the vector name with the index of that element given in **square brackets**. For example, to get the value of the 3rd element in the vector `myvector`, we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvclszXSJ9
```

Extract the 4th element of the vector:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3Rvcls0XSJ9
```

You can extract more than one element by using a vector in the brackets:

First, say I want to extract the 3rd and the 4th element. I can make a vector with 3 and 4 in it:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJudW1zIDwtIGMoMyw0KSJ9
```

Then put that vector in the brackets:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcltudW1zXSJ9
```

We can also do it directly like this, skipping the vector-creation step:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcltjKDMsNCldIn0=
```

In the chunk below extract the 1st and 2nd elements:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteXZlY3RvcltjKDEsMildIn0=
```

13.6 Sequences of numbers

Often we want a vector of numbers in **sequential order**. That is, a vector with the numbers 1, 2, 3, 4, ... or 5, 10, 15, 20, ... The easiest way to do this is using a colon

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIxOjEwIn0=
```

Note that in R `1:10` is equivalent to `c(1:10)`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJjKDE6MTApIn0=
```

Usually to emphasize that a vector is being created I will use `c(1:10)`

We can do any number to any numbers

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJjKDIwOjMwKSJ9
```

We can also do it in *reverse*. In the code below put 30 before 20:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJjKDMwOjIwKSJ9
```

A useful function in *R* is the `seq()` function, which is an explicit function that can be used to create a vector containing a sequence of numbers that run from a particular number to another particular number.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZXEoMSwgMTApIn0=
```

Using `seq()` instead of a `:` can be useful for readability to make it explicit what is going on. More importantly, `seq` has an argument `by = ...` so you can make a sequence of number with any interval between. For example, if we want to create the sequence of numbers from 1 to 10 in steps of 1 (i.e., 1, 2, 3, 4, ... 10), we can type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZXEoMSwgMTAsXG4gICAgYnkgPSAxKSJ9
```

We can change the **step size** by altering the value of the `by` argument given to the function `seq()`. For example, if we want to create a sequence of numbers from 1-100 in steps of 20 (i.e., 1, 21, 41, ... 101), we can type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZXEoMSwgMTAxLFxuICAgIGJ5ID0gMjApIn0=
```

13.7 Vectors can hold numeric or character data

The vector we created above holds numeric data, as indicated by `class()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JiJjbGFzcyhteXZlY3RvcikifQ==
```

Vectors can also hold character data, like the genetic code:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JiJjbGFzcyhteXZlY3RvcikifQ==
```

13.8 Regular expressions can modify character data

We can use **regular expressions** to modify character data. For example, change the Ts to Us

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JiJjbGFzcyhteXZlY3RvcikifQ==
```

Now check it out

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JiJjbGFzcyhteXZlY3RvcikifQ==
```

Regular expressions are a deep subject in computing. You can find some more information about them [here](#).

Chapter 14

Plotting vectors in base R

By: Avril Coghlan

Adapted, edited and expanded: Nathan Brouwer (brouwer@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

14.1 Preface

This is a modification of part of “DNA Sequence Statistics (2)” from Avril Coghlan’s *A little book of R for bioinformatics*. Most of text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

14.2 Plotting numeric data

R allows the production of a variety of plots, including **scatterplots**, **histograms**, **piecharts**, and **boxplots**. Usually we make plots from dataframes with 2 or more columns, but we can also make them from two separate vectors. This flexibility is useful, but also can cause some confusion.

For example, if you have two equal-length vectors of numbers `numeric.vect1` and `numeric.vect2`, you can plot a scatterplot of the values in `myvector1` against the values in `myvector2` using the **base R** `plot()` function.

First, let’s make up some data in put it in vectors:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJudW1lcmljLnZlY3QxIDwtIGMoMTAsIDE1LCAyMiwgMzUsIDQzKVxubnVt.
```

Not plot with the base R `plot()` function:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwbG90KG51bWVyaWMudmVjdDEsIG51bWVyaWMudmVjdDIpIn0=
```

Note that there is a comma between the two vector names. When building plots from dataframes you usually see a tilde (~), but when you have two vectors you can use just a comma.

Also note the order of the vectors within the `plot()` command and which axes they appear on. The first vector is `numeric.vect1` and it appears on the horizontal x-axis.

If you want to label the axes on the plot, you can do this by giving the `plot()` function values for its optional arguments `xlab =` and `ylab =`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwbG90KG51bWVyaWMudmVjdDEsICAgIyBub3RIIGFnYWluIHRoZSBjb21t
```

We can store character data in vectors so if we want we could do this to set up our labels:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteWxhYmVscyA8LSAgYyhcIm51bWVyaWMudmVjdDFcIixcIm51bWVyaWM
```

Then use bracket notation to call the labels from the vector

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwbG90KG51bWVyaWMudmVjdDEsIFxuICAgICBudW1lcmljLnZlY3QyLCB
```

If we want we can use a tilde to make our plot like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwbG90KG51bWVyaWMudmVjdDIgfiBudW1lcmljLnZlY3QxKSJ9
```

Note that now, `numeric.vect2` is on the left and `numeric.vect1` is on the right. This flexibility can be tricky to keep track of.

We can also combine these vectors into a dataframe and plot the data by referencing the data frame. First, we combine the two separate vectors into a dataframe using the `cbind()` command.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZiA8LSBjYmluZChudW1lcmljLnZlY3QxLCBudW1lcmljLnZlY3QyKSJ9
```

Then we plot it like this, referencing the dataframe `df` via the `data = ...` argument.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwbG90KG51bWVyaWMudmVjdDIgfiBudW1lcmljLnZlY3QxLCBkYXRhID0g
```

14.3 Other plotting packages

Base R has lots of plotting functions; additionally, people have written packages to implement new plotting capabilities. The package `ggplot2` is currently the most popular plotting package, and `ggpubr` is a package which makes `ggplot2` easier to use. For quick plots we'll use base R functions, and when we get to more important things we'll use `ggplot2` and `ggpubr`.

Programming in R: functions

Adapted, edited and expanded: Nathan Brouwer (brouwern@gmail.com) under the Creative Commons 3.0 Attribution License (CC BY 3.0).

This is a modification of “DNA Sequence Statistics (1)” from Avril Coghlan’s *A little book of R for bioinformatics*. The text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

- function
- curly brackets

- `function()`

We have been using **built-in functions** such as `mean()`, `length()`, `print()`, `plot()`, etc. We can also create our own functions in R to do calculations that you want to carry out very often on different input data sets. For example, we can create a function to calculate the value of 20 plus the square of some input number:

This function will calculate the square of a number (x), and then add 20 to that value. It stores this in a temporary object called output. The return() statement returns the calculated value. Once you have typed in this function, the function is then available for use. For example, we can use the function for different input numbers (e.g., 10, 25):

You can view the code that makes up a function by typing its name (without any parentheses). For example, we can try this by typing “myfunction”:

eyJ5YW5ndWFnZSI6InIiLCJzYW1wbGUiOiJteWZ1bmN0aW9uIn0=

15.5 Comments in R

When you are typing R, if you want to, you can write comments by writing the comment text after the “#” sign. This can be useful if you want to write some R commands that other people need to read and understand. R will ignore the comments when it is executing the commands. For example, you may want to write a comment to explain what the function `log10()` does:

eyJ5YW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ4IDwtIDFwMFxubG9nMTAoeCkgIyBGaW5kcyB0aGUgbG9nIHRvIHRoZSB

Chapter 16

FASTA Files

Adapted from Wikipedia: https://en.wikipedia.org/wiki/FASTA_format

"In bioinformatics, the FASTA format is a text-based format for representing either nucleotide sequences or amino acid (protein) sequences, in which nucleotides or amino acids are represented using single-letter codes. The format allows for sequence names and comments to precede the sequences. The format originates from the FASTA alignment software, but has now become a near universal standard in the field of bioinformatics.

"The simplicity of FASTA format makes it easy to manipulate and parse sequences using text-processing tools and scripting languages like the R programming language and Python.

"The first line in a FASTA file starts with a ">" (greater-than) symbol and holds summary information about the sequence, often starting with a unique accession number and followed by information like the name of the gene, the type of sequence, and the organism it is from.

"On the next is the sequence itself in a standard one-letter character string. Anything other than a valid character is be ignored (including spaces, tabs, asterisks, etc...).

"A multiple sequence FASTA format can be obtained by concatenating several single sequence FASTA files in a common file (also known as multi-FASTA format).

"Following the header line, the actual sequence is represented. Sequences may be protein sequences or nucleic acid sequences, and they can contain gaps or alignment characters. Sequences are expected to be represented in the standard amino acid and nucleic acid codes. Lower-case letters are accepted and are mapped into upper-case; a single hyphen or dash can be used to represent a gap character; and in amino acid sequences, U and * are acceptable letters.

"FASTQ format is a form of FASTA format extended to indicate information related to sequencing. It is created by the Sanger Centre in Cambridge.

"Bioconductor.org's Biostrings package can be used to read and manipulate FASTA files in R

from <https://zhanglab.dcmf.med.umich.edu/FASTA/>

"FASTA format is a text-based format for representing either nucleotide sequences or peptide sequences, in which base pairs or amino acids are represented using single-letter codes. A sequence in FASTA format begins with a single-line description, followed by lines of sequence data. The description line is distinguished from the sequence data by a greater-than (">") symbol in the first column. It is recommended that all lines of text be shorter than 80 characters in length."

16.1 Example FASTA file

Here is an example of the contents of a FASTA file. (If you are viewing this chapter in the form of the source .Rmd file, the `cat()` function is included just to print out the content properly and is not part of the FASTA format).

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJjYXQoXCI+Z2l8MTg2NjgxMjI4fHJlZnxZUF8wMDE4NjQ0MjQuMXwgcGh5
```

16.2 Multiple sequences in a single FASTA file

Multiple sequences can be stored in a single FASTA file, each on separated by a line and have its own headline.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJjYXQoXCI+TENCTyAtIFByb2xhY3RpbjBwcmVjdXJzb3IgL3B3b3ZpbnVcl
```

16.3 Multiple sequence alignments can be stored in FASTA format

Aligned FASTA format can be used to store the output of **Multiple Sequence Alignment (MSA)**. This format contains

1. Multiple entries, each with their own header line
2. **Gaps** inserted to align sequences are indicated by `.`
3. Each spaces added to the beginning and end of sequences that vary in length are indicated by `~`

In the sample FASTA file below, the **example1** sequence has a gap of 8 near its beginning. The **example2** sequence has numerous `~` indicating that this sequence is missing data from its beginning that are present in the other sequences. The **example3** sequence has numerous `~` at its end, indicating that this sequence is shorter than the others.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJjYXQoXCI+ZXhhbXBsZTEgXG5NS0FMV0FMTEwUEExMVEEdDTEEuLi4
```

16.4 FASTQ Format

Adapted from Wikipedia: https://en.wikipedia.org/wiki/FASTQ_format

"FASTQ format is a text-based format for storing both a biological sequence (usually nucleotide sequence) and its corresponding quality scores. Both the sequence letter and quality score are each encoded with a single ASCII character for brevity.

"It was originally developed at the Wellcome Trust Sanger Institute to bundle a FASTA formatted sequence and its quality data, but has recently become the de facto standard for storing the output of high-throughput sequencing instruments such as the Illumina Genome Analyzer.

"A FASTQ file normally uses four lines per sequence.

- Line 1 begins with a `@` character and is followed by a sequence identifier and an optional description (like a FASTA title line).
- Line 2 is the raw sequence letters.
- Line 3 begins with a `+` character and is optionally followed by the same sequence identifier (and any description) again.
- Line 4 encodes the **quality values** for the sequence in Line 2 of the file, and must contain the same number of symbols as letters in the sequence.

"A FASTQ file containing a single sequence might look like this:"

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJjYXQoXCI+XCJAU0VRX0IEXG5HQVRUVEdHR0dUVENBQUFHQ0FHVEFUQ
```

"Here are the quality value characters in left-to-right increasing order of quality (ASCII):"

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIgIVwiIyQlJicoKSorLC0uLzAxMjM0NTY3ODk6Ozw9Pj9AQUJDREVGR0hJS

FASTQ files typically do not include line breaks and do not wrap around when they reach the width of a normal page or file.

Chapter 17

Downloading DNA sequences as FASTA files in R

This is a modification of “DNA Sequence Statistics” from Avril Coghlan’s *A little book of R for bioinformatics*.. Most of the text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

NOTE: There is some redundancy in this current draft that needs to be eliminated.

17.0.1 Functions

- `library()`
- `help()`
- `length`
- `table`
- `seqinr::GC()`
- `seqinr::count()`
- `seqinr::write.fasta()`

17.0.2 Software/websites

- www.ncbi.nlm.nih.gov
- Text editors (e.g. Notepad++, TextWrangler)

17.0.3 R vocabulary

- `list`
- `library`
- `package`
- `CRAN`
- `wrapper`

17.0.4 File types

- FASTA

17.0.5 Bioinformatics vocabulary

- accession, accession number

- NCBI
- NCBI Sequence Database
- EMBL Sequence Database
- FASTA file

17.0.6 Organisms and Sequence accessions

- Dengue virus: DEN-1, DEN-2, DEN-3, and DEN-4.

The NCBI accessions for the DNA sequences of the DEN-1, DEN-2, DEN-3, and DEN-4 Dengue viruses are NC_001477, NC_001474, NC_001475 and NC_002640, respectively.

According to Wikipedia

“Dengue virus (DENV) is the cause of dengue fever. It is a mosquito-borne, single positive-stranded RNA virus ... Five serotypes of the virus have been found, all of which can cause the full spectrum of disease. Nevertheless, scientists’ understanding of dengue virus may be simplistic, as rather than distinct ... groups, a continuum appears to exist.” https://en.wikipedia.org/wiki/Dengue_virus

17.0.7 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KHJlbnRyZXopXG5saWJyYXJ5KGNvbXBiaW80YWxsKSAifQ

17.1 DNA Sequence Statistics: Part 1

17.1.1 Using R for Bioinformatics

The chapter will guide you through the process of using R to carry out simple analyses that are common in bioinformatics and computational biology. In particular, the focus is on computational analysis of biological sequence data such as genome sequences and protein sequences. The programming approaches, however, are broadly generalizable to statistics and data science.

The tutorials assume that the reader has some basic knowledge of biology, but not necessarily of bioinformatics. The focus is to explain simple bioinformatics analysis, and to explain how to carry out these analyses using *R*.

17.1.2 R packages for bioinformatics: Bioconductor and SeqinR

Many authors have written *R* packages for performing a wide variety of analyses. These do not come with the standard *R* installation, but must be installed and loaded as “add-ons”.

Bioinformaticians have written numerous specialized packages for *R*. In this tutorial, you will learn to use some of the function in the **SeqinR** package to carry out simple analyses of DNA sequences. (**SeqinR** can retrieve sequences from a DNA sequence database, but this has largely been replaced by the functions in the package **rentrez**)

Many well-known bioinformatics packages for *R* are in the Bioconductor set of *R* packages (www.bioconductor.org), which contains packages with many *R* functions for analyzing biological data sets such as microarray data. The **SeqinR** package is from CRAN, which contains R functions for obtaining sequences from DNA and protein sequence databases, and for analyzing DNA and protein sequences.

We will also use functions from the **rentrez** and **ape** packages.

Remember that you can ask for more information about a particular *R* command by using the **help()** function. For example, to ask for more information about the **library()**, you can type:

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KHJlbnRyZXopXG5saWJyYXJ5KGNvbXBiaW80YWxsKSAifQ

You can also do this

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiI/bGlicmFyeSJ9
```

17.1.3 FASTA file format

The FASTA format is a simple and widely used format for storing biological (e.g. DNA or protein) sequences. It was first used by the FASTA program for sequence alignment in the 1980s and has been adopted as standard by many other programs.

FASTA files begin with a single-line description starting with a greater-than sign > character, followed on the next line by the sequences. Here is an example of a FASTA file. (If you're looking at the source script for this lesson you'll see the `cat()` command, which is just a text display function used format the text when you run the code).

```
## >A06852 183 residues MPRLFSYLLGVWLLSQLPREIPGQSTNDFIKACGRELVRWVEICGSVSWGRTALSLEEPQLETGPPAETMPSSITKD
```

17.1.4 The NCBI sequence database

The US National Centre for Biotechnology Information (NCBI) maintains the **NCBI Sequence Database**, a huge database of all the DNA and protein sequence data that has been collected. There are also similar databases in Europe, the European Molecular Biology Laboratory (EMBL) Sequence Database, and Japan, the DNA Data Bank of Japan (DDBJ). These three databases exchange data every night, so at any one point in time, they contain almost identical data.

Each sequence in the NCBI Sequence Database is stored in a separate **record**, and is assigned a unique identifier that can be used to refer to that record. The identifier is known as an **accession**, and consists of a mixture of numbers and letters.

For example, Dengue virus causes Dengue fever, which is classified as a **neglected tropical disease** by the World Health Organization (WHO), is classified by any one of four types of Dengue virus: DEN-1, DEN-2, DEN-3, and DEN-4. The NCBI accessions for the DNA sequences of the DEN-1, DEN-2, DEN-3, and DEN-4 Dengue viruses are

- NC_001477
- NC_001474
- NC_001475

- NC_002640

Note that because the NCBI Sequence Database, the EMBL Sequence Database, and DDBJ exchange data every night, the DEN-1 (and DEN-2, DEN-3, DEN-4) Dengue virus sequence are present in all three databases, but they have different accessions in each database, as they each use their own numbering systems for referring to their own sequence records.

17.1.5 Retrieving genome sequence data using rentrez

You can retrieve sequence data from NCBI directly from *R* using the **rentrez** package. The DEN-1 Dengue virus genome sequence has NCBI accession NC_001477. To retrieve a sequence with a particular NCBI accession, you can use the function `entrez_fetch()` from the **rentrez** package. Note that to be specific where the function comes from I write it as `package::function()`.

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJkZW5ndWVzZXFfZmFzdGEgPC0gcmlVudHJlejo6ZW50cmV6X2ZldGN0KGR
```

Note that the “ ” in the name is just an arbitrary way to separate two words. Another common format would be *dengueseq.fasta*. Some people like *dengueseqFasta*, called **camel case** because the capital letter makes a hump in the middle of the word. Underscores are becoming most common and are favored by developers associated with *RStudio* and the **tidyverse** of packages that many data scientists use. I switch between ”

and "" as separators, usually favoring "_" for function names and "." for objects; I personally find camel case harder to read and to type.

Ok, so what exactly have we done when we made `dengueseq_fasta`? We have an R object `dengueseq_fasta` which has the sequence linked to the accession number "NC_001477." So where is the sequence, and what is it?

First, what is it?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcyhhkZW5ndWVzZXFFZmFzdGEpXG5jbGFzcyhhkZW5ndWVzZXFFZmFzdGEp
```

How big is it? Try the `dim()` and `length()` commands and see which one works. Do you know why one works and the other doesn't?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkaW0oZGVuZ3Vlc2VxX2Zhc3RhKVxubGVuZ3RoKGRlbmd1ZXNlcV9mYXN0YSwgMTAwKSJ9
```

The size of the object is 1. Why is this? This is the genomics sequence of a virus, so you'd expect it to be fairly large. We'll use another function below to explore that issue. Think about this first: how many pieces of unique information are in the `dengueseq` object? In what sense is there only *one* piece of information?

If we want to actually see the sequence we can type just type `dengueseq_fasta` and press enter. This will print the WHOLE genomic sequence out but it will probably run off your screen.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkZW5ndWVzZXFFZmFzdGEifQ==
```

This is a whole genome sequence, but its stored as single entry in a vector, so the `length()` command just tells us how many entries there are in the vector, which is just one! What this means is that the entire genomic sequence is stored in a single entry of the vector `dengueseq_fasta`. (If you're not following along with this, no worries - its not essential to actually working with the data)

If we want to actually know how long the sequence is, we need to use the function `nchar()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJuY2hhcihhkZW5ndWVzZXFFZmFzdGEpIn0=
```

The sequence is 10935 bases long. All of these bases are stored as a single **character string** with no spaces in a single entry of our `dengueseq_fasta` vector. This isn't actually a useful format for us, so below were going to convert it to something more useful.

If we want to see just part of the sequence we can use the `strtrim()` function. Before you run the code below, predict what the 100 means.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzdHJ0cm9udKGRlbmd1ZXNlcV9mYXN0YSwgMTAwKSJ9
```

Note that at the end of the name is a slash followed by an n, which indicates to the computer that this is a **newline**; this is read by text editor, but is ignored by R in this context.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzdHJ0cm9udKGRlbmd1ZXNlcV9mYXN0YSwgNDUpIn0=
```

After the `\n` begins the sequence, which will continue on for a LOOOOOONG way. Let's just print a little bit.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzdHJ0cm9udKGRlbmd1ZXNlcV9mYXN0YSwgNTIpIn0=
```

Let's print some more. Do you notice anything beside A, T, C and G in the sequence?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzdHJ0cm9udKGRlbmd1ZXNlcV9mYXN0YSwgMjAwKSJ9
```

Again, there are the `\n` newline characters, which tell text editors and wordprocessors how to display the file.

Now that we a sense of what we're looking at let's explore the `dengueseq_fasta` a bit more.

We can find out more information about what it is using the `class()` command.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJjbGFzcyhkZW5ndWVzZXFfZmFzdGEpIn0=
```

As noted before, this is character data.

Many things in R are vectors so we can ask `R is.vector()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy52ZWNoY3I0ZGVuZ3Vlc2VxX2Zhc3RhKSJ9
```

Yup, that's true.

Ok, let's see what else. A handy though often verbose command is `is()`, which tells us what an object, well, what it is:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy52ZWNoY3I0ZGVuZ3Vlc2VxX2Zhc3RhKSJ9
```

There is a lot here but if you scan for some key words you will see “character” and “vector” at the top. The other stuff you can ignore. The first two things, though, tell us the `dengueseq_fasta` is a **vector** of the class **character**: that is, a **character vector**.

Another handy function is `str()`, which gives us a peak at the context and structure of an *R* object. This is most useful when you are working in the R console or with dataframes, but is a useful function to run on all *R* objects. How does this output differ from other ways we've displayed `dengueseq_fasta`?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzdHIoZGVuZ3Vlc2VxX2Zhc3RhKSJ9
```

We know it contains character data - how many characters? `nchar()` for “number of characters” answers that:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJuY2hhcihkZW5ndWVzZXFfZmFzdGEpIn0=
```

17.2 OPTIONAL: Saving FASTA files

We can save our data as .fasta file for safe keeping. The `write()` function will save the data we downloaded as a plain text file.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ3cmI0ZShkZW5ndWVzZXFfZmFzdGEsIFxuICAgICAgZmlsZT1cImRlbmd1Z
```

If you do this, you'll need to figure out where *R* is saving things, which requires an understanding of *R*'s **file system**, which can take some getting used to, especially if you're new to programming. As a start, you can see where *R* saves things by using the `getwd()` command, which tells you where on your harddrive *R* currently is using as its home base for files.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnZXR3ZCgpIn0=
```

17.3 Next steps

On their own, FASTA files in R are not directly useful. In the next lesson we'll process our `dengueseq_fasta` file so that we can use it in analyses.

Chapter 18

Downloading DNA sequences as FASTA files in R

This is a modification of “DNA Sequence Statistics” from Avril Coghlan’s *A little book of R for bioinformatics*. Most of the text and code was originally written by Dr. Coghlan and distributed under the Creative Commons 3.0 license.

18.1 Preliminaries

We’ll need the `dengueseq_fasta` FASTA data object, which is in the `compbio4all` package. We’ll also use the `stringr` package for cleaning up the FASTA data, which can be downloaded with `install.packages("stringr")`

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjIGNvbXBiaW80YWxsLCB3aGljaCB0YXMGZGVuZ3Vlc2VxX2Zhc3RhXG5s

18.2 Convert FASTA sequence to an R variable

We can’t actually do much with the contents of the `dengueseq_fasta` we downloaded with the `rentrez` package except read them. If we want to address some biological questions with the data we need is to convert it into a data structure *R* can work with.

There are several things we need to remove:

1. The **meta data** line `>NC_001477.1 Dengue virus 1, complete genome` (metadata is “data” about data, such as where it came from, what it is, who made it, etc.).
2. All the `\n` that show up in the file (these are the **line breaks**).
3. Put each nucleotide of the sequence into its own spot in a vector.

There are functions that can do this automatically, but

1. I haven’t found one I like, and
2. walking through this will help you understand the types of operations you can do on text data.

The first two steps involve removing things from the existing **character string** that contains the sequence. The third step will split the single continuous character string like “AGTTGTTAGTCTACGT...” into a **character vector** like `c("A", "G", "T", "T", "G", "T", "T", "A", "G", "T", "C", "T", "A", "C", "G", "T" ...)`, where each element of the vector is a single character stored in a separate slot in the vector.

18.2.1 Removing unwanted characters

The second item is the easiest to take care of. *R* and many programming languages have tools called **regular expressions** that allow you to manipulate text. *R* has a function called `gsub()` which allows you to substitute or delete character data from a string. First I'll remove all those `\n` values.

The regular expression function `gsub()` takes three arguments: 1. `pattern =` This is what we need it to find so we can replace it. 1. `replacement =` The replacement. 1. `x =` A character string or vector where `gsub()` will do its work.

We need to get rid of the `\n` so that we are left with only A, T and G, which are the actual information of the sequence. We want `\n` completely removed, so the replacement will be "", which is a set of quotation marks with nothing in the middle, which means "delete the target pattern and put nothing in its place."

One thing that is tricky about regular expressions is that many characters have special meaning to the functions, such as slashes, dollar signs, and brackets. So, if you want to find and replace one of these specially designated characters you need to put a slash in front of them. So when we set the pattern, instead of setting the pattern to a slash before an `n` `\n`, we have to give it two slashes `\\n`.

Here is the regular expression to delete the newline character `\n`.

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiIgIyBub3RlOiB3ZSB3YW50IHRvIGZpbmQgYWxsIHRoZSBcXG4sIGJ1dCBu
```

We can use `strtrim()` to see if it worked

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzdHJ0cmVtKGRlbmd1ZXNlcV92ZWN0b3IsIDgwKSJ9
```

Now for the metadata header. This is a bit complex, but the following code is going to take all the that occurs before the beginning of the sequence ("AGTTGTTAGTC") and delete it.

First, I'll define what I want to get rid of in an *R* object. This will make the **call** to `gsub()` a little cleaner to read

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzZXEuAGVhZGVyIDwtIFwiPk5DXzAwMTQ3Ny4xIERlbmd1ZSB2aXJ1cyAx
```

Now I'll get rid of the header with `gsub()`.

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJkZW5ndWVzZXFFdmVjdG9yIDwtIGdzdWIocGF0dGVybiA9IHNlcS5oZWFK
```

See if it worked:

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzdHJ0cmVtKGRlbmd1ZXNlcV92ZWN0b3IsIDgwKSJ9
```

18.2.2 Splitting unbroken strings in character vectors

Now the more complex part. We need to split up a continuous, unbroken string of letters into a vector where each letter is on its own. This can be done with the `str_split()` function ("string split") from the **stringr** package. The notation `stringr::str_split()` means "use the `str_split` function from the **stringr** package." More specifically, it temporarily loads the **stringr** package and gives *R* access to just the `str_split` function. This allows you to call a single function without loading the whole library.

There are several arguments to `str_split`, and I've tacked a `[[1]]` on to the end.

First, run the command

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJkZW5ndWVzZXFFdmVjdG9yX3NwbGl0IDwtIHN0cmVtZ3I6OnN0cl9zcGxpdmVz
```

Look at the output with `str()`

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzdHJ0ZGVuZ3Vlc2VxX3ZlY3Rvcl9zcGxpdmVzCkifQ==
```

We can explore what the different arguments do by modifying them. Change `pattern = ""` to `pattern = "A"`. Can you figure out what happened?

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJlIHJlXj1biB0aGUgY29tbWFuZCB3aXR0IFwicGF0dGVybiAgPSBcIkFcllXu
```

And try it with `pattern = ""` to `pattern = "G"`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIHJlXJ1biB0aGUgY29tbWFnZCB3aXRoIFwicGF0dGVybiAgPSBclkdcllxu
```

Run this code to compare the two ways we just used `str_split` (don't worry what it does). Does this help you see what's up?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJvcHRpb25zKHN0ciA9IHN0ck9wdGlvbnModmVjLmxmlbiA9IDEwKSicbnN0cih
```

So, what does the `pattern = ...` argument do? For more info open up the help file for `str_split` by calling `?str_split`.

Something cool which we will explore in the next exercise is that we can do summaries on vectors of nucleotides, like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ0YWJsZShkZW5ndWVzZXFfdmVjdG9yX3NwbGl0KSJ9
```


Downloading protein sequences in R

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKSJ9

We can use `entrez_fetch()` to download protein sequences.

eyJ5YW5ndWFnZSI6InliLCJzYW1wbGUiOiJpIHNlcXVlbmNIIDE6IFE5Q0Q4M1xubGVwcmFlX2Zhc3RhIDwtIHJlbnRyZX

eyJzYW5ndWFnZSI6InIiLCJzYW1wbGUOiJsZXByYWVfZmFzdGEifQ==

eyJzYW5ndWFnZSI6InRlcjZlYW1wbGUuOiJsZXByYWVfdmVjdG9yIENAgPC0gZmFzdGFfY2xlYW5lc2h5ZXByYWVfZmFz

eyJzYW5ndWFnZSI6InliLCJzYW1wbGUOiJsZW5ndGgobGVwcmFlX3ZlY3RvcilcbmNsYXNzKGxlcHJhZV92ZWNO0b3lpX

Sequence dotplots in R

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

20.1 Preliminaries

20.1.1 Download sequences

eyJzYW5ndWFnZSI6InRlcjZlYXVwGUiOiIjIHBlbmNlIDE6IFE5Q0Q4M1xubGVwcmFlX2ZhczRhIDwtIHJlbmRyZX

To help build our intuition about dotplots we'll first look at some artificial examples. First, we'll see what happens when we make a dotplot comparing the alphabet versus itself. The build-in `LETTERS` object in R contains the alphabet from A to Z. This is a sequence with no repeats.

What we get is a perfect diagonal line.

Now lets' make a sequence where the alphabet gets repeats twice

Note the diagonal lines.

eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJMVRURVJTLjMudGltZXMgPC0gYyhMRVRURVJTLExFVFRFUIMsTEV

Here's another example of repeats.

67

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXEuZmVwZWZ0IDwtIGMoXCJBXCIsXCJDXCIsXCJEXCIsXCJFXCIsXC
```

Make the dotplot:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

20.4 Inversions

See if you can figure out what's going on here.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

20.5 Translocations

See if you can figure out what's going on here.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

20.6 Random sequence

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

20.7 Comparing two real sequences using a dotplot

As a first step in comparing two protein, RNA or DNA sequences, it is a good idea to make a **dotplot**. A **dotplot** is a graphical method that allows the comparison of two protein or DNA sequences and identify regions of close similarity between them. A dotplot is essentially a two-dimensional matrix (like a grid), which has the sequences of the proteins being compared along the vertical and horizontal axes.

In order to make a simple dotplot to represent of the similarity between two sequences, individual cells in the matrix can be shaded black if residues are identical, so that matching sequence segments appear as runs of diagonal lines across the matrix. Identical proteins will have a line exactly on the main diagonal of the dotplot, that spans across the whole matrix.

For proteins that are not identical, but share regions of similarity, the dotplot will have shorter lines that may be on the **main diagonal**, or off the main diagonal of the matrix. In essence, a dotplot will reveal if there are any regions that are clearly very similar in two protein (or DNA) sequences.

We can create a dotplot for two sequences using the `dotPlot()` function in the `seqinr` package.

First, let's look at a dotplot created using only a single sequence. You'd never do this in practice, but it will give you a sense of what dotplots are doing.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

These two sequences are identical, so we have a very distinct diagonal line. But there's also other

Now we'll make a real dotplot of the chorismate lyase proteins from two closely related species, *Mycobacterium leprae* and *Mycobacterium ulcerans*.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0JzZXFpbmI6OmRvdFBsb3Qoc2VxMSwgXG4gICAgICAgICAgICAgICAgc2Vx
```

In the dotplot above, the *M. leprae* sequence is plotted along the x-axis (horizontal axis), and the *M. ulcerans* sequence is plotted along the y-axis (vertical axis). The dotplot displays a dot at points where there is an identical amino acid in the two sequences.

For example, if amino acid 53 in the *M. leprae* sequence is the same amino acid (eg. "W") as amino acid 70 in the *M. ulcerans* sequence, then the dotplot will show a dot the position in the plot where $x = 50$ and $y = 53$.

In this case you can see a lot of dots along a diagonal line, which indicates that the two protein sequences contain many identical amino acids at the same (or very similar) positions along their lengths. This is what you would expect, because we know that these two proteins are **homologs** (related proteins) because they share a close evolutionary history.

Chapter 21

Global proteins alignments in *R*

By: Avril Coghlan.

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

21.1 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKVxuXG5saWJyYXJ5KEJpb3N0cmZ3

21.1.1 Download sequences

As we did in the previous lesson on dotplots, we'll look at two sequences.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIERvd25sb2FkXG4jIyBzZXF1ZW5jZSAxOiBROUNEODNcbmxlcHJhZV9m

21.2 Pairwise global alignment of DNA sequences using the Needleman-Wunsch algorithm

If you are studying a particular pair of genes or proteins, an important question is to what extent the two sequences are similar.

To quantify similarity, it is necessary to **align** the two sequences, and then you can calculate a similarity score based on the alignment.

There are two types of alignment in general. A **global alignment** is an alignment of the *full* length of two sequences from beginning to end, for example, of two protein sequences or of two DNA sequences. A **local alignment** is an alignment of part of one sequence to part of another sequence; the parts the end up getting aligned are the most similar, and determined by the alignment algorithm.

The first step in computing a alignment (global or local) is to decide on a **scoring system**. For example, we may decide to give a score of +2 to a match and a penalty of -1 to a mismatch, and a penalty of -2 to a **gap** due to an **indexl**. Thus, for the alignment:

```
## [1] "G A A T T C"
```

```
## [1] "G A T T - A"
```

we would compute a score of

1. G vs G = match = 2
2. A vs A = match = 2

3. A vs T = mismatch = -1
4. T vs T = match = 2
5. T vs - = gap = -2
6. C vs A = mismatch = 2

So, the scores is $2 + 2 - 1 + 2 - 2 - 1 = 2$.

Similarly, the score for the following alignment is $2 + 2 - 2 + 2 + 2 - 1 = 5$:

```
## [1] "G A A T T C"
```

```
## [1] "G A - T T A"
```

The **scoring system** above can be represented by a **scoring matrix** (also known as a **substitution matrix**). The scoring matrix has one row and one column for each possible letter in our alphabet of letters (e.g. 4 rows and 4 columns for DNA and RNA sequences, 20 x 20 for amino acids). The (i,j) element of the matrix has a value of +2 in case of a match and -1 in case of a mismatch.

We can make a scoring matrix in R by using the `nucleotideSubstitutionMatrix()` function in the **Biostrings** package. **Biostrings** is part of a set of R packages for bioinformatics analysis known as Bioconductor (www.bioconductor.org/).

The arguments (inputs) for the `nucleotideSubstitutionMatrix()` function are the score that we want to assign to a match and the score that we want to assign to a mismatch. We can also specify that we want to use only the four letters representing the four nucleotides (ie. A, C, G, T) by setting `baseOnly=TRUE`, or whether we also want to use the letters that represent **ambiguous cases** where we are not sure what the nucleotide is (e.g. 'N' = A/C/G/T; ambiguous cases occur in some sequences due to sequencing errors or ambiguities).

To make a scoring matrix which assigns a score of +2 to a match and -1 to a mismatch, and store it in the variable `sigma`, we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIG1ha2UgdGhlIG1hdHJpeFxc2lnbWEgPC0gbnVjbGVvdGlkZVN1YnN0aXFI
```

Instead of assigning the same penalty (e.g. -8) to every gap position, it is common instead to assign a **gap opening penalty** to the first position in a gap (e.g. -8), and a smaller **gap extension penalty** to every subsequent position in the same gap.

The reason for doing this is that it is likely that adjacent gap positions were created by the same insertion or deletion event, rather than by several independent insertion or deletion events. Therefore, we don't want to penalize a 3-letter gap (AAA—AAA) as much as we would penalize three separate 1-letter gaps (AA-A-A-AA), as the 3-letter gap may have arisen due to just one insertion or deletion event, while the 3 separate 1-letter gaps probably arose due to three independent insertion or deletion events.

For example, if we want to compute the score for a global alignment of two short DNA sequences 'GAATTC' and 'GATTA', we can use the **Needleman-Wunsch** algorithm to calculate the highest-scoring alignment using a particular scoring function.

The `pairwiseAlignment()` function in the **Biostrings** package finds the score for the optimal global alignment between two sequences using the Needleman-Wunsch algorithm, given a particular scoring system.

As arguments (inputs), `pairwiseAlignment()` takes

1. the two sequences that you want to align,
2. the scoring matrix,
3. the gap opening penalty, and
4. the gap extension penalty.

You can also tell the function that you want to just have the optimal global alignment's score by setting `scoreOnly = TRUE`, or that you want to have both the optimal global alignment and its score by setting `scoreOnly = FALSE`.

21.3. PAIRWISE GLOBAL ALIGNMENT OF PROTEIN SEQUENCES USING THE NEEDLEMAN-WUNSCH ALGORITHM

For example, let's find the score for the optimal global alignment between the sequences 'GAATTC' and 'GATTA'.

First, we'll store the sequences as **character vectors**:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzMSA8LSBcIkdBQVRUQ1wiXG5zMiA8LSBcIkdBVFRBXClifQ==
```

Now we'll align them:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbnMxczIgPC0gQmlvc3RyaW5nczo6cGFpcndpc2VBbGlnbnMxczIgIn0=
```

The output:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbnMxczIgIn0=
```

The above commands print out the optimal global alignment for the two sequences and its score.

Note we set **gapOpening** to be -2 and **gapExtension** to be -8, which means that the first position of a gap is assigned a score of $-8 - 2 = -10$, and every subsequent position in a gap is given a score of -8. Here the alignment contains four matches, one mismatch, and one gap of length 1, so its score is $(4 \cdot 2) + (1 \cdot -1) + (1 \cdot -10) = -3$.

21.3 Pairwise global alignment of protein sequences using the Needleman-Wunsch algorithm

As well as DNA alignments, it is also possible to make alignments of protein sequences. In this case it is necessary to use a scoring matrix for amino acids rather than for nucleotides.

21.3.1 Protein score matrices

There are several well known scoring matrices that come with *R*, such as the **BLOSUM** series of matrices. Different BLOSUM matrices exist, named with different numbers. BLOSUM with high numbers are designed for comparing closely related sequences, while BLOSUM with low numbers are designed for comparing evolutionarily distantly related sequences. For example, **BLOSUM62** is used for **less divergent alignments** (alignments of sequences that differ little), and **BLOSUM30** is used for more divergent alignments (alignments of sequences that differ a lot).

Many *R* packages come with example data sets or data files and you use the **data()** function is used to load these data files. You can use the **data()** function to load a data set of BLOSUM matrices that comes with **Biostrings**

To load the BLOSUM50 matrix, we type:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkYXRhKEJMT1NVTTUwKVxuQkxPU1VNNTAgIyBQcmludCBvdXQgdGhl
```

You can get a list of the available scoring matrices that come with the Biostrings package by using the **data()** function, which takes as an argument the name of the package for which you want to know the data sets that come with it:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkYXRhKHBhY2thZ2U9XCJCcW9zdHJpbmdzXCIPIn0=
```

Another well-known series of scoring matrices are the **PAM** matrices developed by Margaret Dayhoff and her team. These have largely been replaced by BLOSUM but are important for historical reasons because they represent one of the first major bioinformatics, computational biology, and phylogenetics projects ever.

21.3.2 Example protein alignment

Let's find the optimal global alignment between the protein sequences "PAWHEAE" and "HEAGAWGHEE" using the Needleman-Wunsch algorithm using the BLOSUM50 matrix.

First, load the scoring matrix BLOSUM50 and make vectors for the sequence

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIG1hdHJpeFxuZGF0YShCTE9TVU01MClclblxuIyBzZXF1ZW5jZXNcbnMz
```

Now do the alignments.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbnMzczQgPC0gcGFpcndpc2VBbGlnbnM1bnQoczMsIHM0L
```

Look at the results:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbnMzczQgIyBQcmlludCBvdXQgdGhlIG9wdGltYWwgZ2xv
```

We set `gapOpening` to be -2 and `gapExtension` to be -8, which means that the first position of a gap is assigned a score of $-8-2=-10$, and every subsequent position in a gap is given a score of -8. This means that the gap will be given a score of $-10-8-8 = -26$.

21.4 Aligning UniProt sequences

We discussed previously how you can search for UniProt accessions and retrieve the corresponding protein sequences, either via the UniProt website or using the **rentrez** package.

In the examples given above, we learned how to retrieve the sequences for the chorismate lyase proteins from *Mycobacterium leprae* (UniProt Q9CD83) and *Mycobacterium ulcerans* (UniProt A0PQ23), and read them into R, and store them as vectors `lepraeseq` and `ulceransseq`.

You can align these sequences using `pairwiseAlignment()` from the Biostrings package.

As its input, the `pairwiseAlignment()` function requires that the sequences be in the form of a single string (e.g. “ACGTA”), rather than as a vector of characters (e.g. a vector with the first element as “A”, the second element as “C”, etc.). Therefore, to align the *M. leprae* and *M. ulcerans* chorismate lyase proteins, we first need to convert the vectors `lepraeseq` and `ulceransseq` into strings. We can do this using the `paste()` function:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIGNvb3ZlcnQgbGVwcmFLX3ZlY3RvcjB0byBhbiBvYmplY3QgbGVwcmFlc2V
```

Furthermore, `pairwiseAlignment()` requires that the sequences be stored as uppercase characters. Therefore, if they are not already in uppercase, we need to use the `toupper()` function to convert `lepraeseq_string` and `ulceransseq_string` to uppercase:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsZXByYWVzZXFfc3RyaW5nICAgPC0gdG91cHBlcjB0byBhbiBvYmplY3QgbGVwcmFlc2V
```

Check the output

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsZXByYWVzZXFfc3RyaW5nICMgUHJpbnQgb3V0IHRoZSBjb250ZW50IG9n
```

We can now align the the *M. leprae* and *M. ulcerans* chorismate lyase protein sequences using the `pairwiseAlignment()` function:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbkxlcHJhZVVzY2VzYW5zIDwtIEJpb3N0cmllZ3M6OnBha
```

The output:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnbG9iYWxBbGlnbkxlcHJhZVVzY2VzYW5zICMgUHJpbnQgb3V0IHRoZSBv
```

As the alignment is very long, when you type `globalAlignLepraeUlcerans`, you only see the start and the end of the alignment. Therefore, we need to have a function to print out the whole alignment (see below).

21.5 Viewing a long pairwise alignment

If you want to view a long pairwise alignment such as that between the *M. leprae* and *M. ulerans* chorismate lyase proteins, it is convenient to print out the alignment in blocks.

The R function `printPairwiseAlignment()` below will do this for you:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwcmludFBhaXJ3aXNlQWxpZ25tZW50KGdsb2JhbEFsaWduTG9wcmFlVWw
```

The position in the protein of the amino acid that is at the end of each line of the printed alignment is shown after the end of the line. For example, the first line of the alignment above finishes at amino acid position 50 in the *M. leprae* protein and also at amino acid position 60 in the *M. ulcerans* protein. Because there is a difference of $60 - 50 = 10$ bases, there must be 10 insertions in the *M. leprae* to get it to line up. Count the number of dashes in the sequence to see how many there are.

Chapter 22

Local protein alignments in *R*

By: Avril Coghlan.

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

22.1 Preliminaries

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKVxubGlicmFyeShCaW9zdHJpbmdzKSJ
```

22.1.1 Download sequences

As we did in the previous lesson on dotplots, we'll look at two sequences.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIERvd25sb2FkXG4jIyBzZXF1ZW5jZSAxOiBROUNEODNcbmxlcHJhZV9m
```

22.2 Pairwise local alignment of protein sequences using the Smith-Waterman algorithm

You can use the `pairwiseAlignment()` function to find the optimal **local alignment** of two sequences, that is the best alignment of parts (**subsequences**) of those sequences, by using the `type=local` argument in `pairwiseAlignment()`. This uses the **Smith-Waterman algorithm** for local alignment. This is the classic bioinformatics algorithm for finding optimal local alignments. (We'll discuss updated approaches when we get into **database searches** with **BLAST**, the *Basic, Local Alignment Search Tool* that is the workhorse of many day-to-day bioinformatics tasks).

For example, to find the best local alignment between the *M. leprae* and *M. ulcerans* chorismate lyase proteins, we can run:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIGxvYWQgc2NvcmluZyBtYXRyaXhcbmRh dGEoQkxPU1VNNTApXG5cbiM
```

Print out the optimal local alignment and its score

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsb2Nh bEFsaWduTG VwcmFIVWxjZXJhbnMgIn0=
```

As before, we can print out the full alignment with `printPairwiseAlignment()`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwcm ludFBhaXJ3aXNlQWxpZ25tZW50KGxvY2FsQWxpZ25MZXB yYWVv
```

We see that the optimal local alignment is quite similar to the optimal global alignment in this case, except that it excludes a short region of poorly aligned sequence at the start and at the ends of the two proteins.

Chapter 23

Retrieving multiple sequences in R

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKSJ9

By: Avril Coghlan.

Multiple Alignment and Phylogenetic trees <https://a-little-book-of-r-for-bioinformatics.readthedocs.io/en/latest/src/chapter5.html>

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

23.1 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIjYBwYWNrYWdlXG5saWJyYXJ5KGNvbXBiaW80YWxsKVxubGliemFyeSh

23.2 Retrieving a set of sequences from UniProt

Using websites or *R* you can search for DNA or protein sequences in sequence databases such as the **NCBI** database and **UniProt**. Oftentimes, it is useful to retrieve several sequences at once. The *R* function `entrez_fetch()` from the *rentrez* package is useful for this purpose. Other packages can also, such as *sequinr* this but *rentrez* has the cleanest interface.

We'll retrieve the protein sequences for these UniProt accessions

1. P06747: rabies virus phosphoprotein
2. P0C569: Mokola virus phosphoprotein
3. O56773: Lagos bat virus phosphoprotein
4. Q5VKP1: Western Caucasian bat virus phosphoprotein

Rabies virus is the virus responsible for rabies, which is classified by the WHO as a **neglected tropical disease**. Rabies is not a major human pathogen in the USA and Europe, but is problem in Africa. [Mokola virus]([https://en.wikipedia.org/wiki/Mokola_lyssavirus\(\)](https://en.wikipedia.org/wiki/Mokola_lyssavirus())) and rabies virus are closely related viruses that both belong to a group of viruses called the Lyssaviruses. Mokola virus causes a rabies-like infection in mammals including humans.

You can type make a vector containing the names of the sequences. Note that the accessions aren't numbers but are **quoted character strings**:

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIjZXFuYW1lcyA8LSBjKFwiUDA2NzQ3XCIsIFxuICAgICAgICAgICAgICBcll

Confirm that we are working with character data using `is.character()`

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIjYBwYWNrYWdlXG5saWJyYXJ5KGNvbXBiaW80YWxsKVxubGliemFyeSh


```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy5saXN0KHNlcV8xXzJfM180KSJ9
```

The size of the list is the number of elements it contains, not the amount of data in each element. There are 4 sequences, so 4 elements, so `length = 4`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZW5ndGoc2VxXzFfMl8zXzQpIn0=
```

We can access each element of the list by name, like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZXFfMV8yXzNfNCRQMDY3NDcifQ==
```

We can also access it by its index number, like this, using **double-bracket notation**.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZXFfMV8yXzNfNFtbMV1dICNOT1RFOiBkb3VibGUgYnJhY2tldHMifQ==
```

Each element of the list is a vector. We can check this using `is.vector()` like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy52ZWN0b3Ioc2VxXzFfMl8zXzQkUDA2NzQ3KSJ9
```

or using double brackets like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpcy52ZWN0b3Ioc2VxXzFfMl8zXzRbWzFdXS kifQ==
```

Its character data, which we can confirm with `class()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJjbGFzcyhzZXFfMV8yXzNfNCRQMDY3NDcpICMgZG9sbGFyIH NpZ24gbm9
```


Chapter 24

Multiple sequence alignment in R

By: Nathan Brouwer, with some content adapted Coghlan (2011) Multiple Alignment and Phylogenetic trees and under the Creative Commons 3.0 Attribution License (CC BY 3.0). Functions `print_msa()` and `clean_alignment()` adapted from (Coglan 2011).

24.1 Preliminaries

24.1.1 Packages

We'll be using the package `ggmsa` for the first time and you will have to install it using `install.packages("ggmsa")`. You may be asked to re-start R more than once during the installation process.

eyJ5YW5ndWFnZSI6InIiLCJzYW1wbGUiOiJpIG5ldyBwYWNRYWdlc1xuIyMgT25seSBpbnN0YWxsIG9uY2VcbiMgaW5zdD

24.1.2 Functions

The following key functions from `compbio4all` are used in this lesson

- `fasta_cleaner()`
- `entrez_fetch_list()`
- `print_msa()`
- `clean_alignment()`

24.2 Multiple sequence alignment (MSA)

A common task in bioinformatics is to download a set of related sequences from a database, and then to align those sequences using multiple alignment software. This is the first step in almost all phylogenetic analyses using sequence data.

24.3 Make MSA with `msa()`

We'll use a package called `msa` (Bodenhofer et al. 2015). There are several packages that can do multiple sequence alignment in R, but they all require loading an external piece of alignment software that is just accessed via R. The `msa` package actually runs the alignment algorithms entirely in R, making workflows simpler.

24.4 Viewing your MSA

There are several ways to view and explore your MSA

1. Within the R console using `compbio4all::print_msa()`
2. As an R plot using `ggmsa::ggmas()`
3. OPTIONAL: As a PDF using `msa::msaPrettyPrint()`

24.4.1 Viewing a long multiple alignment in the R console.

If you want to view a long multiple alignment within the R console, it is convenient to view the multiple alignment in blocks.

The function `print_msa()` (Coglan 2011) below will do this for you. As its inputs, the function `print_msa()` takes the two things

1. `alignment`: input alignment
2. `chunksize`: the number of columns to print out in each block.

To use `print_msa()` we first need to do a little format conversion:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ2aXJ1c2Fsbl9zZXFpbmIgPC0gbXNhQ29udmVydCh2aXJ1c2FsbiwgdHlwZSA9
```

Then we can print it out like this, making the alignment 60 bases wide:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwcmVudF9tc2EoYWxpZ25tZW50ID0gdmlldXNhbG5fc2VxaW5yLCBcbiAgIC
```

24.4.2 Visualizing alignments as an R plot

A powerful tool for visualizing focal parts of an alignment is `ggmsa`. If you haven't already, download it with `install.packages("ggmsa")` and load it with `library(ggmsa)`.

`ggmsa` prints a sequence alignment out within RStudio. Alignments can be large, so its important to select a subset of the alignment for visualization.

First, let's look at the first 20 bases of our alignment. Note that we are using `virusaln`, NOT `virusaln_seqinr` (sorry for the back and forth between objects.)

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnZ21zYSh2aXJ1c2FsbiwgICAjIHZpcnVzYWxuLCBOT1QgdmlldXNhbG5fc2
```

24.4.2.1 OPTIONAL: File types used by `ggmsa`

The `ggmsa` packages currently only works with certain types of alignment output. We can see what these are with `available_msa()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJhdmFpbGFibGVfbXNhKCkifQ==
```

As you can see there are a number of ways multiple sequence alignments can be represented in *R*. This has to do with the facts that i) There are many pieces of software / algorithms for making MSAs, and many bioinformatics packages that work with them.,

You can see that `AAMultipleAlignment` is listed, which the the format we set previously using the `class()` command.

The `msa` packages has a function `msaConvert()` which can change formats between different ways of representing MSAs which may be useful.

24.4.3 OPTIONAL: Print MSA to PDF

The `msa` package has a fabulous function, `msaPrettyPrint()` for rendering an MSA to PDF. It can take a little bit to run, and in order to view the PDF you need to locate the output. (Again, we'll use `virusaln`, not `virusaln_seqinr`).

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtc2FQcmV0dHlQcmIudCh2aXJ1c2FsbiwgICAgICMgdmlydXNhbG4sIE5PVCI
```

On a Mac usually searching in Finder will locate the file even after it is just created. You can ask R where it is saving things using `getwd()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJnZXR3ZCgpIn0=
```

You can change where R is saving things using the RStudio menu by clicking on Session -> Set Working Directory -> Choose directory...

24.5 Discarding very poorly conserved regions from an alignment

It is often a good idea to discard very **poorly conserved** regions from a multiple sequence alignment before visualizing it or building a phylogenetic tree, as the very poorly conserved regions are likely to be regions that are either **non-homologous** between the sequences being considered (and so do not have any phylogenetic signal), or are homologous but are so **diverged** that they are very difficult to align accurately (and so may add noise to the phylogenetic analysis, and decrease the accuracy of the inferred tree).

To discard very poorly conserved regions from a multiple alignment, you can use the following R function, `clean_alignment()` ((Coglan 2011))

The function `clean_alignment()` takes three arguments (inputs):

1. the input alignment;
2. `minpcnongap`: the minimum percent of letters in an alignment column that must be non-gap characters for the column to be kept; and
3. `minpcid`: the minimum percent of pairs of letters in an alignment column that must be identical for the column to be kept.

For example, if we have a single column (**locus**) with letters “T”, “A”, “T”, “-” (in four sequences), then 75% of the letters are non-gap characters; and the pairs of letters between the three non-gap sequences are

- 1 versus 2: “T,A”,
- 1 versus 3: “T,T”,
- 2 versus 3: “A,T”,

Therefore 33% of the pairs of letters are identical (**PID**) for that position in the alignment.

If you look at the multiple alignment for the virus phosphoprotein sequences (which we printed out using function `print_msa()`, see above), we can see that the last few columns are poorly aligned (contain many gaps and mismatches), and probably add noise to the phylogenetic analysis.

Let’s cleave off anything with less than 30% non-gap and less than 30% PID. NOTE: we’re back to using `virusaln_seqinr`, not `virusaln`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ2aXJ1c2Fsbl9zZXFpbnJfY2x1YW4gPC0gY2x1YW5fYWxpZ25tZW50KGfSaW
```

In this case, we required that at least 30% of letters in a column are not gap characters for that column to be kept, and that at least 30% of pairs of letters in an alignment column must be identical for the column to be kept.

We can print out the filtered alignment by typing:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJwcmIudF9tc2EodmlydXNhbG5fc2VxaW5yX2NsZWFuKSJ9
```

The filtered alignment is shorter, and is missing some of the poorly conserved regions of the original alignment.

Note that it is not a good idea to filter out too much of your alignment, as if you are left with few columns in your filtered alignment, you will be basing your phylogenetic tree upon a very short alignment (little data), and so the tree may be unreliable. Therefore, you need to achieve a balance between discarding the dodgy (poorly aligned) parts of your alignment, and retaining enough columns of the alignment that you will have enough data to base your tree upon.

Chapter 25

Calculating genetic distances between sequences

By: Nathan Brouwer, with some content adapted Coghlan (2011) Multiple Alignment and Phylogenetic trees and under the Creative Commons 3.0 Attribution License (CC BY 3.0). Functions `print_msa()` and `clean_alignment()` adapted from (Coglan 2011).

25.1 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKVxubGlicmFyeShtc2EpXG5saWJyYXJ5

25.2 Introduction

A common first step in performing a **phylogenetic analysis** is to calculate the **pairwise genetic distances** between sequences. The **genetic distance** is an estimate of the evolutionary **divergence** between two sequences, and is usually measured in quantity of evolutionary change, e.g., an estimate of the number of mutations that have occurred since the two sequences shared a **common ancestor**.

We can calculate the genetic distances between protein sequences using the `dist.alignment()` function in the `seqinr` package. The `dist.alignment()` function takes a multiple sequence alignment (MSA) as input. Based on the MSA that you give it, `dist.alignment()` calculates the genetic distance between each **pair** of proteins in the multiple alignment, yielding pairwise distances. For example, to calculate genetic distances between the virus phosphoproteins based on the multiple sequence alignment stored in `virusaln`, we type:

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkYXRhKHZpcnVzYWxuKVxudmlydXNkaXN0IDwtIHNlcWlucjo6ZGlzdC5h

NOTE My result are different from the original results, shown here: I need to check the settings used for the MSA

```
P0C569 O56773 P06747 O56773 0.4142670 P06747 0.4678196 0.4714045 Q5VKP1 0.4828127 0.5067117
0.5034130
```

The genetic distance matrix above shows the genetic distance between each pair of proteins. The sequences are referred to by their **UniProt accessions**. Recall that

- P06747 = rabies virus phosphoprotein
- P0C569 is Mokola virus phosphoprotein
- O56773 is Lagos bat virus phosphoprotein
- Q5VKP1 is Western Caucasian bat virus phosphoprotein.

Based on the genetic **distance matrix** above, we can see that the genetic distance between Lagos bat virus phosphoprotein (O56773) and Mokola virus phosphoprotein (P0C569) is smallest (about 0.414). Similarly, the genetic distance between Western Caucasian bat virus phosphoprotein (Q5VKP1) and Lagos bat virus phosphoprotein (O56773) is the biggest (about 0.507).

The larger the genetic distance between two sequences, the more amino acid changes (such as change from Asp to Met) or **indels** that have occurred since they shared a common ancestor, and the longer ago their **common ancestor** probably lived. (The relationship between number of mutations and time, however, depends on the mutation rate and generation time of the organism).

25.3 Calculating genetic distances between DNA/mRNA sequences

Just like for protein sequences, you can calculate genetic distances between DNA (or mRNA) sequences based on an alignment of the sequences. The RefSeq DNA accession numbers for the proteins we've been using are:

- AF049118 = mRNA sequence for Mokola virus phosphoprotein,
- AF049114 = mRNA sequence for Mokola virus phosphoprotein,
- AF049119 = mRNA sequence for Lagos bat virus phosphoprotein,
- AF049115 = mRNA sequence for Duvnhage virus phosphoprotein.

We can retrieve these DNA sequences using `entrez_fetch_list()`. Some notes about how we'll use this function works:

1. `db` is short for "database"
2. the database is called `nucore` (not `genebank` or `gene`)
3. the argument `rettype` is "tt"; I think it stands for "REtURN type" (I also forget the second "t")

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJiHB1dCBhY2Nlc3Npb25zIGluIHZlY3RvclxuYWNjZXNzaW9uc19tcm5hIDwt
```

We can clean these three sequences using a simple `for()` loop. We set `parse = F` so we get things back as single character string.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJmb3IoaSBpbA8LSBjKEFGMDQ5MTE4ID0gIHZpcnVzX21ybmFfbm
```

Now we need to convert each of these into a named vector of sequences.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtcmFfc2VxX3ZlY3RvciA8LSBjKEFGMDQ5MTE4ID0gIHZpcnVzX21ybmFfbm
```

Finally convert this to a "stringset" using `Biostrings::DNAStringSet()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkbmFfc2VxX3N0cmIuZ3NldCA8LSBCaW9zdHJpbmdzOjppETkFTdHJpbmdT
```

Let's see what we got

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkbmFfc2VxX3N0cmIuZ3NldCJ9
```

Now we can make an alignment use `msa()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ2aXJ1c19tcm5hX2FsbiA8LSBtc2EoaW5wdXRtZXZlID0gZG5hX3NlcV9zdH
```

The output looks like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ2aXJ1c19tcm5hX2FsbiJ9
```

This looks a LOT different than an amino acid alignment, which looked like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkYXRhKHZpcnVzYWxuKVxudmlydXNhbg4ifQ==
```

Why might they be different? First, examine the output above and determine how long the DNA alignment is versus the amino acid alignment? Why are they different, and why is one longer than the other?

The DNA alignment is 1097 columns, while the amino acid alignment is only 306 rows. Note that $306 \times 3 = 918$. 1097 is pretty close to 1097. What's the relevance of multiplying by 3?

25.4 Calculationg genetic distance

You can calculate a genetic distance for DNA or mRNA sequences using the `dist.dna()` function in the `ape` package. `dist.dna()` takes a MSA of DNA or mRNA sequences as its input, and calculates the genetic distance between each pair of DNA sequences in the multiple alignment.

The `dist.dna()` function requires the input alignment to be in a special format known as **DNAbin** format, so we must use the `as.DNAbin()` function to convert our DNA alignment into this format before using the `dist.dna()` function.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIgIyBDb252ZXJ0IHRoZSBhbGlnbm1lbnQgdG8gXCJETkFiaW5cIiBmb3JtYXN0
```

The output of `as.DNAbin()` gives us a short summery of the alignment

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJ2aXJ1c19tcm5hX2Fsbl9iaW4ifQ==
```

Now to make and view the alignment:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIgIyBDYWxjdWxhdGUgdGhlIGdlbmV0aWMgZGlzdGFuY2UgbWF0cm14XG50
```

NOTE: my results for this alignment are the same as the original by Coghlan. I'm not sure why my amino acid alignment produces divergent results but the DNA is the same.

Chapter 26

Unrooted neighbor-joining phylogenetic trees

NOTE: the code for this chapter works as intended but there are some differences between my results and what is reported by the original author of the chapter. This is likely to do with different alignment software, though it could just be a typo.

By: Avril Coghlan. Multiple Alignment and Phylogenetic trees <https://a-little-book-of-r-for-bioinformatics.readthedocs.io/en/latest/src/chapter5.html>

Adapted, edited and expanded: Nathan Brouwer under the Creative Commons 3.0 Attribution License (CC BY 3.0).

26.1 Preliminaries

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGNvbXBiaW80YWxsKVxubGlicmFyeShzZXFpbnIpIn0=

You will need to install the `ape` package if you do not have it already using `install.packages("ape")`.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsaWJyYXJ5KGFwZSkifQ==

26.1.1 Key functions

- `compbio4all::unrooted_NJ_tree` (Coghlan 200x)

26.1.2 Key vocab

- clade
- bootstrap
- resample
- rooted vs. unrooted tree
- outgroup

26.2 Building an unrooted phylogenetic tree for protein sequences

Once we have a **distance matrix** that gives the **pairwise distances** between all our protein sequences, we can build a **phylogenetic tree** based on that distance matrix. One method for using this is the **neighbor-joining algorithm**.

If we have the distance matrix already made we can make the tree like this using `ape::nj()`. The distance matrix is saved in `compbio4all` as `virus_mrna_dist`. Load this with `data()`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjIjIGxvYWQgdGhIIGRpc3RhbmNIIG1hdHJpeFxuZGF0YSh2aXJ1c19tcm5hX2
```

26.2.1 Build tree with `unrooted_NJ_tree()`

Coghlan (2011) wrote a function to simplify the steps of making an NJ tree. The R function `unrooted_NJ_tree()` is a **wrapper** for functions from the `ape` package which builds a phylogenetic tree based on an alignment of sequences, using the NJ algorithm.

The `unrooted_NJ_tree()` function takes an alignment of sequences its input, calculates **pairwise distances** between the sequences based on the alignment behind the scenes, and then builds a phylogenetic tree based on the pairwise distances. It returns the phylogenetic tree, and also makes a plot of that tree. It also gives us information about to what extent the data in the original MSA support the evolutionary relationships shown in the tree.

The alignment is saved in `compbio4all` as `virusa1n` and can be loaded with the `data()` command.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjIjIGxvYWQgdGhIIGRpc3RhbmNIIG1hdHJpeFxuZGF0YSh2aXJ1c19tcm5hX2
```

Take a look at the structure of the data

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjIjIGxvYWQgdGhIIGRpc3RhbmNIIG1hdHJpeFxuZGF0YSh2aXJ1c19tcm5hX2
```

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjIjIGxvYWQgdGhIIGRpc3RhbmNIIG1hdHJpeFxuZGF0YSh2aXJ1c19tcm5hX2
```

Note that you need to specify that the type of sequences that you are using are protein sequences when you use `unrooted_NJ_tree()`, by setting `type=protein`.

We can see that Q5VKP1 (Western Caucasian bat virus phosphoprotein) and P06747 (rabies virus phosphoprotein) have been grouped together into a **clade** on the tree, and that O56773 (Lagos bat virus phosphoprotein) and P0C569 (Mokola virus phosphoprotein) are grouped together.

This is consistent with what we saw above in the genetic distance matrix, which showed that the genetic distance between Lagos bat virus phosphoprotein (O56773) and Mokola virus phosphoprotein (P0C569) is relatively small.

26.3 Bootstrap values indicate support for clades

In the plot, the numbers in blue boxes are **bootstrap values** for the nodes in the tree. A bootstrap value for a particular node in the tree gives an idea of the confidence that we have in the clade (group) defined by that node in the tree. If a node has a high bootstrap value (near 100%) then we are very confident that the clade defined by the node is correct, while if it has a low bootstrap value (near 0%) then we are not so confident.

Note that the fact that a bootstrap value for a node is high does not necessarily guarantee that the clade defined by the node is correct, but just tells us that it is quite likely that it is correct given the data and analysis we're using.

The bootstrap values are calculated by making many (for example, 100) random **resamples** of the alignment that the phylogenetic tree was based upon. Each resample of the alignment consists of a certain number x (e.g., 200) of randomly sampled *columns* from the alignment. Each resample of the alignment (e.g., 200 randomly sampled columns) forms a sort of fake alignment of its own, and a phylogenetic tree can be based upon the *resample*. We can make 100 random resamples of the alignment, and build 100 phylogenetic trees based on the 100 resamples. These 100 trees are known as the **bootstrap trees**. For each clade (grouping) that we see in our original phylogenetic tree, we can count in how many of the 100 bootstrap trees it appears. This is known as the **bootstrap value** for the clade in our original phylogenetic tree.

For example, if we calculate 100 random resamples of the virus phosphoprotein alignment, and build 100 phylogenetic trees based on these resamples, we can calculate the bootstrap values for each clade in the virus phosphoprotein phylogenetic tree.

NOTE: I am currently not able to reproduce these results:

In this case, the bootstrap value for the node defining the clade containing Q5VKP1 (Western Caucasian bat virus phosphoprotein) and P06747 (rabies virus phosphoprotein) is 25%, while the bootstrap value for node defining the clade containing of Lagos bat virus phosphoprotein (O56773) and Mokola virus phosphoprotein (P0C569) is 100%. The bootstrap values for each of these clades is the percent of 100 bootstrap trees that the clade appears in.

Therefore, we are very confident that Lagos bat virus and Mokola virus phosphoproteins should be grouped together in the tree. However, we are not so confident that the Western Caucasian bat virus and rabies virus phosphoproteins should be grouped together.

26.4 Branch lengths indicate divergence between sequences

The lengths of the branches in the plot of the tree are proportional to the amount of evolutionary change (estimated number of mutations) along the branches. In this case, the branches leading to Lagos bat virus phosphoprotein (O56773) and Mokola virus phosphoprotein (P0C569) from the node representing their common ancestor are slightly shorter than the branches leading to the Western Caucasian bat virus (Q5VKP1) and rabies virus (P06747) phosphoproteins from the node representing their common ancestor.

This suggests that there might have been more mutations in the Western Caucasian bat virus (Q5VKP1) and rabies virus (P06747) phosphoproteins since they shared a common ancestor, than in the Lagos bat virus phosphoprotein (O56773) and Mokola virus phosphoprotein (P0C569) since they shared a common ancestor.

26.5 Unrooted trees lack an outgroup

The tree above of the virus phosphoproteins is an **unrooted** phylogenetic tree as it does not contain an **outgroup** sequence; that is, a sequence of a protein that is known to be more distantly related to the other proteins in the tree than they are to each other.

As a result, we cannot tell which direction evolutionary time ran in along the internal branches of the tree. For example, we cannot tell whether the node representing the common ancestor of (O56773, P0C569) was an ancestor of the node representing the common ancestor of (Q5VKP1, P06747), or the other way around.

In order to build a **rooted** phylogenetic tree, we need to have an outgroup sequence in our tree. In the case of the virus phosphoproteins, this is unfortunately not possible, as there is not any protein known that is more distantly related to the four proteins already in our tree than they are to each other.

However, in many other cases, an outgroup - a sequence known to be more distantly related to the other sequences in the tree than they are to each other - is known, and so it is possible to build a rooted phylogenetic tree.

We discussed above that it is a good idea to investigate whether discarding the poorly conserved regions of a multiple alignment has an effect on the phylogenetic analysis. In this case, we made a filtered copy of the multiple alignment and stored it in the variable `viru$seqinr_clean` (see above). We can make a phylogenetic tree based this filtered alignment, and see if it agrees with the phylogenetic tree based on the original alignment:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJkYXRhKHZpcnVzYWxuX3NlcWlucl9jbGVhbilcbmNsZWFuZWWR2aXJ1c2Fsb
```

As in the phylogenetic tree based on the raw (unfiltered) multiple alignment, O56773 and P0C569 are still grouped together, and Q5VKP1 and P06747 are still grouped together. Thus, filtering the multiple alignment does not have an effect on the tree. The bootstrap value, however, have changed.

If we had found a difference in the trees made using the unfiltered and filtered multiple alignments, we would have to examine the multiple alignments closely, to see if the unfiltered multiple alignment contains a lot of very poorly aligned regions that might be adding noise to the phylogenetic analysis (if this is true, the tree based on the filtered alignment is likely to be more reliable).

Chapter 27

A complete bioinformatics workflow in R

By: Nathan L. Brouwer

This lesson walks you through an entire workflow for a bioinformatics, including

1. obtaining FASTA sequences
2. cleaning sequences
3. creating alignments
4. creating distance a distance matrix
5. building a phylogenetic tree

We'll examine the Shroom family of genes, which produces Shroom proteins essential for tissue formation in many multicellular eukaryotes, including neural tube formation in vertebrates. We'll examine shroom in several very different organisms, including humans, mice and sea urchins. There is more than one type of shroom in vertebrates, and we'll also look at two different Shroom genes: shroom 1 and shroom 2.

This lesson draws on skills from previous sections of the book, but is written to act as an independent summary of these activities. There is therefore a review of key aspects of R and bioinformatics throughout it.

27.1 Software Preliminaires

27.1.1 Vocab

- argument
- function
- list
- named list
- vector
- named vector
- for() loop
- R console

27.1.2 R functions

- library()
- round()
- plot()
- mtext()

- `nchar()`
- `rentrez::entrez_fetch()`
- `combio4all::entrez_fetch_list()`
- `combio4all::print_msa()` (Coghlan 2011)
- `Biostrings::AAStringSet()`
- `msa::msa()`
- `msa::msaConvert()`
- `msa::msaPrettyPrint()`
- `seqinr::dist.alignment()`
- `ape::nj()`

A few things need to be done to get started with our R session.

27.1.3 Download necessary packages

Many R sessions begin by downloading necessary software packages to augment R's functionality.

If you don't have them already, you'll need the following packages from CRAN:

1. `ape`
2. `seqinr`
3. `rentrez`
4. `devtools`

The CRAN packages can be loaded with `install.packages()`.

You'll also need these packages from Bioconductor:

1. `msa`
2. `Biostrings`

For install packages from Bioconductor, see the chapter at the beginning of this book on this process.

Finally, you'll need this package from GitHub

1. `combio4all`

To install packages from GitHub you can use the code `devtools::install_github("brouwern/combio4all")`

27.1.4 Load packages into memory

We now need to load up all our bioinformatics and phylogenetics software into R. This is done with the `library()` command.

To run this code just click on the sideways green triangle all the way to the right of the code.

NOTE: You'll likely see some red code appear on your screen. No worries, totally normal!

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIGdpdGh1YiBwYWNRYWdlc1xubGlicmFyeShjb21wYmlvNGFsbClcblxuIyBI

27.2 Downloading macro-molecular sequences

We're going to explore some sequences. First we need to download them. To do this we'll use a function, `entrez_fetch()`, which accesses the **Entrez** system of database (ncbi.nlm.nih.gov/search/). This function is from the `rentrez` package, which stands for "R-Entrez."

We need to tell `entrez_fetch()` several things

1. `db = ...` the type of entrez database.
2. `id = ...` the **accession** (ID) number of the sequence

3. `rettype = ...` file type what we want the function to return.

Formally, these things are called **arguments** by *R*.

We'll use these settings:

1. `db = "protein"` to access the Entrez database of protein sequences
2. `rettype = "fasta"`, which is a standard file format for nucleic acid and protein sequences

We'll set `id = ...` to sequences whose **accession numbers** are:

1. NP_065910: Human shroom 3
2. AAF13269: Mouse shroom 3a
3. CAA58534: Human shroom 2
4. XP_783573: Sea urchin shroom

There are two highly conserved regions of shroom 3 1. ASD 1: aa 884 to aa 1062 in hShroom3 1. ASD 2: aa 1671 to aa 1955 in hShroom3

Normally we'd have to download these sequences by hand through pointing and clicking on GeneBank records on the NCBI website. In *R* we can do it automatically; this might take a second.

All the code needed is this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJlIEh1bWFnIHNoYm9vbSAzICChILiBzYXBpZW5zKVxuaFNocm9vbTMgPC0g
```

The output is in FASTA format; we'll use the `cat()` to do a little formatting for us:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJlYXQoaFNocm9vbTMpIn0=
```

Note the initial `>`, then the header line of `NP_065910.3 protein Shroom3 [Homo sapiens]`. After that is the amino acid sequence. The underlying data also includes the **newline character** `\n` to designate where each line of amino acids stops.

We can get the rest of the data by just chaining the `id = ...` argument:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJlIEh1bWFnIHNoYm9vbSAzYSAoTS4gbXVzY3VsdXMpXG5tU2hyb29tM2EgP
```

I'm going to check about how long each of these sequences is - each should have an at least slightly different length. If any are identical, I might have repeated an accession name or re-used an object name. The function `nchar()` counts of the number of characters in an *R* object.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJlY2h0cihoU2hyb29tMylcbm5jaGFyKG1TaHJvb20zYSIcbm5jaGFyKHNTaHJ
```

27.3 Prepping macromolecular sequences

“90% of data analysis is data cleaning” (-Just about every data analyst and data scientist on twitter)

We have our sequences, but the current format isn't directly usable for us yet because there are several things that aren't sequence information

1. metadata (the header)
2. page formatting information (the newline character)

We can remove this non-sequence information using a function I wrote called `fasta_cleaner()`, which is in the `compbio4all` package. The function uses **regular expressions** to remove the info we don't need.

ASIDE: If we run the name of the command with out any quotation marks we can see the code:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0OiJmYXN0YV9jbGVhbmVyIn0=
```

End ASIDE

Now use the function to clean our sequences; we won't worry about what `pare = ...` is for.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJoU2hyb29tMyAgPC0gZmFzdGFfY2x1YW5lcihoU2hyb29tMywgIHBhcnNIID0
```

Now let's take a peek at what our sequences look like:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJoU2hyb29tMyJ9
```

27.4 Aligning sequences

We can do a **global alignment** of one sequence against another using the `pairwiseAlignment()` function from the **Bioconductor** package **Biostrings** (note that capital "B" in **Biostrings**; most R package names are all lower case, but not this one).

Let's align human versus mouse shroom:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJhbGlnbi5oMy52cy5tM2EgPC0gQmlvc3RyaW5nczo6cGFpcndpc2VBbGlnbm1l
```

We can peek at the alignment

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJhbGlnbi5oMy52cy5tM2EifQ==
```

The **score** tells us how closely they are aligned; higher scores mean the sequences are more similar. Its hard to interpret the number on its own so we can get the **percent sequence identity (PID)** using the `pid()` function.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJCW9zdHJpbmdzOjpwYWQoYWxpZ24uaDMudnMubTNhKSJ9
```

So, *shroom3* from humans and *shroom3* from mice are ~71% similar (at least using this particular method of alignment, and there are many ways to do this!)

What about human shroom 3 and sea-urchin shroom?

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJhbGlnbi5oMy52cy5oMiA8LSBCaW9zdHJpbmdzOjpwYWlyd2lzZUFsaWdubV
```

First check out the score using `score()`, which accesses it directly without all the other information.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzY29yZShhbGlnbi5oMy52cy5oMikifQ==
```

Now the percent sequence alignment with `pid()`:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJCW9zdHJpbmdzOjpwYWQoYWxpZ24uaDMudnMuaDIpIn0=
```

So Human shroom 3 and Mouse shroom 3 are 71% identical, but Human shroom 3 and human shroom 2 are only 34% similar? How does it work out evolutionary that a human and mouse gene are more similar than a human and a human gene? What are the evolutionary relationships among these genes within the shroom gene family?

27.5 The shroom family of genes

I've copied a table from a published paper which has accession numbers for 15 different Shroom genes.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzaHJvb21fdGFibGUgPC0gYyhcIkNBQTc4NzE4XCIGLCBcIlguIGxhZXZpcyE
```

I'll do a bit of formatting; you can ignore these details if you want

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJlIGNvb3ZlcnQgdG8gbWF0cm14XG5zaHJvb21fdGFibGVfbWF0cm14IDwtIG1
```

Take a look:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzaHJvb21fdGFibGUifQ==
```

Instead of getting one sequence at a time we can download several by accessing the “accession” column from the table

We can give this whole set of accessions to `entrez fetch()`:

We can look at what we got here with `cat()` (I won't display this because it is very long!)

The current format of these data is a single, long set of data. This is a standard way to store, share and transmit FASTA files, but in *R* we'll need a slightly different format.

eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzaHJvb21zX2xpc3QgPC0gZW50cmV6X2ZldGNoX2xpc3QoZGIgPSBcInByb3

evJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJsZW5ndGgc2hyb29tc19saXN0KSJ9

evJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJmb3IoaSBpbjAxOmxiYmnd0aChzaHJvb21zX2xpc3QpKXtcbiAgc2hyb29tc19sa

evJsYW5ndWFnZSI6InLiLCJzYW1wbGUiOiIjIG1ha2UgYSB2ZWNo3IgdG8gc3RvcnUgb3V0cHV0XG5zaHJvb21zX3ZlY3.

evJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzaHJvb21zX3ZlY3Rvc19zcycA8LSBCaW9zdHJpbmdzOjpBQVN0cmIuZ1NldCI

We must **align** all of the sequences we downloaded and use that **alignment** to build a **phylogenetic tree**. This will tell us how the different genes, both within and between species, are likely to be related.

We'll use the software **msa**, which implements the **ClustalW** multiple sequence alignment algorithm. Normally we'd have to download the ClustalW program and either point-and-click our way through it or use the **command line***, but these folks wrote up the algorithm in R so we can do this with a line of R code. This will take a second or two.

27.7.2 Viewing an MSA

Once we build an MSA we need to visualize it.

27.7.2.1 Viewing an MSA in R

We can look at the output from `msa()`, but its not very helpful

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduIn0=
```

A function called `print_msa()` (Coghlan 2011) which I've put into `combio4all` can give us more informative output by printing out the actual alignment into the R console.

To use `print_msa()` We need to make a few minor tweaks though first. These are behind the scenes changes so don't worry about the details right now. We'll change the name to `shrooms_align_seqinr` to indicate that one of our changes is putting this into a format defined by the bioinformatics package `seqinr`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduKSA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

I won't display the output from `shrooms_align_seqinr` because its very long; we have 14 shroom genes, and shroom happens to be a rather long gene.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduKSA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

27.7.2.2 Displaying an MSA as an R plot

I'm going to just show about 100 amino acids near the end of the alignment, where there is the most overlap across all of the sequences. This is set with the `start = ...` and `end = ...` arguments. Note that we're using the `shrooms_align` object.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduKSA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

27.7.2.3 Saving an MSA as PDF

We can take a look at the alignment in PDF format if we want. In this case I'm going to just show about 100 amino acids near the end of the alignment, where there is the most overlap across all of the sequences. This is set with the `y = c(...)` argument.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduKSA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

You can see where R is saving things by running `getwd()`

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2FsaWduKSA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

On a Mac you can usually find the file by searching in Finder for the file name, which I set to be "shroom_msa.pdf" using the `file = ...` argument above.

27.8 Genetic distance.

Next need to first get an estimate of how similar each sequences is. The more amino acids that are identical to each other, the more similar.

Instead of similarity, we usually work in terms of *difference* or **genetic distance** (a.k.a. **evolutionary distance**). This is done with the `dist.alignment()` function.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2Rpc3QgPC0gc2VxaW5yOjpkZXN0LmFsaWduWVudChzaHJv
```

We've made a matrix using `dist.alignment()`; let's round it off so its easier to look at using the `round()` function.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2Rpc3Rfcm91bmRlZCA8LSBcIkFBTXVsdGlwbGVbBbGlnbm1lbn
```

Now let's look at it

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0IjZaHJvb21zX2Rpc3Rfcm91bmRlZCJ9
```

27.9 Phylogenetic trees (finally!)

We got our sequence, built multiple sequence alignment, and calculated the genetic distance between sequences. Now we are - finally - ready to build a phylogenetic tree.

First, we let R figure out the structure of the tree. There are **MANY** ways to build phylogenetic trees. We'll use a common one used for exploring sequences called **neighbor joining** algorithm via the function `nj()`. Neighbor joining uses genetic distances to cluster sequences into **clades**.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0cmVlIDwtIG5qKHNoem9vbXNfZGZldCkifQ==
```

27.9.1 Plotting phylogenetic trees

Now we'll make a quick plot of our tree using `plot()` (and add a little label using a function called `mtext()`).

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0cmVlIDwtIG5qKHNoem9vbXNfZGZldCkifQ==
```

This is an ****unrooted tree***. For the sake of plotting we've also ignored the evolutionary distance between the sequences.

To make a rooted tree we remove `type = "unrooted"`.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0cmVlIDwtIG5qKHNoem9vbXNfZGZldCkifQ==
```

We can include information about branch length by setting `use.edge.length = ...` to T.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0cmVlIDwtIG5qKHNoem9vbXNfZGZldCkifQ==
```

Some of the branches are now very short, but most are very long, indicating that these genes have been evolving independently for many millions of years.

Let's make a fancier plot. Don't worry about all the steps; I've added some more code to add some annotations on the right-hand side to help us see what's going on.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGU0cmVlIDwtIG5qKHNoem9vbXNfZGZldCkifQ==
```


Part I

Appendices

Appendix 01: Getting access to R

27.10 Getting Started With R and RStudio

- R is a piece of software that does calculations and makes graphs.
- RStudio is a GUI (graphical user interface) that acts as a front-end to R
- You can use R directly, but most people use a GUI of some kind
- RStudio has become the most popular GUI

The following instructions will lead you click by click through downloading R and RStudio and starting an initial session. If you have trouble with downloading either program go to YouTube and search for something like “Downloading R” or “Installing RStudio” and you should be able to find something helpful, such as “How to Download R for Windows”.

27.10.1 RStudio Cloud

TODO: Add RStudio cloud

27.10.2 Getting R onto your own computer

To get R on to your computer first go to the CRAN website at <https://cran.r-project.org/> (CRAN stands for “comprehensive R Archive Network”). At the top of the screen are three bullet points; select the appropriate one (or click the link below)

- Download R for Linux
- Download R for (Mac) OS X
- Download R for Windows

Each page is formatted slightly differently. For a current Mac, click on the top link, which as of 8/16/2018 was “R-3.5.1.pkg” or click this link. If you have an older Mac you might have to scroll down to find your operating system under “Binaries for legacy OS X systems.”

For PC select “base” or click this link.

When its downloaded, run the installer and accept the defaults.

27.10.3 Getting RStudio onto your computer

RStudio is an R interface developed by a company of the same name. RStudio has a number of commercial products, but much of their portfolio is freeware. You can download RStudio from their website www.rstudio.com/. The download page (www.rstudio.com/products/rstudio/download/) is a bit busy because it shows all of their commercial products; the free version is on the far left side of the table of products. Click on the big green DOWNLOAD button under the column on the left that says “RStudio Desktop Open Source License” (or click on this link).

This will scroll you down to a list of downloads titled “Installers for Supported Platforms.” Windows users can select the top option RStudio 1.1.456 - Windows Vista/7/8/10 and Mac the second option RStudio

1.1.456 - Mac OS X 10.6+ (64-bit). (Versions names are current of 8/16/2018).

Run the installer after it downloads and accept the default. RStudio will automatically link up with the most current version of R you have on your computer. Find the RStudio icon on your desktop or search for “RStudio” from your task bar and you’ll be read to go.

27.10.4 Keep R and RStudio current

Both R and RStudio undergo regular updates and you will occasionally have to re-download and install one or both of them. In practice I probably do this about every 6 months.

Getting started with R itself (or not)

Vocabulary

- console
- script editor / source viewer
- interactive programming
- scripts / script files
- .R files
- text files / plain text files
- command execution / execute a command from script editor
- comments / code comments
- commenting out / commenting out code
- stackoverflow.com
- the rstats hashtag

R commands

- `c(...)`
- `mean(...)`
- `sd(...)`
- `?`
- `read.csv(...)`

This is a walk-through of a very basic R session. It assumes you have successfully installed R and RStudio onto your computer, and nothing else.

Most people who use R do not actually use the program itself - they use a GUI (graphical user interface) “front end” that make R a bit easier to use. However, you will probably run into the icon for the underlying R program on your desktop or elsewhere on your computer. It usually looks like this:

ADD IMAGE HERE

The long string of numbers have to do with the version and whether is 32 or 64 bit (not important for what we do).

If you are curious you can open it up and take a look - it actually looks a lot like RStudio, where we will do all our work (or rather, RStudio looks like R). Sometimes when people are getting started with R they will accidentally open R instead of RStudio; if things don't seem to look or be working the way you think they should, you might be in R, not RStudio

27.10.4.1 R's console as a scientific calculator

You can interact with R's console similar to a scientific calculator. For example, you can use parentheses to set up mathematical statements like

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiI1KigxKzEpIn0=
```

Note however that you have to be explicit about multiplication. If you try the following it won't work.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiI1KDErMSkifQ==
```

R also has built-in functions that work similar to what you might have used in Excel. For example, in Excel you can calculate the average of a set of numbers by typing “=average(1,2,3)” into a cell. R can do the same thing except

- The command is “mean”
- You don't start with “=”
- You have to package up the numbers like what is shown below using “c(…)”

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKGMoMSwyLDMpKSJ9
```

Where “c(…)” packages up the numbers the way the mean() function wants to see them.

If you just do the following R will give you an answer, but its the wrong one

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKDEsMiwwKSJ9
```

This is a common issue with R – and many programs, really – it won't always tell you when somethind didn't go as planned. This is because it doesn't know something didn't go as planned; you have to learn the rules R plays by.

27.10.4.2 Practice: math in the console

See if you can reproduce the following results

Division

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIxMC8zIn0=
```

The standard deviation

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzZChjKDU5MTAsMTUpKSAjIG5vdGUgdGhlIHVzZSBvZiBcImMoLi4uKVwi
```

27.10.4.3 The script editor

While you can interact with R directly within the console, the standard way to work in R is to write what are known as **scripts**. These are computer code instructions written to R in a **script file**. These are save with the extension **.R** but area really just a form of **plain text file**.

To work with scripts, what you do is type commands in the script editor, then tell R to **excute** the command. This can be done several ways.

First, you tell RStudio the line of code you want to run by either * Placing the cursor at the end a line of code, OR * Clicking and dragging over the code you want to run in order highlight it.

Second, you tell RStudio to run the code by * Clicking the “Run” icon in the upper right hand side of the script editor (a grey box with a green error emerging from it) * pressing the control key (“ctrl”) and then then enter key on the keyboard

The code you've chosen to run will be sent by RStudio from the script editor over to the console. The console will show you both the code and then the output.

You can run several lines of code if you want; the console will run a line, print the output, and then run the next line. First I'll use the command mean(), and then the command sd() for the standard deviation:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKGMoMSwyLDMpKVxuc2QoYyglLDIsMykpIn0=
```

27.10.4.4 Comments

One of the reasons we use script files is that we can combine R code with **comments** that tell us what the R code is doing. Comments are preceded by the hashtag symbol `#`. Frequently we'll write code like this:

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKGMoMSwyLDMpKSJ9
```

If you highlight all of this code (including the comment) and then click on “run”, you’ll see that RStudio sends all of the code over console.

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

Comments can also be placed at the *end* of a line of code

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKGMoMSwyLDMpKSAjTm90ZSAgdGhlIHVzZSBvZiBjKC4uLikifQ=
```

Sometimes we write code and then don’t want R to run it. We can prevent R from executing the code even if its sent to the console by putting a “`#`” *infront* of the code.

If I run this code, I will get just the mean but not the sd.

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJtZWFuKGMoMSwyLDMpKVxuI3NkKGMoMSwyLDMpKSJ9
```

Doing this is called **commenting out** a line of code.

27.11 Help!

There are many resource for figuring out R and RStudio, including

- R’s built in “help” function
- Q&A websites like **stackoverflow.com**
- twitter, using the hashtag `#rstats`
- blogs
- online books and course materials

27.11.1 Getting “help” from R

If you are using a function in R you can get info about how it works like this

```
eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiI/bWVhbiJ9
```

In RStudio the help screen should appear, probably above your console. If you start reading this help file, though, you don’t have to go far until you start seeing lots of R lingo, like “S3 method”, “na.rm”, “vectors”. Unfortunately, the R help files are usually not written for beginners, and reading help files is a skill you have to acquire.

For example, when we load data into R in subsequent lessons we will use a function called “read.csv”

Access the help file by typing “`?read.csv`” into the console and pressing enter. Surprisingly, the function that R give you the help file isn’t what you asked for, but is `read.table()`. This is a related function to `read.csv`, but when you’re a beginner thing like this can really throw you off.

Kieran Healy as produced a great cheatsheet for reading R’s help pages as part of his forthcoming book. It should be available online at <http://socviz.co/appendix.html#a-little-more-about-r>

27.11.2 Getting help from the internet

The best way to get help for any topic is to just do an internet search like this: “R read.csv”. Usually the first thing on the results list will be the R help file, but the second or third will be a blog post or something else where a usually helpful person has discussed how that function works.

Sometimes for very basic R commands like this might not always be productive but its always work a try. For but things related to stats, plotting, and programming there is frequently lots of information. Also try searching YouTube.

27.11.3 Getting help from online forums

Often when you do an internet search for an R topic you'll see results from the website www.stackoverflow.com, or maybe www.crossvalidated.com if its a statistics topic. These are excellent resources and many questions that you may have already have answers on them. Stackoverflow has an internal search function and also suggests potentially relevant posts.

Before posting to one of these sites yourself, however, do some research; there is a particular type and format of question that is most likely to get a useful response. Sadly, people new to the site often get “flamed” by impatient pros.

27.11.4 Getting help from twitter

Twitter is a surprisingly good place to get information or to find other people knew to R. Its often most useful to ask people for learning resources or general reference, but you can also post direct questions and see if anyone responds, though usually its more advanced users who engage in twitter-based code discussion.

A standard tweet might be “Hey #rstats twitter, am knew to #rstats and really stuck on some of the basics. Any suggestions for good resources for someone starting from scratch?”

27.12 Other features of RStudio

27.12.1 Adjusting pane the layout

You can adjust the location of each of RStudio 4 window panes, as well as their size.

To set the pane layout go to 1. ”Tools” on the top menu 1. ”Global options” 1. “Pane Layout”

Use the drop-down menus to set things up. I recommend 1. Lower left: “Console” 1. Top right: “Source” 1. Top left: “Plot, Packages, Help Viewer” 1. This will leave the “Environment...” panel in the lower right.

27.12.2 Adjusting size of windows

You can clicked on the edge of a pane and adjust its size. For most R work we want the console to be big. For beginners, the “Environment, history, files” panel can be made really small.

27.13 Practice (OPTIONAL)

Practice the following operations. Type the directly into the console and execute them. Also write them in a script in the script editor and run them.

Square roots

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJzcXJ0KDQyKSJ9
```

The date Some functions in R can be executed within nothing in the parentheses.

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiJkYXRlKCKifQ==
```

Exponents The \wedge is used for exponents

```
eyJsYW5ndWFnZSI6InliLCJzYW1wbGUiOiI0Ml4yIn0=
```

A series of numbers A colon between two numbers creates a series of numbers.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiIxOjQyIn0=

logs The default for the `log()` function is the natural log.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsb2coNDIyIn0=

`log10()` gives the base-10 log.

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJsb2cxM0MikifQ==

exp() raises e to a power

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJleHAoMy43Mzc2NykidQ==

Multiple commands can be nested

eyJsYW5ndWFnZSI6InIiLCJzYW1wbGUiOiJzcXJ0KDQyKV4yXG5sb2coc3FydCg0MileMilcbmV4cChsb2coc3FydCg0MileM